

دانشگاه صنعتی شهرود

دانشکده: کامپیوتر و فن آوری اطلاعات
گروه: هوش مصنوعی

پایان نامه دوره‌ی کارشناسی ارشد مهندسی کامپیوتر - هوش مصنوعی

انتخاب بهینه‌ی کلمات کلیدی برای موتورهای جستجو

حمزه هدھدکیان

استاد راهنما:

دکتر مرتضی زاهدی

استاد مشاور:

دکتر حمید حسن بور

تیرماه 1390

تقدیم:
نے مادر کے
بaba
حضرت پدر کے
آرش
صوری

با تشکر

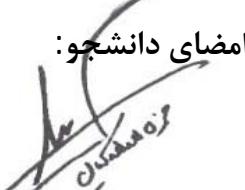
از اساتید دانشکده مهندسی کامپیوتر، که بدون راهنمایی های بی دریغ ایشان رسیدن به
پایان ممکن نبود

تعهدنامه

اینجانب حمزه هدھدکیان دانشجوی دوره کارشناسی ارشد رشته هوش مصنوعی دانشکده کامپیوتر و فناوری اطلاعات دانشگاه صنعتی شاهرود نویسنده پایان نامه تحت عنوان انتخاب بهینه کلمات کلیدی برای موتور های جستجو تحت راهنمایی های دکتر مرتضی زاهدی متعهد می شوم:

- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است.
- در استفاده از نتایج پژوهش های محققان دیگر به مرجع مورد استفاده استناد شده است.
- مطالب مندرج در پایان نامه تا کنون توسط خود یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است.
- کلیه حقوق این اثر متعلق به دانشگاه صنعتی شاهرود می باشد و مقالات مستخرج با نام «دانشگاه صنعتی شاهرود» و یا «Shahrood University of Technology» به چاپ خواهد رسید.
- حقوق معنوی تمام افرادی که در به دست آمدن نتایج اصلی پایان نامه تاثیرگذار بوده اند در مقالات مستخرج از پایان نامه رعایت می گردد.
- در کلیه مراحل انجام این پایان نامه، در مواردی که از موجودات زنده (یا بافت های آنها) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است.
- در کلیه مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است اصل رازداری، ضوابط و اصول اخلاقی رعایت شده است.

تاریخ 1390/9/19

امضای دانشجو:

گرمه مقدمه‌دان

مالکیت نتایج و حق نشر

- کلیه حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه های رایانه ای، نرم افزارها و تجهیزات ساخته شده) متعلق به دانشگاه صنعتی شاهرود می باشد. این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود.

دانشگاه صنعتی شاھروود

دانشکده: مهندسی کامپیوٹر و فناوری اطلاعات

پایان نامه ارشد آقای حمیده هدیده‌گیان

تحت عنوان

الانتخاب بینهایی کلیمات کلیدی درای محتواهای جستجو

در تاریخ ۱۴ مرداد ۹۷ توسط کمیته تخصصی زیر جبهت احمد مهرک کارشناسی ارشد مورد ارزیابی و با

درجه عالی مورد پذیرش قرار گرفت.

ردیف	اسماد ملکی	اصدار	اصادر راهنمایی
۱	دکتر حمیده هدیده‌گیان	دریافت	دکتر حمیده هدیده‌گیان
۲	دکتر حمیده هدیده‌گیان	دریافت	دکتر حمیده هدیده‌گیان

ردیف	اصادر	اصادر	اصادر
۱	دکتر حمیده هدیده‌گیان	دکتر حمیده هدیده‌گیان	دکتر حمیده هدیده‌گیان
۲	دکتر حمیده هدیده‌گیان	دکتر حمیده هدیده‌گیان	دکتر حمیده هدیده‌گیان
۳	دکتر حمیده هدیده‌گیان	دکتر حمیده هدیده‌گیان	دکتر حمیده هدیده‌گیان
۴			
۵			

چکیده

به طور کلی، روش‌های ارائه شده جهت استخراج خودکار کلمات کلیدی، سعی در بدست آوردن نتایج بهتر در معیارهایی مانند بازخوانی و دقت دارند. اگرچه این معیارها، میزان کارایی روش استخراج کلمات کلیدی در نقش یک انسان را نشان می‌دهد اما با توجه به نقش غیرقابل انکار موتورهای جستجو در دنیای امروز، به نظر می‌رسد که در انتخاب کلمات کلیدی علاوه بر توجه به معیارهای رایج بازیابی اطلاعات باید به افزایش میزان دسترسی پذیری متن توسط موتورهای جستجو نیز توجه ویژه‌ای شود.

در این تحقیق روشی جدید برای استخراج خودکار کلمات کلیدی ارائه شده است که همزمان با افزایش دسترسی پذیری متن، امتیاز مناسبی در معیارهای بازخوانی و دقت نیز کسب می‌کند. روش ارائه برای استخراج کلمات کلیدی از دو تابع امتیاز دهنده استفاده می‌کند: تابع امتیاز دهنده که کلمات کلیدی و تابع ارزیابی میزان دسترسی پذیری. تابع اول سعی در بالا بردن بازخوانی و دقت دارد در حالیکه تابع دوم در طول فرایند آموزش با استفاده از الگوریتم ژنتیک و بازخوردهای موتورهای جستجو به بهینه سازی ضرایب خصوصیات در تابع اول می‌پردازد.

همچنین در این پژوهه با بهره گیری از کلمات برجسته ساز، فرایند پس‌پردازشی ارائه شده که با گزینش نهایی کلمات کلیدی از میان لیست کلمات کاندید، منجر به بهبود کارایی روش در معیار دقت می‌شود. آزمایشات نشان می‌دهد که با به کارگیری تعداد تکرار مناسب در طی فرایند آموزش و ایجاد موازنۀ منطقی در کسب هریک از سه معیار یاد شده می‌توان به نتایج مطلوبی در هر سه معیار دست یافت.

واژگان کلیدی: استخراج خودکار کلمات کلیدی، بهینه‌سازی نتایج موتورهای جستجو، کلمات پر تکرار فارسی، الگوریتم ژنتیک، کلمات برجسته ساز فارسی، بهبود معیار دقت، افزایش دسترسی پذیری

لیست مقالات مستخرج از این پایان نامه

- [1] Hamzeh Hodhodkian, Morteza Zahedi,"**Improving Precision in Automatic Keyword Extraction Using Attention Attraction Strings**", Arabian Journal for Science and Engineering, July 2011
- [2] Hamzeh Hodhodkian, Morteza Zahedi," **An Efficient Approach for Keyword Selection; Improving Accessibility of Web Contents by General Search Engine**", International Journal of Web & Semantic Technology (IJWesT), July 2011

1.....	فصل اول: مفاهیم اولیه و تعاریف پایه.....
2.....	2-1 مقدمه
2.....	2-1 استخراج کلمات کلیدی زیر وظیفه ای از متن کاوی.....
4.....	3-1 کلمات کلیدی
4.....	1-3-1 تعریف کلمات کلیدی
5.....	2-3-1 کلمات کلیدی و فرایندهای متن کاوی
6.....	3-3-1 کلمات کلیدی و موتور های جستجو
7.....	4-1 چالش های موجود در پردازش متون فارسی.....
9.....	5-1 فرایند کلی استخراج خودکار کلمات کلیدی
11.....	1-5-1 تئوری های مسئله.....
14	فصل دوم: مروری بر کارهای گذشته.....
15	1-2 مقدمه
16	2-2 پایگاه های داده.....
17	3-2 فرایندهای پیش پردازش
18.....	1-3-2 یکسان سازی
21.....	2-3-2 تشخیص مرز واژه ها و جملات یا توکنیزه سازی
25.....	3-3-2 حذف کلمات پر تکرار.....
26.....	4-3-2 ریشه یابی کلمات

34	روش های استخراج کلمات کلیدی 4-2
34	1-4-2 تقسیم بندی تکنیکی روش ها
36	2-4-2 روش های امتیاز دهی به کلمات
41	2-4-3 پارامترهای ارزیابی کلمات استخراج شده
45	5-2 مروری بر تحقیقات انجام شده
48	فصل سوم: آشنایی با اصول اولیه SEO
49	1-3 مقدمه
49	2-3 ساختار یک موتور جستجو
50	1-2-3 تارگذار
50	2-2-3 خزنه
50	3-2-3 شاخص گذار
50	4-2-3 پایگاه داده
51	5-2-3 سیستم رتبه‌بندی
51	3-3 SEO اصول اولیه
51	1-3-3 تگ عنوان (<title>)</title>
52	2-3-3 متا تگ های توضیح
52	3-3-3 آدرس ها (URL)
53	4-3-3 کلمات

54	5-3-3 متون پیوندی
55	6-3-3 تصاویر
56	7-3-3 تگ های تیتر و زیر تیتر
57	فصل چهارم: پایگاه داده ها
58	1-4 مقدمه
58	2-4 نقطه شروع
60	3-4 انتخاب اسناد
62	4-4 پاک سازی تگ ها برای پایگاه داده
63	5-4 انتخاب کلمات کلیدی
67	فصل پنجم: فرایند های پیش پردازش
68	1-5 مقدمه
68	2-5 یکسان سازی
70	3-5 حذف کلمات پر تکرار
71	4-5 تعیین مرز جملات
72	5-5 تعیین مرز کلمات
73	فصل ششم: استخراج کلمات کلیدی
74	1-6 مقدمه

75	2-6 بخش آموزش
77	1-2-6 تابع امتیاز دهی
79	2-2-6 تابع ارزیابی
83	3-2-6 الگوریتم ژنتیک اجرا شده
86	3-6 فاز تست
86	1-3-6 استخراج نتایج
87	2-3-6 پس پردازش برای بهبود میزان دقت با استفاده از کلمات برجسته ساز
92	3-3-6 ارزیابی نهایی
93	فصل هفتم: نتایج عملی
94	1-7 مقدمه
94	2-7 بررسی تاثیرگذاری اندازه پایگاه داده
97	3-7 میزان تاثیر گذاری هریک از خصوصیات بر دسترسی پذیری
98	4-7 نتایج فرایند پس پردازش با استفاده از کلمات برجسته ساز
104	فصل هشتم: نتیجه گیری و پیشنهاداتی برای آینده
105	1-8 مقدمه
105	2-8 نتیجه گیری
106	3-8 پیشنهاداتی برای آینده
106	1-3-8 ایجاد پایگاه داده بزرگتر برای اطمینان بخشی به نتایج بدست آمده

107 2-3-8 استفاده از یک روش ریشه یابی

107 3-3-8 بررسی تاثیرات ایجاد یک مرز دینامیک امتیاز

فهرست شکل ها

شکل 1-1 شمای کلی سیستم استخراج کلمات کلیدی 10
شکل 2-1 نمودار تعداد رخداد کلمات بر اساس رتبه آن ها (تئوری لان) 13
شکل 1-3 نمونه ای از دو آدرس ناخوانا (شکل بالا) و ناخوانا (شکل پایین) 53
شکل 2-3 نمونه ای از یک عکس دارای alt (1) و عکس بدون alt (2) 55
شکل 1-4 روند اتخاذ شده برای ایجاد پایگاه داده برای این پروژه 59
شکل 2-4 نمودار توزیع آماری اسناد در پایگاه داده همشهری از دید موضوع 60
شکل 3-4 توزیع موضوعی اسناد پایگاه داده ایجاد شده در این پروژه 61
شکل 4-4-تعداد اسناد پایگاه داده براساس طول بر حسب کلمه 62
شکل 5-4 روند پاک سازی یک سند پایگاه داده همشهری 63
شکل 6-4 شکل رکورد ذخیره سازی اطلاعات و کلمات کلیدی یک سند 64
شکل 7-4 ساختار رکورد جدول ذخیره کننده کلمات کلیدی 64
شکل 8-4- تعداد اسناد موجود در پایگاه داده بر اساس تعداد کلمات کلیدی انتخاب شده برای آن 65
شکل 9-4 نمونه ای از یک سند پایگاه داده ایجاد شده 66
شکل 1-6 شمای کلی روش ارائه شده در این پروژه برای استخراج کلمات کلیدی 74
شکل 2-6 شبیه کد الگوریتم ژنتیک 84
شکل 3-6 روند استخراج کلمات برجسته ساز 89
شکل 4-6 فرایند گزینش نهایی کلمات کلیدی 91
شکل 1-7 نحوه تغییرات پارامترهای کارایی در اثر تغییرات تعداد اسناد پایگاه داده آموزش 96
شکل 2-7 مقایسه رشد γ (ضریب NTFIDF) با β (ضریب NSL) در طول فرایند آموزش 97

98 شکل 7-3 مقایسه رشد γ (ضریب NPLen) با α (ضریب NTFIDF) در طول فرایند آموزش

102 شکل 7-4 نمودار تغییرات میزان دقت بر حسب اندازه پایگاه داده

فهرست جداول

جدول 1-2 نمونه ای از دگرگونی کلمات در هنگام پیوند.....	28
جدول 2-2 هشت گروه فعل های فارسی	29
جدول 2-3 تعدادی از توابع شباهت مرسوم	40
جدول 2-4 افزار های ممکن برای مجموعه اسناد	43
جدول 6-1 نمونه ای از برچسب های HTML به همراه توضیح و مثال.....	81
جدول 7-1 تاثیرات اندازه پایگاه داده بر ملاک ارزیابی کارایی.....	96
جدول 7-2 نتایج بدست آمده برای فرایند پس پردازش با حجم مختلف پایگاه داده طی فاز آموزش برای روش 1	99
جدول 7-3 نتایج بدست آمده برای فرایند پس پردازش با حجم مختلف پایگاه داده طی فاز آموزش برای روش 2	100
جدول 7-4 نتایج بدست آمده برای فرایند پس پردازش با حجم مختلف پایگاه داده طی فاز آموزش برای روش 3	100
جدول 7-5 نتایج بدست آمده برای فرایند پس پردازش با حجم مختلف پایگاه داده طی فاز آموزش برای روش ارائه شده در این پژوهه	101
جدول 7-6 میزان تغییرات میانگین دقیق در اندازه های مختلف پایگاه داده	101
جدول 7-7 میزان تغییرات در متوسط ملاک های کارایی در زمان استفاده از 400 سند	103

فصل اول:

مفاهیم اولیه و تعاریف پایه

1-1 مقدمه

در این فصل به دنبال تعریف جایگاه پروژه انجام شده در تحقیقات امروزی و دلایل نیاز به آن هستیم. در ابتدای فصل با تعریف متن کاوی و بازیابی اطلاعات سعی در مشخص کردن جایگاه روش‌های استخراج کلمات کلیدی داریم. سپس با تعریف مفهوم کلمات کلیدی و ارائه ساختار کلی روش‌های استخراج کلمات کلیدی به بررسی تئوری‌های مرتبط با این فرایند و اشکالات و چالش‌های پیش رو با آن می‌پردازیم.

2-1 استخراج کلمات کلیدی زیر وظیفه‌ای از متن کاوی

متون دیجیتال از دید یک دستگاه کامپیوتر رشته‌ای ساده از کاراکتر‌های در حالیکه در یک ساختار انسانی این رشته کاراکترها حاوی ارتباطات و معانی فراوانی در خود هستند. آنچه که مسلم است با توجه به حجم تولید اطلاعات در دنیای کنونی، نیاز شدیدی به خودکار سازی فرایند پردازش متون توسط کامپیوترا احساس می‌شود و به همین علت نیازمند روش‌ها و الگوریتم‌هایی برای پیش پردازش یا پردازش این متون جهت تبدیل آن‌ها به الگوهایی مناسب با ساختار‌های کامپیوترا هستیم.

مفهوم متن کاوی بصورت "یافتن قوانین جذاب (یه معنای مخفی در ذات متن، که بصورت بالقوه مفید هستند) در حجم وسیعی از پایگاه داده متنی" تعریف می‌شود که می‌تواند به "پیدا کردن مفهوم یا چیکده اطلاعات موجود در یک سند از روی متن ظاهری آن" نیز گسترش یافته و تعبیر شود [1].

پروسه‌های متن کاوی را می‌توان در سه شاخه، بازیابی اطلاعات، استخراج اطلاعات و کشف دانش بررسی کرد.

بازیابی اطلاعات اصولاً با بازیابی مستندات و مدارک مرتبط است. به این معنا که در این روش دانشی کشف نمی‌شود بلکه تنها بسته‌ای از کلمات از میان مجموعه‌ای از مستندات بیرون کشیده می‌شود که با توجه به نیاز مطرح شده از سوی کاربر، بیشترین ارتباط را با مفهوم درخواست یا جستجو شده داشته باشد.

اما برای تعریف استخراج اطلاعات، در نظر بگیرید که موسسه‌ای بسیار موفق دارید و از این‌رو تعداد زیادی ایمیل در روز دریافت می‌کنید. شما می‌خواهید سابقه‌ی این ایمیل‌ها را ثبت کنید. اینکه چه کسانی آن‌ها را فرستاده‌اند، تاریخ فرستادن چه روزی بوده است، عنوان و متن آن چه بوده است و نظایر آن. این امر با بیرون کشیدن این اطلاعات از تک تک ایمیل‌ها و پر کردن یک پایگاه داده از این اطلاعات میسر است. می‌توانید برای این کار با تعریف یا کشف یک قالب از داده‌هایی که به آن سرو کار دارید به این هدف بررسید. مثلاً می‌شود برنامه‌ای داشت که به طور اتوماتیک به دنبال کلمه‌ی "عنوان" در سند بگردد و آنچه را که بعد از آن آمده است را به عنوان یک فیلد در پایگاه داده پر کند. هر چند یافتن این قالب در داده‌های غیر ساخته یافته‌ی دیگر ممکن است به هیچ وجه ساده نباشد، اما زمانی که شما این کار را به پایان برده‌ید، با داده‌های کاملاً ساخت یافته مواجه‌ید که از درون ایمیل‌ها بدست آمده است. اما از سویی دیگر این نتیجه باز هم همان اطلاعاتی است که داشته‌اید. به این معنی که هیچ چیز جدیدی از آنچه دارید کشف نشده است و همان را که قبلاً می‌دانستید هنوز هم می‌دانید.

برای روشن شدن مفهوم کشف دانش همان مثال قبل را در نظر بگیرد با این تفاوت که اینبار به دو مفهوم A و B که از مجموعه برخورد کرده ایم که رابطه‌ای منطقی دارند. مثلاً اینکه "مناطق بارانی" و "کشت برنج" دارای رابطه‌ای به این صورتند که کشت برنج به مناطق بارانی نیاز دارد ($A \rightarrow B$). به علاوه فرض کنید که مفهوم B نیز با مفهوم C به همین شکل دارای ارتباط است. مثلاً "مناطق شمال ایران" و

“مناطق بارانی” این رابطه را با هم دارند که مناطق شمال ایران منطقه‌ی بارانی هستند ($B \rightarrow C$). سیستمی که کمی باهوش‌تر باشد می‌تواند از این اطلاعات استخراج شده که در میان اطلاعات ما بوده‌اند. اطلاعاتی نو بدست آورد به این معنی که به رغم عدم وجود این رابطه بطور مستقیم سیستم کشف می‌کند که “مناطق شمال ایران” قابل “کشت برنج” است ($A \rightarrow C$). این مثال نشان می‌دهد که کشف دانش به چه معناست.

با وجود آنکه هدف اصلی آنالیز متون فهم محتوای متنی و حتی کشف دانش است اما این مسئله هنوز در دنیای الگوریتم‌ها و کامپیوترهای امروزه‌ای مسئله‌ای پیچیده به نظر می‌رسد که ما را بر آن می‌دارد که در سطوح پایین‌تر به دنبال زیر وظایفی باشیم که البته هریک از این زیر وظایف به خودی خود ارزش زیادی دارند.

یکی از این زیر وظایف استخراج کلمات کلیدی است که نقش بسزایی در عملی شدن سایر فرایندها دارد. در ادامه به تعریف کلمات کلیدی می‌پردازیم و سپس بعد از ارائه مختصری از خصوصیات و مشکلات زبان فارسی به بحث‌های تئوریک در مورد مفاهیم استخراج کلمات کلیدی خواهیم پرداخت. و در پایان فصل با ارائه یک شکل کلی از فرایند استخراج کلمات کلیدی به تشریح ساختار پایان نامه می‌پردازیم.

1-3-1 کلمات کلیدی

1-3-1-1 تعریف کلمات کلیدی

به مجموعه‌ای از کلمات یا عبارات کوتاه که بیان کننده منظور اصلی یک متن باشند کلمات کلیدی آن متن گفته می‌شود [2].

مجموعه کلمات کلیدی در حالت ایده آل در حین آنکه مفهوم اصلی متن را به طور کامل توضیح می دهد باید بتواند متن را از سایر متون به صورت منحصر به فرد جدا کند.

مسئله تعیین کلمات کلیدی به دو شکل تعریف می شود:

1. منصب کردن کلمات کلیدی^۱: در این حالت به یک متن تعدادی از کلمات کلیدی که از یک دامنه محدود انتخاب شده اند منصب می شود. به عنوان مثال سایت خبری تابناک بخش های خبری ویژه ای برای مسئله فلسطین در نظر گرفته است. حال هر خبری که در آن ارتباطی با فلسطین وجود داشته باشد یکی از کلمات کلیدی که به این خبر منصب خواهد شد کلمه فلسطین خواهد بود. این دامنه محدود منجر می شود که مسئله انتصاب کلمات کلیدی به طور کامل به یک مسئله کلاسیندی تبدیل شود که در آن یک متن به عنوان یک نمونه عضوی از یک یا چند کلاس از کلاس های ایجاد شده توسط کلمات کلیدی محدود خواهد شد.

2. استخراج کلمات کلیدی^۲: انتخاب (یا پیشنهاد) کلماتی از یک متن که نشان دهنده مفهوم کلی آن باشد این پرسه بیشتر شبیه به یک رتبه بندی است که رتبه بالاتر به معنای قدرت بیشتر عبارت در تشریح متن مورد بررسی است.

2-3-1 کلمات کلیدی و فرایندهای متن کاوی

کلمات کلیدی کوچکترین واحدی نحوی هستند که می توانند معنای سند را با خود حمل کنند از این‌رو در بسیاری از فرایندهای خودکار کاربردی متن کاوی از قبیل شاخص‌گذاری^۳، خلاصه‌سازی^۴، کلاسه‌بندی^۵،

¹ Keywords assignment

² Keywords extraction

³ Indexing

⁴ summarization

⁵ classification

خوش بندی^۱، فیلترینگ^۲، تشخیص و ردیابی موضوع^۳ و نیز مصورسازی اطلاعات^۴ پایه تحقیقات بر کلمات کلیدی بنا نهاده می شود که باعث می شود استخراج کلمات کلیدی هسته‌ی تکنولوژیکی تمامی فرایندهای متن کاوی به شمار آید.^[3]

از آنجا که به طور معمول در بسیاری از اسناد کلمات کلیدی مشخص نشده و نیز انجام این فرایند با استفاده از نیروی انسانی زمانبر بوده و هزینه زیاد و خطاها احتمالی در پی دارد در نتیجه تحقیقات زیادی برای اتوماتیک سازی روش‌های استخراج کلمات کلیدی انجام گرفته است. اما از آنجا که زبان فارسی ویژگی‌های منحصر به فرد خود را داراست می طلبد که به مقوله استخراج کلمات کلیدی در متون فارسی نگاه ویژه‌ای معطوف شود.

3-3-1 کلمات کلیدی و موتورهای جستجو

موتورهای جستجو در دنیای امروز نقش غیرقابل انکاری در ارائه دانش به کاربران دارند. در واقع حتی گاهی این مسئله مطرح می شود که اگر یک مطلب توسط یک موتور جستجو دیده نشود شاید هرگز دیده نشود. همچنین بررسی‌ها نشان می دهد نحوه رتبه بندی این موتورها بسته به کلمات کلیدی به کار برده شده در مطلبی است که موتور جستجو آن را به عنوان پاسخ یک عبارت جستجو^۵ ارائه می‌دهد. در نتیجه منطقی به نظر می رسد که در انتخاب کلمات کلیدی یک متن، غیر از توجه به این که چقدر کلمات انتخابی به تعریف مفهوم اصلی متن می پردازند، به این مسئله نیز توجه کنیم که این کلمات از نظر یک موتور جستجو چه قدر اهمیت دارند.

¹ Clustering

² filtering

³ Topic detection and tracking

⁴ Information visualization

⁵ Query

به همین علت فرایند استخراج کلمات کلیدی ارائه شده در این پروژه به میزان اهمیت یک کلمه از نظر موتور جستجو نیز توجه کرده و فرایند ارزیابی کلمات کلیدی استخراج شده را بر پایه نتایج موتورهای جستجو قرار داده است.

4-1 چالش‌های موجود در پردازش متون فارسی

زبان فارسی که در گستره‌ی جغرافیایی ایران، افغانستان و برخی نقاط تاجیکستان مورد استفاده قرار می‌گیرد یک زبان خاورمیانه ایست که در دامنه لغات و حتی گاهی قواعدش متأثر از زبان‌های دیگر مثل عربی، ترکی، و بعضی زبان‌های اروپایی مانند انگلیسی و فرانسوی است. از این رو آنالیز ریخت‌شناسی این زبان به ارتباط با زبانی‌ها زیادی نیاز دارد [4]

زبان فارسی از دیدگاه ویژگی‌های زبانی یک زبان پیوندی و ضمیر انداز است و در شکل رسمی آن ترتیب فاعل - مفعول - فعل برای جملاتش قید می‌شود این در حالیست که استثنایات مجازی که به خاطر فرایندهایی مثل نامکانی، بهم ریختگی، حرکت جهت برجسته سازی، تاخیر، شکافت و شبه شکافت و غیره بوجود می‌آیند این زبان را به یک زبان بدون ترتیب نزدیک کرده است [5].

به خاطر تفاوت‌های ویژه طبیعت این زبان، کار کردن بر روی سیستم‌های بازیابی اطلاعات در زبان فارسی نیازمند توجهات ویژه‌ای است که در زبان‌های دیگر نظیر انگلیسی به این مسئله نیاز نیست [6].

تحقیقاتی که در زمینه پردازش و درک متون فارسی صورت می‌گیرد هم با مشکلاتی که در سایر زبان‌ها وجود دارد سروکار دارد هم با مشکلاتی که مخصوص زبان فارسی است. برخی از این مشکلات ناشی از تحقیقات کم در رابطه شناخت قواعد زبان فارسی و برخی دیگر به علت محدودیت‌های سیستم‌های هوش مصنوعی است [5] که بطور خلاصه در فهرست زیر ذکر شده است:

1. عدم وجود منابع زبانی مناسب و کافی برای زبان فارسی مانند

- واژه های تک زبانه و چند زبانه محاسباتی
- واژه های معنایی و متصل به هستان شناسی¹
- هستان شناسی های لغوی
- هستان شناسی جامع عمومی و تخصصی
- مجموعه متون² عمومی و تخصصی ساده یا برچسب خورده (با برچسب های اجزاء کلام، کسره ی اضافه، نقش های موضوعی، مفاهیم و روابط مفهومی و غیره)
- مجموعه مدون قوانین ساخت واژه ای و دستوری پوشای

2. ساختارهای متفاوت فایلی اسناد فارسی مانند نسخه های گوناگون Microsoft word ، زرنگار،

و غیره Pe2

3. استانداردهای گوناگون دخیل در کدگذاری اسناد مختلف فارسی: این استانداردها و جزئیات آن ها

به حدی متنوع و متفاوت است که بحث در مورد آن ها خسته کننده می نماید. اما به طور عمومی با وجود اشکالاتی که استاندارد³ نرم افزار هایی مثل Word و VisualStudio وجود دارد استاندار کدگذاری این نرم افزارها در حال رایج شدن است.

4. عدم وجود استانداردهای شیوه ی نگارش این مسئله منجر به مشکلاتی از قبیل:

- مشکل تشخیص مرز کلمات به خاطر شیوه های گوناگون ممکن در نوشتار
- ایجاد تفاوت در صورت های کلمات به خاطر عدم وجود صورت های ممکن برای انواع پسوندها و پیشوندها

¹ Anthology

² از مجموعه متون به عنوان ترجمه کلمه corps از عبارت پیکره نیز استفاده می شود.

³ بزرگترین تفاوت استاندارد این نرم افزارها با استاندارد 3342 موسسه استاندارد ایران رعایت نکردن ترتیب برای چهار حرف مخصوص زبان فارسی («پ»، «ز»، «گ»، «ج» می باشد).

5. حذف اطلاعات گویشی مانند صوت‌ها که باعث ایجاد مشکلاتی از قبیل:

- اشکال در تشخیص کسره اضافه: کسره اضافه بزرگترین عامل تشخیص مرز گروه‌های اسمی است که متساقنه با حذف صوت‌ها، حذف شده و ایجاد مشکل می‌کند.

- مسئله‌ی ابهام در معانی در کلماتی مثل رَوَد و رُود می‌شود.

6. افعال و کلمات مرکب یا کلماتی که هم آبی دارند

7. مسئله‌ی هم نگاره‌ها در کلماتی مثل «شیر»

8. معناشناسی و مشکلات تحلیل معنایی

۱-۵ فرایند کلی استخراج خودکار کلمات کلیدی

به طور کلی تمامی روش‌های استخراج کلمات کلیدی را می‌توان بر ساختاری شبیه به شکل ۱-۱ انطباق داد. (مسیر‌های نقطه چین به معنای اختیاری بودن ادامه مسیر می‌باشد)

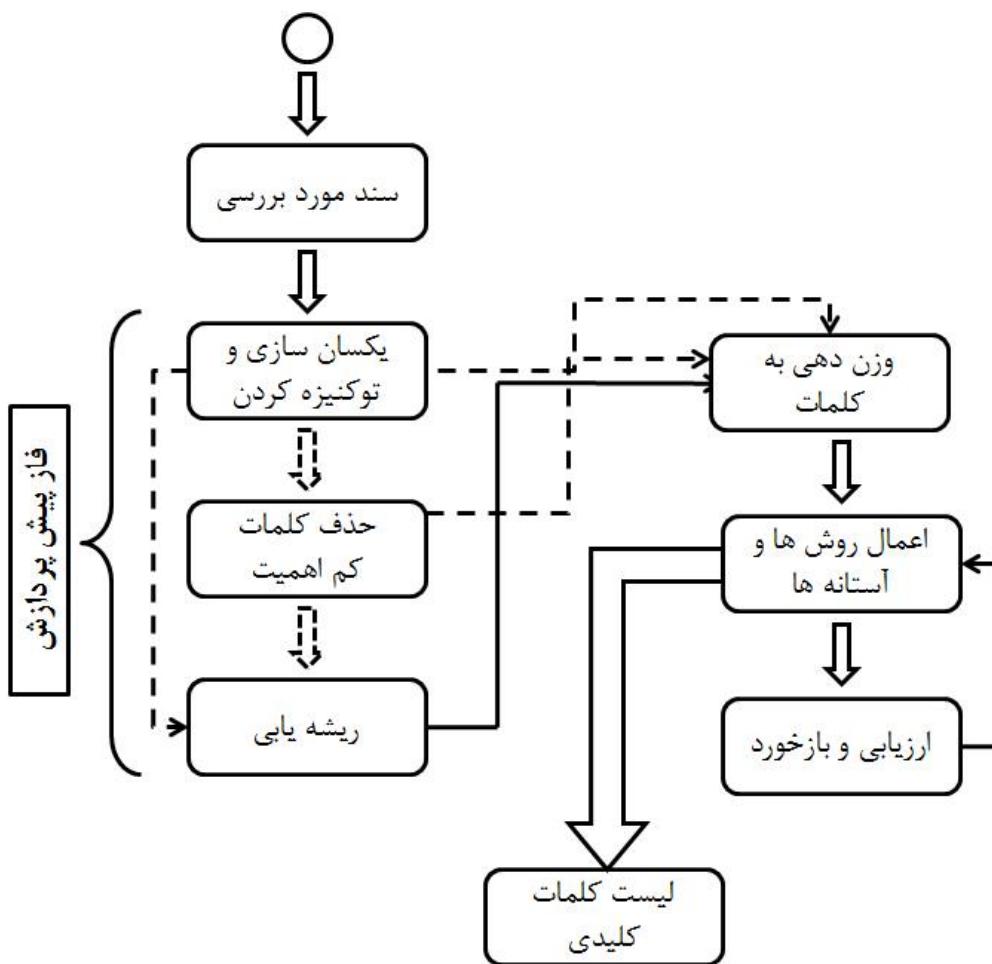
برای استخراج مجموعه کلمات کلیدی ابتدا سند مورد بررسی با استفاده از قواعد احتمال یا علائم نگارشی و یا سایر روش‌ها به بسته‌های کلمات یا توکن‌ها تبدیل می‌شود سپس کلماتی که معنای خاصی ندارند و بیشتر نقش دستوری دارند یا به خاطر عمومی بودن و تکرار زیادشان تمایزی بین این متن با متون دیگر ایجاد نمی‌کنند حذف می‌شوند.^۱

البته در بعضی زبان‌ها مثل زبان فارسی که نوع رمزگذاری‌ها تفاوت دارد و یا حتی چند شکل درست برای یک ترکیب می‌تواند وجود داشته باشد پیش از هر کاری بهتر است فرایندی به نام یکسان‌سازی^۲

¹ کلمات مذکور را کلمات پرتکرار می‌نامند که در فصل دوم و چهارم به طور مفصل مورد بحث قرار می‌گیرند.

² Unification

اجرا شود. فرایند یکسان سازی، حالات متفاوت عبارات یا رمزگذاری های متفاوت را به یک نوع یکسان تبدیل می کند.



شکل 1-1 شماتیکی سیستم استخراج کلمات کلیدی

ریشه یابی¹ در واقع تبدیل کردن مجموعه کلمات با ریشه یکسان به یک کلمه یکسان است. در ریشه یابی هیچ تقیدی برای ریشه قائل نیستیم و اینکه یک ریشه طبق تعاریف رسمی دستور زبان بدست آید یا نه اهمیتی ندارد.

¹ stemming

همانطور که در شکل مشخص است استفاده از مراحل پیش پردازش بسته به نوع روش به کار برده شده است. در واقع در برخی روش‌ها ریشه یابی انجام نمی‌شود یا ممکن است برای دقت بالاتر کلمات پرتکرار نیز حذف نشوند.

پس از آنکه لیست کلمات استخراج شد به کلمات وزن داده می‌شود. روش‌های وزن دهنده کلمات یکی از بزرگترین متمایز کننده های روش‌های مختلف محسوب می‌شود. با اعمال یک آستانه مثلاً یک امتیاز خاص یا یک تعداد کلمه خاص لیستی از کلمات کلیدی بدست می‌آید. به طور معمول این کلمات کلیدی با استفاده از معیارهایی که هر روش تعریف می‌کند مورد ارزیابی قرار می‌گیرد تا لیست نهایی استخراج شود. در روش‌هایی که از یادگیری ماشینی کمک می‌گیرند معمولاً بازخوردی از ارزیابی برای بهینه‌سازی و ارتقا امتیاز کلمات نیز گرفته می‌شود.

1-5-1- تئوری‌های مسئله

دمسئله استخراج کلمات کلیدی بر پایه تئوری لان¹ و قانون ZIPF بنا نهاده شده است که به مرور این دو تئوری می‌پردازیم

تئوری ZIPF: این تئوری بیان می‌دارد که اگر تعداد وقوع یک کلمه F بوده و در لیست صعودی مرتبه شده کلمات بر حسب تعداد، رتبه r را داشته باشد آنگاه

$$F^*r=K$$

¹ Lane

که K یک عدد ثابت است. یعنی اگر کلمه‌ای مثل «را» رتبه ۱ را در لیست تکرار کلمات با ۲۰۰۰ تکرار داشته باشد توقع داریم که کلمه‌ای مثل «می شود» که رتبه ۲۰ را در لیست تکرار کلمات دارد ۱۰۰ بار تکرار شده باشد. یک نتیجه واضح این قانون اینست که تعداد کمی از کلمات عمومی اند و زیاد تکرار می‌شوند، در حالیکه که تعداد متوسطی از کلمات با تعداد مناسبی تکرار می‌شوند و تعداد زیادی از کلمات با تکرار بسیار کم وجود دارد.

تئوری لان

این تئوری بیان می‌دارد که تعداد وقوع یک کلمه در یک متن میزان اهمیت کلمه را نشان می‌دهد. همچنین میزان اهمیت یک جمله را می‌توان براساس وضعیت نسبی اهمیت کلمات آن جمله تعیین کرد. بر اساس این تئوری کلمات تکرار شونده را می‌توان برای استخراج کلمات و جملات تشریح کننده سند به کار گرفت.

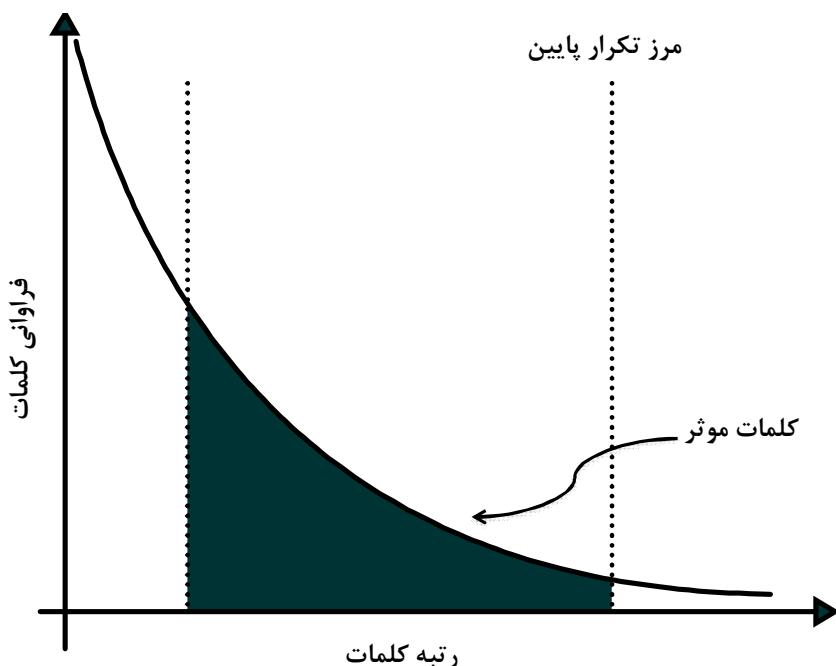
در نمودار شکل ۲-۱ f تعداد رخداد کلمات مختلف و r رتبه آن است و به نوعی نشان دهنده قانون ZIPF می‌باشد.(رابطه ۱-۱)

لان با اعمال یک آستانه گیر بالا و پایین بر روی این منحنی، توانست کلمات کم اهمیت را جدا کند. کلماتی که از آستانه‌ی بالا تجاوز کرده باشند، به عنوان کلمات بی ارزش و پر تکرار درنظر گرفته می‌شوند، و کلماتی که از آستانه پایین کم تر باشند، کلمات نادر و کمیاب نامیده می‌شوند که در محتویات مقاله یا سند شرکت چندانی ندارند.

$$F * r \xrightarrow{\approx} \text{cons tant}$$

$$1-1$$

همچنین لان با استفاده از این فرض که کلمات پر تکرار بین این دو مرز بیشترین تفکیک را بین سند موجود و سایر اسناد ایجاد می کند به این نتیجه رسید که ارزش کلمات از میانه به سمت مرزها کمتر شده تا به صفر می رسد. البته آستانه های تعریف شده با استفاده از سعی و خطا بدست آمده و هیچ پیش گویی نمی توان در آن مورد انجام داد. اما جالب اینجاست که این ایده، پایه اکثر کارهای بعدی بازیابی اطلاعات شده است همانطور که خود لان نیز از آن برای ایجاد یک چکیده نویس خودکار استفاده کرده است.



شکل 1-2 نمودار تعداد رخداد کلمات بر اساس رتبه آن ها (ئئوری لان)

فصل دوم:

مژوی بر کارهای گذشته

1-2 مقدمه

همانطور که در فصل اول بصورتی کوتاه بررسی کردیم، استخراج کلمات کلیدی شامل دو بخش پیش پردازش و انتخاب کلمات است.

روش های ارائه شده برای انتخاب کلمات کلیدی را می توان از جنبه های مختلفی تقسیم بندی کرد:

به عنوان مثال مانند تمامی روش های هوش مصنوعی می توان برای فرایند انتخاب خودکار کلمات کلیدی حالت با ناظر و بدون ناظر در نظر گرفت به این معنا که اگر برای هر سند مورد بررسی یک مجموعه کلمات کلیدی مشخص کنیم و سیستم را بر اساس این کلمات مشخص شده آموزش دهیم سیستم یک سیستم یادگیری با ناظر خواهد بود و در غیر اینصورت سیستم یک سیستم یادگیری بدون ناظر بشمار می آید.

همچنین می توان روش ها را برپایه نوع وزن دهی به آماری و غیر آماری تقسیم کرد.

اما تقسیم بندی از نظر نوع تکنیک که در بخش 4-2 به آن می پردازیم جامع ترین نوع تقسیم بندی در بسیاری از جهات محسوب می شود

در این فصل به ارائه خلاصه ای از کارهای انجام شده مرتبط با تحقیق می پردازیم این فصل شامل سه بخش پایگاه های داده، فرایندهای پیش پردازش و روش های استخراج کلمات کلیدی است. اگرچه تاکید ما در این فصل بر روش های استفاده شده بر روی متون فارسی است اما مروری بر تحقیقاتی که نتایج یا بخشی از روش های آن ها نیز در پژوهش آمده است نیز صورت گرفته است.

2-2 پایگاه های داده

وجود یک مجموعه متنی استاندارد پیش نیاز هر گونه تحقیق بر روی ساختارهای بازیابی اطلاعات محسوب می‌باشد. این در حالیست که کارهایی که اخیراً در بدست آوردن اطلاعات از زبان فارسی شده معمولاً از مشکل فقدان مجموعه متنی استاندارد رنج می‌برد در این بخش از فصل دوم به کارهای صورت گرفته بر روی جمع آوری متنی استاندارد به عنوان پایگاه داده ای برای استفاده در تحقیقات بازیابی اطلاعات می‌پردازیم.

پایگاه داده مطرح شده در [7] حاوی 25 مگابایت از اسناد مجلس است و به علت موضوعیت تک بعدی و سایز کم می‌تواند نتایج کار را تحت الشاعع قرار دهد.

مجموعه متنی شیراز [8] شامل 10 مگابایت متنی دو زبانی است که برای تست ترجمه ماشینی در دانشگاه نیو مکزیکو طراحی شده است.

پایگاه داده FLDB¹ ساخته شده از متنی منتخبی از متنی فارسی امروزی رسمی و غیر رسمی است و یک سری ورودی‌های دیکشنری و لیستی شامل حدود 3 میلیون عدد از کلمات است [9]. اگرچه بزرگی FLDB در متنی خوب ساختیافته فارسی خوب است اما به اندازه کافی اطلاعات مناسب برای سیستم های بازیابی اطلاعات ندارد.

یکی دیگری از مجموعه سند‌های فارسی محک است [10]. این مجموعه برای ارزیابی روش‌های استخراج اطلاعات ایجاد می‌شود اما دارای 3007 سند است که آن را برای سیستم‌های بزرگ استخراج اطلاعات نامناسب می‌سازد.

¹ Farsi Language Data Base

همشهری [11] پایگاه داده ایست که بر اساس مقالات روزنامه با توجه به مشخصات دقیق¹ درست شده است. به علاوه اطلاعات آماری در مورد اسناد و جستجوهایی مربوط به آن ها در این مقاله آمده است. این مجموعه بزرگترین مجموعه زبان فارسی جمع آوری شده تا کنون می باشد. پایگاه داده کنونی همشهری 1600000 سند دارد و به همراه 65 جستجو و نتایج تفضیلی آن برای استفاده محققین بر روی وب² بصورت رایگان قرار گرفته است.

جامعه آماری مناسب در موضوعات مختلف و نیز ساختار مناسب ارائه باعث شد (تصورت XML) که مجموعه همشهری برای ساخت پایگاه داده خاص این پروژه استفاده شود. که در فصل چهارم به این موضوع می پردازیم.

3-2 فرایندهای پیش پردازش

همانطور که در شکل 1-1 نشان داده شده معمولاً پیش از اعمال روش های استخراج کلمات کلیدی فرایندهایی برای پیش پردازش متن و ایجاد زمینه برای کسب نتایج بهتر صورت می گیرد.

به طور کلی انواع پیش پردازش را به چهار دسته تقسیم بندی کرد:

1. یکسان سازی

2. توکنیزه سازی (تشخیص مرز کلمات و جملات)

3. حذف کلمات پر تکرار

¹ کنفرانس بین المللی بازیابی اطلاعات و پردازش زبان های طبیعی که به منظور ارزیابی سیستم های بازیابی اطلاعات برگزار می شود.

² <http://ece.ut.ac.ir/DBRG/hamshahri/>

4. ریشه یابی

کارایی روش های پیش پردازش می تواند به دو صورت مطلق و در خلال فرایند کلی بررسی شود

سنجدش کارایی پیش پردازش بصورت مطلق: در این حالت یک فرایند با ارائه تابع مناسب و قابل درک برای سیستم های انسانی نتایج پیش پردازش را مورد بررسی قرار می دهد به عنوان مثال فرض کنید یک سیستم پیش پردازش بدون استفاده از فهرست ثابت کلمات پر تکرار به شناسایی و حذف پویای آنان می پردازد در این صورت تعداد کلماتی که یک کارشناس آنها را در متن بی ارزش و پر تکرار می داند و از متن حذف نشده باشد به تعداد کلماتی که حذف شده اند می تواند نسبت خط را تعیین کند.

سنجدش کارایی پیش پردازش در خلال فرایند کلی: معمولاً زمانی که تابع ارزیابی مناسبی برای یک فرایند پیش پردازش قابل تعریف نباشد می توان میزان کارایی پیش پردازش را براساس نتایج فرایند کلی در حضور و عدم حضور روش بررسی کرد. به عنوان مثال فرض کنید در حالت قبل بحث داشتن کارشناس برای کلمات بی ارزش مسئله را به داوری انسانی مرتبط می کند که هم خطابذیر و هم هزینه بر است در این صورت میزان موفقیت روش با استفاده از فرایند پیش پردازش نسبت به میزان موفقیت روش در عدم استفاده از فرایند پیش پردازش می تواند ملاک خوبی برای سنجدش کارایی روش باشد.

1-3-2 یکسان سازی

در زبان فارسی به علت شرایط خاصی که وجود دارد گاهی با وجود آنکه دو واژه کاملاً یکسان هستند، به علت نوع رمزگذاری نویسه ها یا تعدد حالت های نوشتاری ممکن، از دید یک پردازشگر متن دو واژه

متفاوت تلقی می شوند. به عنوان مثال، پیشوند «می» و پسوند «ها» در ابتدا و انتهای واژه‌ها، ممکن است

به سه صورت مختلف دیده شوند:

جدا بدون فاصله	جدا با فاصله	چسبان
کتاب‌ها	کتاب‌‌ها	کتابها
می‌رود	می‌‌رود	میرود

و یا واژه‌ها «مسئول» و «مجموعه» به صورت‌های زیر در اسناد مختلف دیده می شوند:

مسئول	مسئول	مسئول
مجموعه	مجموعه	مجموعه‌ی

یا حروف «ی» و «ک» که در متون مختلف ممکن است داری شکل یکسان اما رمزگذاری متفاوت باشند.

برای جلوگیری از بروز این مشکل، نیازمند یکسان سازی متون مورد پردازش هستیم.

به عنوان یک مثال از یکسان سازی، در مورد پسوند «ها» می‌توان با تبدیل شکل‌های «جدا با فاصله» و

«جدا بدون فاصله» به شکل چسبان تمامی اشکال‌ها را به حالتی یکسان تبدیل کرد.

پروسه یکسان سازی دارای دو بخش است:

1-3-2 یکسان سازی تکنیکی

این پروسه شامل یکسان سازی نوع رمزگذاری بر اساس یک استاندارد مشخص و همچنین یکسان سازی

با یک ساختار فایلی یکسان است. انواع استاندارد برای رمز گذاری نویسه‌ها وجود دارد که برخی از آن‌ها

توسط نرم افزارهایی کاربردی انتخاب شده و برخی دیگر توسط موسسه استاندارد. اگرچه انتخاب موسسه استاندارد از نظر رسمی بودن اعتبار خوبی دارد اما، در واقع نمی‌توان گفت کدام استاندارد بر آن یکی ترجیح دارد چراکه به هر حال با وجود تمامی مشکلات استاندارد انتخاب شده توسط نرم افزارهای معترض کارایی این نرم افزارها در کفه‌ی ارزیابی سنگینی زیادی نشان می‌دهد.

اما در بحث یکسان سازی ساختار فایل باید گفت با وجود قالب‌های گوناگون رایج مورد استفاده در زبان فارسی بهترین قالب نمایش XHTML است. این قالب به علت آنکه از سوی وب مورد پشتیبانی قرار می‌گیرد و همچنین مستقل از سیستم عامل است مناسب ترین قالب به شمار می‌رود. اما به هر حال برای یکسان سازی باید برنامه‌ای برای تبدیل قالب‌های دیگر به XHTML، یا پیدا کرد یا نوشت.

2-1-3-2 یکسان سازی بر اساس دانش زبان

این یکسان سازی شامل چگونگی نوشتنهای است. این مسئله که «ها» علامت جمع به کلمه بچسبد یا نه و یا آنکه "ی" در پایان کلماتی که به «ه» ختم می‌شوند چگونه باید بیاید در این یکسان سازی حل می‌شود. اگرچه استفاده از قواعد استاندارد و رایج برای یکسان سازی ارجح است اما در صورتیکه سیستم یکسان سازی شما به عنوان یک بخش داخلی از یک کل به شمار می‌رود می‌توانید از هر قاعده‌ای استفاده کنید.¹.

در مرجع [12] مجموعه‌ای نسبتاً جامع و موثر از تبدیل‌ها برای یکسان سازی متون فارسی ارایه شده است. نکته‌ای که در یکسان سازی باید دانست آنست که هرچه فرایند یکسان سازی جامع‌تر و مفصل‌تر باشد نتیجه کار بهتر است و استثنایات بیشتر تحت پوشش قرار می‌گیرند اما در مقابل پیچیدگی زمانی مسئله بالاتر می‌رود.

¹ فرایند یکسان سازی در صورت استفاده از قواعد استاندارد نرمال سازی (normalization) نیز نامیده می‌شود.

2-3-2 تشخیص مرز واژه ها و جملات یا توکنیزه سازی

توکنیزه سازی یا تقسیم کردن رشته های کاراکتری متن به کلمات یکی از فرایندهای پیش پردازش اصلی در اکثر روش ها به شمار می رود. همچنین از آنجا که برخی از روش های امتیاز دهی به کلمات مورد سیستم های استخراج کلمات کلیدی، از مکان کلمه در جمله یا مکان جمله نسبت به متن نیز بهره می گیرند فرایند تعیین مرز جملات نیز گاهی به عنوان یک فرایند پیش پردازش در نظر گرفته می شود.

1-2-3-2 تعیین مرز جملات

در بیشتر مواقع، تعیین مرز جمله ها از طریق بررسی عالیم جدا کننده انجام می شود. عالیمی که برای تعیین مرز جمله از آنها استفاده می شود، عبارتند از: «.»، «؟»، «!»، «؟؟»، «؟؟؟» و «:». باید توجه داشت که جستجو برای یافتن این عالیم به تنها یکی کافی نیست و در برخی موارد مشکلاتی را به همراه دارد. به عنوان مثال، در زبان فارسی بعضی از واژه های اختصاری به صورت چند حرف که با نقطه از هم جدا شده اند، ظاهر می شوند (مانند «ه.ق.» که مخفف «هجری قمری» است). در صورتی که تعداد این حالات مشکل ساز در پیکره زیاد باشد، باید از روش های پیچیده تری استفاده شود، که قادر به شناسایی این موارد باشند. همچنین ممکن است در متونی که بصورت غیر رسمی نوشته شده یا قواعد نوشتاری در آن ها رعایت نشده (این امر در متون تحت وب اتفاقی دور از انتظار نیست). باید با استفاده از ساختارهای نحوی زبان جملات را مشخص کرد. لیست زیر تعدادی از این موارد را نشان می دهد..

1. در آدرس های وب علائم «.» به کار می رود.
2. علائم اختصاری مانند «ق.م» به معنای قبل از میلاد که البته منجر به اشکال در تعیین مرز جملات نیز می شوند.

D.J.& J. H.Martin در کتاب خود پیشنهاد می‌کند که برای حل مسئله ابهام نقطه پایان یا نقطه درون یک کلمه از کلاس‌بند دودویی¹ به عنوان یک راه حل رایج استفاده شود[13]. به این ترتیب که برای تمامی نقطه‌های یک متن بردار خصوصیات تعیین و سپس با استفاده از یک سیستم یادگیری ماشین با آموزش تفاوت نقاط پایان جمله با سایر نقاط این تعیین تفاوت را به یک کلاس‌بند تشخیص دهنده دو کلاس مذکور² و اگزار می‌گردد. همچنین الگوریتم ساده‌ای را نیز به عنوان یک راه حل با صحت عملکرد نسبتاً مناسب برای زبان انگلیسی ارائه کرده است که در بخش فرایندهای پیش‌پردازش با اعمال یکسری تغییرات ساده برای زبان فارسی آن را به کار گرفته ایم.

2-2-3-2 تعیین مرز کلمات

کلمه به عنوان کوچکترین واحد معنادار متن در اکثر تحقیقات به عنوان پایه ارزیابی‌ها و نتیجه گیری‌ها قرار می‌گیرد. به طور مرسوم کلمات با استفاده از علائم نوشتاری از قبیل فضای خالی، علامت پرش، علامت خط جدید، «، «، «<»، «، «[»، «-»، «_» و «/». شناسایی شده و از یکدیگر تفکیک می‌شوند. اما در زبان نوشتاری فارسی به علت ویژگی‌های خاصی که وجود دارد، تنها به این علائم تکیه کردن کافی نیست چراکه در این صورت ممکن است بخش‌هایی از یک کلمه، به عنوان کلمه‌ای مستقل و ترکیب چند کلمه مستقل، به عنوان یک کلمه واحد در خروجی بدست آید. لیستی از مواردی که در این تحقیق با آن مواجه شده‌ایم بصورت زیر است:

1. کلماتی شامل حروف «ر»، «ژ»، «ز»، «د»، «ذ»، «ا»، «و»

از آنجا که این حروف از سمت چپ به حرف دیگری نمی‌چسبند در دو حالت منجر به

اشتباهاتی توسط نگارنده می‌شوند:

¹ Binary classifier

² نقاط پایان جمله و کلاس نقاط با کابردی غیر از پایان جمله

أ. حضور در انتهای کلمه: به خاطر فاصله‌ای که بصورت ذاتی با کلمه بعد ایجاد می‌کنند گاهی باعث می‌شوند نگارنده تایپ فاصله بین کلمات را از یاد ببرد.

به عنوان مثال عبارت: «اصرار بر جستجو»، «اصراربرجستجو» اگرچه در هر دو حالت به راحتی خوانده می‌شود اما حالت دوم از دید یک پردازشگر متنهای ۱۲ حرفی محسوب می‌شود.

ب. حضور در میان کلمه: منجر به عدم تشخیص فاصله اشتباهات تایپ شده میان کلمات توسط نگارنده می‌شوند.

به عنوان مثال: «تابناک»، «تابنا ک»

2. افعال و کلمات چند بخشی

زبان فارسی مملو از کلمات و ترکیباتیست که با وجود یکپارچگی در معنا با استفاده از علامت فاصله به چند بخش تقسیم می‌شوند در ادامه برخی از این نوع کلمات فهرست شده است:

- أ. کلمه‌هایی که دارای صفت جمع مثل «ها»، «ان» و غیره هستند.
- ب. افعال مضارع که از کلمه «می» استفاده می‌کنند.
- ج. افعال مرکب از قبیل «آتش زدن»
- د. کلمات چند بخشی از قبیل «ساده زیست»
- ه. کلمات مرکبی که با یک دیگر هم آیی دارند: از قبیل «امام خمینی» یا «محمود فرشچیان».

3. نشانه‌ی دستوری / و «»

اگرچه نشانه دستوری «/» در زبان فارسی به عنوان کلمه جانشین به کار می‌رود و دو کلمه را جایگزین هم معرفی می‌کند اما در موارد دیگر معنای متفاوتی نیز دارد. (مثلاً به جای عبارت چه مرد چه زن در یک فرم شکل مردانه به کار برد می‌شود). به عنوان نمونه:

أ. تاریخ در فارسی با علامت «/» مشخص می‌شود.

ب. در آدرس‌های کامپیوتری علامت «/» به کار می‌رود.

مشکل «،» در مورد اعداد در صورتی که ارقام آن‌ها با علامت «،» از هم جدا شود اتفاق می‌افتد مانند 999,222 در اینصورت تقسیم این عدد به دو عدد کوچتر معنای عدد را بهم می‌ریزد

برای حل مشکل واژه‌های مرکب روش‌های مختلفی پنهانهاد شده است. در مرجع [14] یک روش آماری مبتنی بر پیکره برای تعیین واژه‌هایی که باهم‌آبی دارند، معرفی شده است. در این روش، با توجه به فراوانی رخداد دو واژه به صورت متوالی و هم چنین فراوانی هر کدام به تنها، تصمیم‌گیری می‌شود که آیا دو واژه باهم‌آبی معنی‌داری دارند یا خیر. باید توجه داشت که باهم‌آبی دو واژه لزوماً به معنای آن نیست که آن دو واژه با هم یک واژه مرکب را تشکیل می‌دهند. تجربه نشان می‌دهد که در بیشتر کاربردهای پردازش متن، در نظر گرفتن واژه‌هایی که باهم‌آبی دارند به عنوان یک واژه واحد، مفید واقع می‌شود.

به طور کلی از فعالیت‌های انجام شده در زمینه تعیین مرز کلمات و گروه‌های کلمات می‌توان به مرجع [15] اشاره نمود که به تشخیص انتهای کلمات و فاصله گذاری میان آن‌ها می‌پردازد. هم چنین مرجع [16] در مطالعه‌ای به بررسی نحوه‌ی تشخیص کسره‌ی اضافه در متن با استفاده از روش‌های آماری مبتنی بر پیکره‌های زبانی پرداخته است. تشخیص کسره‌ی اضافه محدود کمک بسیاری به حل مشکل ابهام در شناسایی مرز گروه‌های اسمی می‌نماید. از سوی دیگر در بعضی کارها این مرحله با مراحل دیگر

ادغام و یا در طی مراحل دیگر انجام می شود. مثلا در برخی تحلیل گر ساخت واژی معرفی شده و در حین تحلیل ساخت واژی، مرز کلمات نیز تعیین و فاصله های زائد درون کلمات حذف می شوند.

مقاله [17] با استفاده از یک روش آماری مبتنی بر پیکره برای شناسایی واژه های به هم چسبیده معرفی شده است. این روش، با استفاده از معیار احتمال شرطی متقارن تعیین می کند که آیا یک عنصر متنی، یک واژه ای کامل است یا مجموعه ای از واژه های به هم چسبیده است. روش پیشنهادی روی یک مجموعه ای 18000 تایی از اخبار ورزشی مورد آزمایش قرار گرفته و تعداد نمونه هایی که به درستی از هم جدا شدند به کل نمونه های جداد شده نسبت قابل قبولی (بیش از 80 درصد در مورد نمونه های جداد شده به دو بخش) را نشان داده است.

3-3-2 حذف کلمات پر تکرار

کلمات پر تکرار در یک فرایند متن کاوی به کلماتی تلقی می شود که در تمامی متن ها به کرات یافت می شوند و معمولاً معنی مستقل خاصی نداشته یا برای ایجاد قواعد دستوری به کار می روند. به طور مثال کلماتی مثل "را" "در" "به" یا افعالی مثل "است" و "می شود" یا در زبان انگلیسی کلماتی شبیه به

A, ABOVE, ACROSS, AFTER, AGAIN, AGAINST, ALL, ALMOST, ALONE

را می توان جزو کلمات پر تکرار محسوب کرد.

عدم حذف کلمات پر تکرار در یک متن به علت فرکانس بالای این کلمات منجر به کم اثر شدن کلمات غیرمعمول با فرکانس نسبتاً خوب در متن می شود. از طرف دیگر با توجه به این نکته که کلمات مذکور فارغ از زبان مورد بررسی ، حجم بزرگی از یک متن را تشکیل می دهند حذف این کلمات تاثیر مقید زیادی بر کارایی سیستم چه از نظر پچیدگی های زمانی و چه از نظر حافظه داشته باشد [18].

از جمله کارهایی که در زمینه حذف کلمات پر تکرار شده است می‌توان به کار تقوا [19] که در آن با اشاره به انواع فعل‌های بی ارزش مثل «بود» پرداخته است یا مرجع [20] که با استفاده از روش‌های آماری و تحلیلی لیستی 928 کلمه‌ای از کلمات پر تکرار فارسی را ارائه داده است اشاره کرد. همچنین مقاله ذکر شده در مرجع [21] به این مسئله پرداخته است که حذف کلمات پر تکرار فارغ از دامنه می‌تواند مشکلاتی را ایجاد کند. این مشکل در مرجع [20] نیز به نوعی بررسی شده است و نشان داده شده که مثلاً در تحقیقی که برروی پایگاه داده همشهری انجام شد کلمه با ارزشی مثل «ایران» نیز در صورت نگاه آماری صرف جزو کلمات پر تکرار تلقی می‌شود.

در پایان شایان ذکر است که البته گاهی بعضی استثنایات در مسئله حذف کلمات کلیدی مشکلاتی را به وجود می‌آورد. مثلاً تا چند سال قبل اگر جمله مشهور شکسپیر "be or not to be" را در موتورهای جستجو وارد می‌کردید نتیجه جستجو این بود "لغات استفاده شده شما به علت تکرار زیاد، توسط موتور جستجو حذف شده است" و این نشان می‌دهد که برخورد فارغ از معنی با کلمات پر تکرار می‌تواند مشکل ساز باشد اگرچه در اکثر تحقیقات این طور استثنایات قربانی افزودن بر سرعت و کمتر شدن پیچیدگی خواهند شد اما در این پروژه همانگونه که در بخش 3-5 خواهیم دید با اضافه کردن این قید که «حذف کلمات پر تکرار به شرط عدم حضور در یک توالی بیش از سه کلمه‌ای از کلمات پر تکرار، قابل انجام است.» از استثنایاتی به این شکل تا حد زیادی جلوگیری شده است.

4-3-2 ریشه یابی کلمات

ریشه یابی یکی از کارهای رایج در استخراج کلمات کلیدی است چرا که باعث انتخاب کلماتی می‌شود که به اشکال متفاوت تکرار زیادی داشته‌اند. هدف از ریشه یابی حذف اضافات از کلمه و رسیدن به ریشه ای

یکسان برای کلمات هم ریشه است. اکثر ریشه یاب هایی که برای فارسی پیشنهاد شده از روش های مبتنی بر حذف پسوندها و پیشوندها استفاده می کنند. در این قسمت ابتدا به بررسی مشکلات ریشه یابی در زبان فارسی و همچنین راه حل های محتمل برای آنان می پردازیم سپس به برخی از کارهایی که تاکنون در زمینه ریشه یابی انجام شده اند اشاره می کنیم.

1-4-3-2 دشواری های ریشه یابی در زبان فارسی

1. مشکلات ناشی از خطاهای یکسان سازی

• مشکلات ناشی از مسائلی که استاندارد یکسانی ندارند مانند «ها»ی نشانه جمع یا «ی»

(بعد از «ه») (مانند کتیبه هی)

• مشکلات ناشی از فاصله گذاری ها که در تعیین مرز جملات و کلمات نیز به آن اشاره

شد.

واضح است که برای حل این مشکلات باید هرچه بیشتر یکسان سازی را دقیقتر و جامع تر انجام

داد.

2. کلمات مرکب: در زمینه کلمات مرکب (به ویژه فعل مرکب که تاثیر بیشتری در فرایندهای ریشه-

یابی دارد.) با وجود کوشش های فراوان هنوز ابهامات و ناهمانگی های فراوانی وجود دارد که به

نظر می رسد برای رفع آن لازم است فهرست کاملی از افعال مرکب به همراه شکل پیشنهادی آن

از سوی فرهنگستان ارائه شود.

3. دگرگونی در کلمه ها هنگام پیوند:

برخی از کلمات فارسی در هنگام ادغام با پسوندها و پیشوندها دچار تغییر شکل می شوند. در

زمان ریشه یابی لازم است که برای این حالات (یا حداقل برای حالاتی که بیشتر دیده می شوند)

نیز راه حلی ارائه شود یا بصورت لیستی از استثنایات در برنامه ریشه یاب منظور شوند.

جدول 2-1 نمونه ای از این دگرگونی ها را نمایش می دهد.

جدول 2-1 نمونه ای از دگرگونی کلمات در هنگام پیوند

زنده + م ← زنده ام	زنده + ان ← زندگان
ن + افتاد ← نیفتاد	گو+م ← گویم
زنده+ها ← زنده ها	ن + آزمایش ← نیازما

4. کلمات زبان های بیگانه با ساختار ظاهری شبیه به کلمات فارسی

در مجموع واژگان فارسی کلمات زیادی از زبان های دیگر وجود دارد که ساختار ظاهری آنان شبیه به کلمات فارسی است. مثلا کلمه رئالیست در یک نگاه رشته ای صرف می تواند بصورت ساختار «رئالی + ست» دیده شود که در این صورت با یک ریشه یابی اشتباه حروف «ست» که بخش از اصل کلمه هستند به عنوان پسوند حذف خواهند شد.

5. شناسایی ریشه فعل ها

یکی از ابتدایی ترین مشکلات در زمینه شناسایی ریشه فعل ها خود مسئله شناسایی فعل است. در واقع در زبان فارسی فرایند تعیین بخش های سخن¹ فرایند مشکلی است. به عنوان مثال قواعدی مثل «فعل شناسه می گیرد» نمی تواند پاسخ گوی تشخیص نیاز به تشخیص فعل ها از اسمی باشد چراکه گاهی اسمی و صفات نیز شناسه می گیرند. به عنوان مثال کلمه «خوشحالیم» به جای «خوشحال هستیم» به کار گرفته می شود.

¹ Part of speech

«در فعل های ساده پس از حذف «ن» از مصدر، بن ماضی باقی می مند و از جهت تغییری که از

بن ماضی به مضارع انجام می گیرد، آنها را می توان در هشت گروه جای داد. جدول 2-2

تغییرهای فعل ها در گروه های هشت گانه نشان می دهد.» [22]

تعداد کمی از فعل های فارسی از قاعده های این هشت گروه پیروی نمی کنند. بن گذشته و

غیر گذشته این فعلها جداگانه نوشته شدند. مرجع [23] با نقل جدول 2-2 ذکر می کند که این

روش بر مجموعه بزرگی از کلمات به کارگرفته شده و بن گذشته و غیر گذشته بیشتر فعل های

ساده‌ی فارسی به خوبی شناسایی گردیده اند.

جدول 2-2 هشت گروه فعل های فارسی

مثال	نحوه ساخت بن مضارع	۱. گروه پیش از گذشته	۲. گروه پس از گذشته	۳. گروه
نالیدن نال نالید	نال + ید + ن بن مضارع است پس از حذف «ید» باقیمانده	ید	نالیدن بن مضارع است پس از حذف «ن» باقیمانده	1
خوردن خور خورد	خور + د + ن مضارع است پس از حذف «د» باقیمانده بن	د	خوردن بن مضارع است پس از حذف «ن» باقیمانده	2

آزمودن					
آزما	آزمود - ن	پس از حذف «و»، «و» به «ا» تبدیل می شود	- ود	ودن	3
آزمود					
افتادن					
افت	افت + اد - ن	پس از حذف «اد» باقیمانده بین مضارع خواهد بود	- اد	ادن	4
افتاد					
ساختن					
ساز	ساخ + ت - ن	«ت» حذف و «خ» تبدیل به «ز» می شود	- خت	ختن	5
ساخت					
آراستن					
آرا	آرا + ست + ن	پس از حذف «ست» باقیمانده بن مضارع است	- ست	ستن	6
آراست					
کاشتن					
کار	کاش + ت - ن	پس از حذف «ت»، «ش» تبدیل به «ر» می شود	- شت	شتن	7
کاش					

تافتن					
تاب	تاف + ت - ن	پس از حذف «ت»، «ف» تبديل به «ب» می شود	- فت	فتون	8
تافت					

2-4-3-2 روش های ریشه یابی

روش های ریشه یابی را می توان به دو دسته کلی تقسیم کرد:

روش های ساختاری که با توجه به ساختار کلمات و قواعد زبان کار می کنند. این روش ها با فرمالیزه-

کردن قواعد زبان امکان استفاده از ماشین برای تشخیص ریشه یک کلمه در زبان را فراهم می کنند.

معمولاً در این روش ها یک ماشین پذیرنده متناهی¹ ساخته می شود و وندهای مختلف مسیرهای موجود

در این ماشین را ایجاد می کنند. دو اشکال بزرگ برای این روش ها وجود دارد

1. روش های ساختاری نیازمند دانش کافی در مورد ساختار زبان است.

2. روش های ساختاری در صورت نیاز به تغییر نیاز به کد نویسی مجدد دارند و معمولاً در این روش

ها نگهداری و به گسترش برنامه هزینه بر خواهد بود

روش های غیر ساختاری: در این روش ها یک مجموعه بزرگ از کلمات با انواع ساخت ها جمع آوری می -

شود و واضح است که عملکرد ریشه یاب وابستگی مستقیم به جامعیت این مجموعه کلمات دارد. در این

روش با استفاده از تحلیل های آماری وندهایی که در کلمات تکرار شده اند شناسایی می شوند. در این

روش ها با استفاده از دو مرحله آموزش و ریشه یابی با استفاده از مرحله آموزش سعی می شود از برخورد

¹ DFA: Deterministic Finite Accepters

مستقیم با قواعد زبانی جلوگیری شود. این روش‌ها به علت غیر ساختاری بودن و وابستگی بسیار کم به دانش زبانی در سیستم‌های چند زبانه کاربرد بهتری دارند همچنین بطور عمومی با توجه به پیچیدگی‌های طراحی کمتر از نتایج مطلوبی برخوردارند.

اما روش‌های غیر ساختاری از دو مشکل بزرگ رنج می‌برند:

2. با وجود پیچیدگی‌های کم تر طراحی پیچیدگی زمانی و حافظه‌ای بیشتری دارند

3. مجموعه کلمات فارسی برای این روش‌ها وجود ندارد

با توجه به اینکه ریشه‌یابی فرایندی است که در هر دو روش تا حدودی از زبان مورد تحقیق تاثیر می‌پذیرد. کارهای انجام شده در این زمینه را به دو بخش زبان انگلیسی و زبان فارسی تقسیم کرده ایم

3-4-3-2 ریشه‌یابی برای زبان انگلیسی

در راه حل‌های ساختاری الگوریتم Porter [24] معروف‌ترین الگوریتم ریشه‌یابی انگلیسی است. روش Porter بر پایه زبان‌شناسی و دسته‌بندی کلمه‌ها به کمک واژه‌ها و هجاه‌ها بنا نهاده شده است. پس از آن وندهای کلمات درون مجموعه کلمات بطور خودکار برداشته می‌شوند. از الگوریتم‌های خوب دیگر در این زمینه می‌توان به کار B. Lovins [25] نیز اشاره کرد.

در میان الگوریتم‌های غیرساختاری که به نوعی الگوریتم‌های ریشه‌یابی جدیدی محسوب می‌شوند الگوریتم Bacchne [26][27] را می‌توان یکی از کاملترین الگوریتم‌ها بر شمرد.

4-3-2 ریشه یابی برای فارسی

از آنجا که تفسیر علمی گرامر و قواعد فارسی قدمت چندانی ندارد، سیستم های ریشه یابی فارسی نیز معمولاً کارهایی تازه هستند و با توجه که بیشتر در روش های ساختاری به زبان فارسی پرداخته شده است اشکالاتی نیز در این تحقیقات وجود دارد.

ریشه یاب «بن» [28] الگوریتمی برای ریشه یابی در زبان فارسی با استفاده از الگوریتمی شبیه به الگوریتم Porter ریشه یاب ارائه کرده است. در این الگوریتم پایان کلمات به عنوان وند در نظر گرفته شده است و در یک جدول که در حقیقت کار یک DFA¹ را بر عهده دارد به جستجوی آن می پردازد و پس از یافتن پسوند احتمالی کلمه آن را حذف می کند. در صورتیکه پایان یک کلمه با دو پسوند همخوانی داشته باشد پسوند بزرگتر به عنوان پسوند انتخاب و حذف می شود. همچنین در این تحقیق حد سه نویسه برای حداقل طول کلمات بدون پسوند در نظر گرفته شده است تا از تخریب کلماتی که حروف پایانی شبیه به پسوند دارند جلوگیری کند. (به عنوان مثال کلمه «تاوان» در صورت حذف «ان» به عنوان پسوند جمع به کلمه ای بی معنی تبدیل می شود).

تقوا [4] در روشنی کاملاً یکسان با [28] به طراحی یک DFA برای ریشه یابی کلمات فارسی پرداخته است. DFA طراحی شده در این تحقیق از 70 حالت برای حل مسئله کمک می گیرد و بسیاری از پسوندها را شامل می شود.

از دیگر تحقیقات ریشه یابی کلمات که پایه آن DFA است می توان به مرجع [29] اشاره کرد.

اما در مرجع [30] در روشنی تازه با استفاده از الگوریتمی که بر پایه الگوریتم Bacchne بنا نهاده شده است به ریشه یابی کلمات فارسی با استفاده از روش های آماری پرداخته است. مجموعه مستندات برای

¹ Deterministic Finite Accepters

آموزش در این سیستم از سایت خبرگزاری دانشجویان ایران ISNA انتخاب شده است. سیستم ارائه شده در این مقاله با بررسی این مستندات به عنوان مجموعه آموزش لیستی از اشتقاق‌ها و پسوندها و پیشوندها را برای خود می‌سازد و به هر کدام از آن‌ها یک وزن نسبت می‌دهد. در هنگام ریشه‌یابی برای یک کلمه حالت مختلفی از تقسیم آن کلمه در نظر گرفته می‌شود و با توجه به لیست‌های گفته شده بهترین ریشه برای آن استخراج می‌شود. نتایج الگوریتم طبق ادعای نویسنده با نتایج سایر روش‌ها برابری می‌کند.

4-2 روش‌های استخراج کلمات کلیدی

در این بخش از فصل دوم به بررسی روش‌های استخراج کلمات کلیدی می‌پردازیم.

1-4-2 تقسیم‌بندی تکنیکی روش‌ها

از نظر تکنیکی روش‌های استخراج کلمات کلیدی را می‌توان به چهار بخش تقسیم کرد:

1-1-4-2 روش‌های آماری

این روش‌ها نیازی به مجموعه داده‌های آموزش ندارند و فقط با استفاده از اطلاعات آماری کلمات به تشخیص کلمات کلیدی در یک سند می‌پردازند.

مرجع [31] به عنوان یک روش آماری با استفاده از N-gram‌ها به شاخص گذاری اتوماتیک پرداخته است. N-gram‌ها یک زیرتوالی n عنصری از یک توالی داده شده محسوب می‌شوند که بسته به کاربرد عناصر زیر توالی درخواست شده می‌توانند پدیده‌های، حروف، واژه‌ها، یا کلمات باشند. مدل n-gram یک مدل احتمالی برای پیش‌بینی عنصر بعدی در یک توالی به شمار می‌رود که در تحلیل آماری NLP و

توالی های زنگنه کاربرد فراوان دارد. اگر سایز n برابر با 1 باشد Unigram، برابر با 2 Bigram سه عنصری trigram و برای سایر حالات همان n-gram خوانده می شود. بعضی از مدل های زبانی ساخته شده از روی n-gram به عنوان مرتبه 1 n-1 مدل مارکوف شناخته می شوند.

مرجع [32] با استفاده از همآیی¹ کلمات یا اصطلاحاً کلمات هم رخداد عمل استخراج کلمات کلیدی را انجام داده است.

همچنین مرجع [33] از PAT-tree² برای استخراج کلمات کلیدی بهره گرفته است. PAT-tree یک ساختمان داده بسیار قدرتمند است که بطور گسترده در تحقیقات بازیابی اطلاعات مورد استفاده قرار گرفته است. به طور مفهومی این ساختمان داده با یک درخت جستجوی رقمی³ برابر می کند اما نوعاً ساختار ساده تری نسبت به یک درخت جستجوی رقمی دارد.

2-1-4-2 روش های مبتنی بر ساختارهای زبانی

این روش ها کلمات کلیدی را بر استفاده از خصوصیات زبانی آن ها انتخاب می کنند. این روش ها شامل سه دسته لغوی مانند مقاله ارائه شده در مرجع [34]، نحوی مانند مرجع [35] و آنالیز معنا محور مانند مرجع [36] می باشند.

3-1-4-2 روش های برپایه یادگیری ماشین

این روش ها با استفاده از کلمات کلیدی استخراج شده توسط انسان به ایجاد یک مدل می پردازند و از این مدل ایجاد شده برای استخراج کلمات کلیدی در اسناد جدید بهره می گیرد.

¹ Co-occurrence

² PATrica tree که به نام suffix tree (درخت پسوندی) نیز شناخته می شود.

³ Digital search tree

این روش‌ها به توجه به فراوانی و اهمیت یادگیری ماشینی در هوش مصنوعی به طرز گستردۀ ای مورد استفاده قرار گرفته‌اند. از این میان می‌توان به استفاده از قانون بیز^۱ [37]، ماشین‌های بردار پشتیبان^۲ [38] بقچه کلمات^۳ [35] اشاره کرد. همچنین برخی از ابزار‌های استخراج کلمات کلیدی مانند GenEx^۴ [39] و KEA^۵ [40] نیز با استفاده از روش‌های این دسته ساخته شده‌اند.

4-1-4-2 روش‌های ترکیبی

روش‌های ترکیبی با استفاده از ترکیب روش‌های ذکر شده قبلی بدست آمده‌اند. در واقع روش‌هایی که با استفاده از توابع اکتشافی بر اساس مکان و طول و خصوصیات نمایشی کلمات کلیدی یا با استفاده از تگ‌های HTML [41] به استخراج کلمات کلیدی دست می‌زنند نیز در این دسته رده بندی می‌شود.

در پایان باید گفت به طور کلی در اکثر روش‌های معروف تعداد کلمات استخراج شده به عنوان کلمه کلیدی 10 الی 15 کلمه می‌باشد. همچنین اکثر روش‌های استخراج کلمات کلیدی مبتنی بر پردازش زبان طبیعی^۶ از دیکشنری برای مشخص کردن ریشه کلمات و بخش‌های گفتار استفاده می‌کنند.

2-4-2 روش‌های امتیاز دهی به کلمات

وزن دهی به کلمات یکی از مهمترین مراحل استخراج کلمات کلیدی است. در این مرحله به هریک از کلمات موجود در لیست بدست آمده از مرحله پیش‌پردازش، امتیازی داده می‌شود که میزان کلمه بودن کلمه را تعیین می‌کند. سپس با اعمال یک آستانه که می‌تواند براساس امتیاز کلمات یا تعداد مشخص کلمات در رتبه‌های مختلف باشد کلمات کلیدی انتخاب می‌شوند.

¹ Naïve Bayes

² SVM (support vector machine)

³ Bagging

⁴ Keyword Extraction Algorithim

⁵ نام روش از ترکیب دو عبارت extraction و genetic algoritm گرفته شده است.

⁶ NLP (Natural Language Processing)

1-2-4-2 وزن دهی غیر آماری

در این روش وزن دهی به کلمات بر اساس پارامترهای انسانی که اهمیت یک کلمه در متن را مشخص می‌کند صورت می‌گیرد. لیست زیر به نمونه‌هایی از این پارامترها اشاره می‌کند [36].

- مکان قرارگیری کلمه در متن (چیکده، عنوان، و ...)
- مفهوم هر کلمه، که بیانگر ارتباط کلمه با کلمه‌های دیگر است. (متراff، متضاد)
- کاربرد خاص کلمه: مثلاً اسمی خاص دارای بار معنایی بیشتری هستند.

2-2-4-2 وزن دهی آماری

بر عکس روش‌های وزن دهی غیر آماری که تکیه بر متون ساخت یافته یا لغتنامه‌ها دارند، روش‌های وزن دهی آماری تنها بر حضور کلمات تکیه دارند درنتیجه به نوعی مستقل از زبان محسوب شده و می‌توانند در سطح وسیعی از اسناد مورد استفاده قرار گیرند.

وزن دهی آماری بر دو پایه فراوانی مطلق که وزن دهی را بر اساس تعداد حضور کلمه در یک سند منفرد و فراوانی نسبی که ارزش دهی را براساس حضور کلمه در اسناد مختلف انجام می‌دهد، بنا نهاده شده است. روش‌های فراوانی مطلق در سیستم‌های برخط¹ که نمی‌تواند به حجم مقبولی محدود شود کارایی بیشتری دارد.

روش‌های زیادی برای وزن دهی آماری تعریف شده که برخی از آن‌ها مرور می‌کنیم:

روش وزن دهی² TFIDF

¹ Online

² Term Frequency * Invert Document Frequency

این روش به خاطر سادگی محاسبات و نیز نتایج قابل قبول یکی از پرکاربردترین روش های وزن دهی

آماری به کلمات بشمار می رود[42]. برای کلمه i در سند k ام TFIDF بصورت رابطه تعريف می شود:

$$TFIDF(w_{ik}) = freq(w_{ik}, D_k) * \log_2 \left(\frac{N}{n} + 1 \right) \quad 1-2$$

که در آن

$freq(w_{ik}, D_k)$ به معنای تعداد تکرار کلمه i در سند k

N : تعداد کل اسناد

n : تعداد اسنادی که کلمه i در آن وجود دارد.

سیگنال و نویز

در این روش هرچه احتمال دیده شدن کلمه بیشتر باشد بار اطلاعاتی کمتری برایش در نظر گرفته می-

شود. از آنجا که کلمات با ارزش تر در تعداد سند کمتری ظاهر شده اند نویز کمتری هم دارند.

نویز را بصورت رابطه 2-2 تعريف می کنیم:

$$NOISE_k = \sum_{i=1}^n \frac{FREQ_{ik}}{TOTFREQ_k} * \log_2 \frac{TOTFREQ_k}{FREQ_{ik}} \quad 2-2$$

پارامتر سیگنال به گونه ای تعريف می شود که ارزش بیشتر را با معکوس نویز نشان دهد. (رابطه 3-2)

$$signal_k = \log_2 (TOTFREQ_k) - NOISE_k \quad 3-2$$

حالا وزن کلمه k در سند k ام بصورت 4-2 محاسبه خواهد شد.

$$WEIGHT_{ik} = FREQ_{ik} * signal_k \quad 4-2$$

میزان ایجاد تمایز

در این پارامتر کلماتی که موجب بیشترین تمایز متن با سایر متون می شوند به عنوان کلمات کلیدی برچسب می خورند. میزان تمایز بر اساس معکوس میزان ایجاد شباهت، که معیارهای گوناگونی برای آن تعریف شده است بدست می آید. باهم به بررسی معیارهای شباهت می پردازیم

فرض کنید دو متن را بصورت دو بردار

$$\begin{aligned} X &= (x_1, x_2, \dots, x_t) \\ Y &= (y_1, y_2, \dots, y_t) \end{aligned}$$

در نظر می گیریم. که درایه های x_i و y_j فراوانی کلمات در استناد مربوطه باشد. در این صورتتابع متوسط مشابهت (رابطه 5-2) بیان کننده میزان شباهت کل استناد مجموعه به یکدیگر است. این تابع هر استناد به یکدیگر شبیه تر باشند مقدار بیشتری می گیرد.

$$AVGSIM = \frac{1}{n(n-1)} * \sum_{i=1}^n \sum_{j=1}^n similarity(D_i, D_j) \quad i \neq j \quad 5-2$$

آنچه که در این تابع به نام شباهت¹ تعریف شده است تابعی است که میزان مشابهت دو سند به یکدیگر را مشخص می کند. برای این تابع تعاریف مختلفی ارائه شده است که در جدول 3-2 تعدادی از آن ها را مشاهده می کنید.

¹ similarity

جدول 2-3 تعدادی از توابع شباهت مرسوم

$similarity(X, Y)$	معیار مشابهت
$\sum_{i=1}^t x_i y_i$	ضرب داخلی
$2 * \frac{\sum_{i=1}^t x_i y_i}{\sum_{i=1}^t x_i^2 + \sum_{i=1}^t y_i^2}$	ضریب Dice
$\frac{\sum_{i=1}^t x_i y_i}{\sqrt{\sum_{i=1}^t x_i^2} * \sqrt{\sum_{i=1}^t y_i^2}}$	ضریب کسینوسی
$\frac{\sum_{i=1}^t x_i y_i}{\sum_{i=1}^t x_i^2 + \sum_{i=1}^t y_i^2 - \sum_{i=1}^t x_i y_i}$	ضریب jaccard
$\frac{\sum_{i=1}^t x_i y_i}{\min\left(\sum_{i=1}^t x_i^2, \sum_{i=1}^t y_i^2\right)}$	ضریب همپوشانی

میزان تمایز برای کلمه K را بصورت $AVGSIM_k$ تعریف می کنیم که به معنای متوسط مشابهت در زمان عدم حضور کلمه در کلیه اسناد است. اگر K کلمه ای عمومی باشد آنگاه حذف آن موجب تمایز بیشتر اسناد با یکدیگر شده و $AVGSIM_k$ کاهش می یابد اما اگر k یک کلمه خاص و در اصطلاح متمايز

کننده باشد حذف آن منجر به شباهت بیشتر اسناد به هم شده و $A VGSIM$ هم افزایش می‌یابد با استفاده از تعاریف ذکر شده مقدار تمايز و وزن کلمه K در سند i با روابط 6-2 و 7-2 محاسبه می‌شود.

$$Discvalue_k = (AVGSIM)_k - AVGSIM \quad 6-2$$

$$WEIGHT_{ik} = FREQ_{ik} * Discvalue_k \quad 7-2$$

3-4-3 پارامترهای ارزیابی کلمات استخراج شده

کلمات استخراج شده از یک متن توسط یک روش بازیابی اطلاعات را از دو منظر می‌توان ارزیابی کرد:

- در برگیری به این معنا که نشان می‌دهد چه مقدار از مفاهیم یک متن ممکن است با کلمات ارائه شده بازیابی شود افزایش تعداد کلمات مناسب این ملاک را افزایش می‌دهد.
- تعیین کننده: به معنای میزان دقیق بودن کلمات است و به این معنا که کلمات با ارتباط کمتر با اصل متن کمتر ارائه شده باشند..

برای ارزیابی کلمات به دو شکل عمل می‌شود:

1. داوری مبتنی بر کارشناس انسانی:

أ. نتایج بر اساس کلمات کلیدی استخراج شده توسط یک کارشناس ارزیابی شود.

این نوع ارزیابی بسیار رایج است و در این ارزیابی نتایج نهایی این پروژه نیز به همین

روش انجام گرفته است اما اشکالات زیر را می‌توان بر این نوع ارزیابی وارد دانست.

1. کلمات کارشناس همیشه از روی خود متن استخراج نمی‌شوند به این معنا که

یک کارشناس گاهی کلماتی را در لیست کلمات کلیدی ارائه می‌کند که

مفهومی از متن را در خود دارد اما این کلمه مستقیماً در خود متن آورده نشده

است. این نوع کلمات به هیچ وجه نمی تواند با استفاده از روش های استخراج

کلمات کلیدی رایج استخراج شود.

2. معمولاً به پایگاه داده‌ای که برای متون آن کلمات کلیدی استخراج شده باشد و

در عین حال از جامعیت موضوعی برخوردار باشد دسترسی کمتری وجود دارد

در واقع بیشتر پایگاه داده‌هایی که برای متون آن کلمات کلیدی استخراج شده

باشد محدود به اسناد رسمی از قبیل صورت جلسه‌ها و یا اسنادی مانند مقالات

علمی باشد.

3. همانطور که در بخش قبل گفته شد کلمات کلیدی برای متونی خاص با منظور

خاص استفاده می شود بهمین خاطر معمولاً از نویسنده خواسته می شود تعداد

کلمات استخراجی را محدود نگه دارد.

ب. نتایج به یک کارشناس برای ارزیابی سپرده می شود.:

در این حالت کارشناس مقادیر را از نظر ارزشی به کلمات انتخاب شده می دهد به عنوان

مثال (خوب ، متوسط ، بد). این نوع داوری اگرچه مشکلات حالت قبلی را ندارد اما

هزینه‌بر و زمانبر بوده و نیز می تواند دچار عدم عدالت در داوری باشد.

2. داوری مبتنی بر سیستم‌های بازیابی اطلاعات:

در این حالت برای ارزیابی کلمات کلیدی، باید یک مجموعه از متون به همراه مجموعه ای از

پرس و جوهای مرتبط با اسناد و هم چنین تمام اسناد مرتبط با پرس و جوها وجود داشته باشد .

پرس و جوها به مجموعه اعمال می شوند، و اسناد با توجه به مشابهت پرس و جو و اسناد ،

بازیابی می شوند.

در حالت کلی نتایج بدست آمده از داوری ها، افزارهای گوناگونی را به وجود می‌آورد که براساس این افزارها (جدول 4-2) ملاک های ارزیابی روش های استخراج کلمات کلیدی تعریف می شوند.

جدول 4-2 افزارهای ممکن برای مجموعه اسناد

	نتایج بدست آمده	نتایج بدست نیامده	
پاسخ های بدست آمده	$A \cap B$	$\bar{A} \cap B$	B
پاسخ های بدست نیامده	$A \cap \bar{B}$	$\bar{A} \cap \bar{B}$	\bar{B}
	A	\bar{A}	

در جدول 4-2 افزار اول مجموعه نتایج را به دو حالت نتایج بدست آمده از داوری (A) و نتایج بدست نیامده از داوری (\bar{A})، و افزار دوم نتایج را به دو حالت پاسخ های بدست آمده از سیستم (B) و پاسخ های ایجاد نشده در سیستم (\bar{B}) تقسیم می کند

پارامترهایی که برای ارزیابی نتایج سیستم در تحقیقات به کار گرفته می شوند بصورت زیر از جدول 4-2 استخراج می شوند (توجه کنید که در رابطه های ذکر شده تعداد عناصر موجود در هر مجموعه به کار رفته است) :

1-3-4-2 میزان بازخوانی

پارامتر میزان بازخوانی¹ تعیین می کند که چه مقدار از نتایجی که سیستم ارزیابی تعیین کرده بوده بدست آمده است. به عنوان مثال اگر در بحث استخراج کلمات کلیدی این پارامتر تعیین می کند که چند درصد از کلمات کلیدی مورد انتظار را سیستم بدست آورده است. پارامتر میزان بازخوانی بصورت رابطه 8-2 بدست می آید. پارامتر میزان بازخوانی را می توان میزان صحت کارکرد سیستم دانست

¹ recall

$$R = \frac{A \cap B}{A} \quad 8-2$$

2-3-4-2 میزان دقต

پارامتر میزان دقت¹ تعیین می کند که سیستم طراحی شده چقدر توانایی به ارائه صرف نتایج مورد نظر کاربر دارد. به عنوان مثال پارامتر دقت در یک سیستم استخراج کلمات کلیدی تعیین می کند که چه تعداد از کلمات ارائه شده دقیقا کلماتی است که کاربر انتظار داشته است در حالیکه پارامتر بازخوانی که تعیین می کرد چه تعداد از کلمات کاربر ارائه شده است. پارامتر میزان دقت با استفاده از رابطه 9-2 بدست می آید.

$$P = \frac{A \cap B}{B} \quad 9-2$$

شایان ذکر است که پارامتر های میزان دقت و میزان بازخوانی در عمل تا حدودی با یکدیگر رابطه معکوس دارند به این معنا که در صورتیکه بخواهیم پارامتر دقت را افزایش دهیم باید مجموعه ی پاسخ سیستم را کوچک نگه داریم تا از ورود کلمات اضافه جلوگیری کنیم در حالیکه برای افزایش پارامتر میزان بازخوانی باید مجموعه ی پاسخ را افزایش دهیم تا شанс حضور بیشتری از کلمات مورد انتظار کاربر را ایجاد کرده باشیم. به همین خاطر در ارزیابی دو سیستم بازیابی اطلاعات، سیستمی بر دیگری برتری دارد که در هر دو پارامتر نتایج بهتری کسب کند.

3-3-4-2 پارامتر Fmeasure

این پارامتر به منظور در نظر گرفتن هر دو معیار بازخوانی و دقتم بصورت همزمان است . این پارامتر بصورت میانگین هارمونیک بازخوانی و دقتم تعریف می شود (رابطه 10-2) و هدف سیستم های بازیابی اطلاعات را می توان بیشینه کردن این پارامتر عنوان کرد.

¹ precision

$$F_{measure} = \frac{2 * P * R}{P + R} \quad 10-2$$

در پایان این بخش مثالی را برای تمامی این پارامترها حل می کنیم.

اگر کلمات کلیدی اعلامی یک کارشناس برای یک متن ورزشی فرضی «امیر قلعه نویی – باشگاه استقلال – قهرمانی لیگ برتر – فرهاد مجیدی» باشد و یک سیستم استخراج کلمات کلیدی بصورت اتوماتیک کلمات «امیر قلعه نویی – قهرمانی لیگ برتر – استقلال – سقوط ابومسلم مشهد – فرهاد مجیدی» را کلیدی اعلام کند در این صورت ارزیابی سیستم بصورت $precision = \frac{2}{5} = 40\%$ و $recal = \frac{2}{4} = 50\%$ داشت. خواهد بود.

همچنین با توجه به رابطه 10-2 خواهیم داشت.

$$F_{measure} = \frac{2 * 0.4 * 0.5}{0.4 + 0.5} = 0.44$$

5-2 مروری بر تحقیقات انجام شده

استخراج کلمات کلیدی به روش آماری، به استثنای برخی فرایندهای پیش پردازش تا حدی زیادی مستقل از زبان است به همین علت در ادامه تعدادی از کارهای انجام شده در سایر زبان‌ها آورده شده است. برای زبان فارسی تنها توانسته ایم از یک کار در حوزه نمایه سازی نام ببریم

روش‌های یادگیری با ناظر به علت امکان ارزیابی ملموس و نیز انطباق اصول آن با روش‌های رایج هوش مصنوعی از راهکارهای پر استفاده در استخراج کلمات کلیدی محسوب می‌شوند. این روش‌ها نیازمند تعدادی متن که کلمات کلیدی آن‌ها مشخص شده باشد، به عنوان پایگاه داده هاست. سیستم در طول

فاز آموزش سعی در یادگیری نحوه انتخاب کلمات کلیدی می کند و در فاز تست نتایج بر اساس پارامترهای بازخوانی و دقت مشخص می شود. از آنجا که روش استفاده شده در این تحقیق نیز بر این اساس پایه ریزی شده برخی از مقالاتی را که تاثیر بیشتری در این تحقیق داشته اند و برپایه این روش‌ها طراحی شده اند، مرور می کنیم.

در مرجع [3] نویسنده با استفاده از روش آماری CRF¹ و همچنین نگاه کردن به مسئله استخراج کلمات کلیدی به عنوان یک مسئله برچسب گذاری رشته‌ها به استخراج کلمات کلیدی به روش با ناظر، در متون چینی پرداخته است.

خصوصیاتی که برای امتیاز دهی کلمات استفاده می شود به نوبه خود می تواند تاثیر زیادی بر نتایج کار بگذارد. مرجع [43] یکی از مقالاتی است که برای استخراج کلمات کلیدی با ناظر به معروفی تعدادی خصوصیت پرداخته و میزان تاثیر گذاری هریک از این خصوصیات را در انتخاب کلمات کلیدی با استفاده از تست ANOVA بررسی کرده است.

زمانی که نخواهیم از یک مجموعه سند به عنوان پایگاه داده اولیه استفاده کنیم مجبوریم از خصوصیات در انتخاب کلمات کلیدی استفاده کنیم که تنها از یک سند منحصر به فرد بدست می آید. این خصوصیات می توانند مبتنی بر دانش آماری موجود در همان سند یا مبتنی بر ویژگی‌های غیر آماری مانند برچسب های HTML (در اسناد تحت وب) و ... باشد. مرجع [41] همانطور که قبلاً هم ذکر شد با استفاده از برچسب‌های HTML به استخراج کلمات کلیدی پرداخته است.

مرجع [44] با استفاده از چهار خصوصیت تعداد تکرار عبارت نرمال شده با تقسیم شدن بر تعداد تکرار پر تکرارترین کلمه (NTF)، معکوس تعداد اسناد شامل کلمه (IDF)، توزیع نرمال تعداد تکرار پاراگراف و

¹ Conditional Random Fields

اینکه آیا متن در عنوان ، تیتر یا زیرتیتر حضور دارد مسئله استخراج کلمات کلیدی را به عنوان یک مسئله تشخیص کلاس قلمداد کرده است و سپس با استفاده از یک شبکه عصبی پرسپترون چند لایه به حل مسئله پرداخته است.

مرجع [40] برای ارزیابی اینکه یک کلمه کلیدی است یا خیر؟ از یک درخت تصمیم¹ استفاده کرده است. همچنین این روش از 9 خصوصیت آماری تعداد کلمه موجود در عبارت، مکان اولین وقوع عبارت در متن، تعداد وقوع عبارت در سند، طول نسبی عبارت، و همچنین سه خصوصیت نحوی «آیا عبارت حاوی اسم خاص است؟»،«آیا کل عبارت در پایان به یک صفت ختم می شود؟» و «آیا عبارت حاوی یک فعل عمومی است؟» برای حل مسئله بهره می گیرد.

پروژه KEA [39] که پیش از این به آن اشاره شد با استفاده از سه خصوصیت TF و IDF و مکان اولین وقوع کلمه و یک کلاس‌بند بیزین تصمیم می گیرد که یک کلمه آیا کلیدی است یا نه. همچنین روش GenEx روش دیگریست که در این منبع معرفی شده و از الگوریتم ژنتیک استفاده می کند. این مقاله با بررسی شرایط توضیح می دهد که هریک از این دو روش در یک شرایط خاص نتایج بهتری نسبت به دیگری بدست می آورد.

در نمایه ساز «سینا»[42] به عنوان یک نمایه ساز متون فارسی ابتدا کلمات پر تکرار بر اساس یک لیست ثابت 180 کلمه ای حذف شده اند. سپس با استفاده از یک روش مبتنی بر حذف پسوندها و پیشوندها به فرایند ریشه یابی انجام شده است. این نمایه ساز برای فرایند استخراج کلمات کلیدی بعد از تست چهار روش TFIDF، LNU، NTC و ITN با توجه به نتایج LNU و TFIDF را برای امتیازدهی گزینش و روش خود را بر 450 متن شامل چکیده مقالات مرتبط با کامپیوتر را به عنوان پایگاه داده اعمال کرده است.

¹ Decision tree

فصل سوم:

آشنایی با اصول اولیه SEO

1-3 مقدمه

تابع ارزیابی کلمات کاندید که در فصل ششم مورد استفاده می‌گیرد بر اساس عملکردهای موتورهای جستجو طراحی شده است. برای شرح بهتر این تابع فصل جاری را به اصول بهینه سازی ساختار سایت‌ها برای دستیابی به رتبه‌های بالاتر یا اصطلاحاً SEO¹ پرداخته ایم. بیشتر مطالب این فصل بر اساس اسناد منتشر توسط موتور جستجوی گوگل به خصوص سند ذکر شده در مرجع [45] می‌باشد و مسلمان فقط نکاتی که در طراحی تابع ارزیابی کلمات کاندید لحاظ شده است بررسی شده‌اند. در واقع آنچه که در ادامه آمده است پنجره‌ای کوچک به دنیای بزرگ و پیچیده SEO محسوب می‌شود.

2-3 ساختار یک موتور جستجو

به طور خلاصه یک موتور جستجو شامل پنج قسمت است:

1. خزنه² یا کشف کننده

2. عنکبوت³ یا تارگذار

3. شاخص گذار⁴

4. پایگاه داده

5. الگوریتم رتبه‌بندی

¹ Search Engine Optimization

² Crowler

³ Spider

⁴ indexer

1-2-3 تارگذار

در یک موتور جستجو، تارگذار در واقع کار یک کاربر اتوماتیک را انجام می دهد. این بخش از موتور جستجو به صفحات مختلف سایت ها سر می زند و محتوای صفحات را در اختیار سایر بخش ها می گذارد.

2-3 خزنده

خزنده به عنوان تصمیم گیرنده برای تارگذار عمل می کند و با فرمانهایش برای تارگذار مشخص می کند که پیوندهای یک صفحه را دنبال کرده یا در همانجا متوقف می شود.

3-2-3 شاخص گذار

اطلاعات جمع آوری شده توسط تارگذار، بوسیله شاخص گذار بررسی شده و بصورت کلمات و امتیاز کلمات ذخیره می شود. در واقع شاخص گذار وظیفه تبدیل سیستم انسانی متن به یک سیستم عددی را دارد همچنین اکثر شاخص گذار ها برای کاهش حجم داده ها از کلمات پرتکرار صرفنظر می کنند به همین خاطر توصیه می شود در صورتیکه می خواهید از یک توالی کلمات پرتکرار استفاده کنید آن را در میان علامت نقل قول («» یا "") ارائه کنید

4-2-3 پایگاه داده

داده های تحلیل شده شاخص گذار به پایگاه داده فرستاده شده و در اینجا بعد از گروه بندی ذخیره می شود. هرچه پایگاه داده یک موتور جستجو بزرگتر باشد توانایی پاسخگویی و جامعیت پاسخ موتور جستجو بیشتر است. در واقع وقتی پرسشی از یک موتور جستجو مطرح می شود موتور جستجو نه کل وب! که پایگاه داده اش را برای ارسال پاسخ جستجو می کند.

5-2-3 سیستم رتبه‌بندی

زمانی که کاربر کلماتی را برای جستجو به موتور جستجوگر تحویل می‌دهد موتور جستجو صفحاتی را که به موضوع مرتبط باشد مشخص می‌کند سپس آن‌ها را به ترتیب از بیشترین ارتباط تا کمترین ارتباط مرتب می‌کند به عنوان پاسخ ارائه می‌کند.

اگرچه الگوریتم رتبه‌بندی یک موتور جستجو از کاربران مخفی نگه داشته می‌شود، اما کاربران با روش‌های مختلف سعی در کشف آن و ارائه مطالب سایت بصورتی که رتبه بالاتری کسب کنند، دارند.

3-3 اصول اولیه SEO

بهینه سازی برای موتورهای جستجو یا SEO به طور معمول به ایجاد تغییرات در بخش‌هایی از سایت گفته می‌شود که منجر به تاثیرات مطلوب در دستیابی پذیری بیشتر توسط موتورهای جستجو می‌شود.

در ادامه به مواردی که حضور کلمات کلیدی در آن تاثیر زیادی در میزان دسترسی‌پذیری سایت توسط موتورهای جستجوی می‌شود اشاره می‌کنیم.

1-3-3 تگ عنوان (`<title>`)

عنوان یک صفحه وب هم به کاربران و هم به یک موتور جستجو می‌گوید که موضوع اصلی و خاص صفحه مورد نظر چیست. این تگ در میان `<head>` قرار می‌گیرد. در حالت ایده‌آل شما باید برای

هر صفحه از سایتتان یک عنوان منحصر به فرد قرار دهید. از آنجا که عنوان انتخابی شما در خلاصه

توضیح¹ نتایج جستجو نشان داده می شود این مسئله نیز در انتخاب باید موثر باشد.

طول عنوان نیز مهم است بطوریکه کلمات کلیدی موجود در یک عنوان طولانی کمتر از کلمات کلیدی

در یک عنوان کوتاه ارزش دارند.

2-3-3 متا تگ های توضیح

این تگ ها که در تگ <head> آورده می شوند اطلاعات مفیدی برای موتورهای جستجو فراهم می کنند.

چرا که بر عکس عنوان که بهتر است چند کلمه باشد یک متا تگ را می توان در اندازه یک پاراگراف کوتاه

تنظیم کرد. همچنین به نظر می رسد موتورهای جستجو به مطالب این تگ ها اهمیت ویژه ای می دهند

بطوریکه در بیشتر موارد خلاصه توضیحات را بر پایه این تگ قرار می دهند، به همین علت بهتر است از

گذاشتن یک لیست صرف کلمات یا عباراتی که توضیح کاملی در مورد سایت ندارند خودداری شود. حتی

در حالت ایدهآل سعی شود متا تگ های توضیحی نیز بصورت منحصر به فرد نوشته شود. (البته در

مواردی که یک سایت صفحات زیادی دارد می توان با استفاده از اتوماتیک سازی برای این کار استفاده

کرد).

3-3-3 آدرس ها (URL)

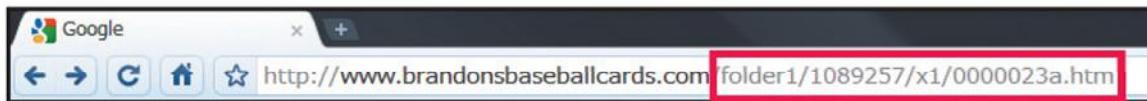
استفاده از نامهای با معنی برای فایل ها و فولدرهای یک سایت نه تنها به عنوان یک توصیه مناسب برای

طراحان سایت مطرح می شود بلکه منجر به کشف بهتر سایت توسط خزنه های موتورهای جستجو نیز

می شود. و در واقع موتورهای جستجو به سایت هایی که در آن به آدرس های با معنی ارجاع وجود دارد

امتیاز بیشتری می دهند. شکل 3-1 دو آدرس خوانا و ناخوانا را نمایش می دهد.

¹ snippet



1. یک آدرس بد برای سایتی که جهت ارائه کارت های بیس بال طراحی شده است.



2. یک آدرس مناسب. کلمات مشخص شده به کاربران و موتورهای جستجو کمک می کند تا فهم راحت تری از مفهوم ارائه شده در سایت داشته باشند.

شكل 3-1 نمونه ای از دو آدرس ناخوانا (شکل بالا) و ناخوانا (شکل پایین)

هرچه عمق آدرس کمتر و ساختار درختی و نحوه نامگذاری ها مناسب تر باشد هم سایت شما برای موتور جستجو دسترسی پذیرتر است و هم احتمال انتخاب شدن توسط کاربران به علت ذکر آدرس در خلاصه توضیحات بیشتر خواهد بود.

بهتر است سعی شود تا صفحاتی که از ریشه به آن ها می رسیم با صفحات شاخه ها متفاوت باشد به این معنی که داشتن آدرس domain.com\p1.htm هم زمان با sub.domain.com\p1.htm تاثیر بدی بر رتبه بندی خواهد داشت. همچنین انتظار می رود که آدرس ها با حروف کوچک نوشته شوند.

4-3-3 کلمات

به کلماتی فکر کنید که کاربران ممکن است برای یافتن مطلب شما جستجو کنند. مطمئناً کاربری که در یک موضوع خبره است کلماتی متفاوت از کسی که اطلاعات کمی دارد جستجو می کند به عنوان مثال کاربری که در مورد استخراج کلمات کلیدی تحقیق می کند ممکن است عبارت «بهبود بازخوانی» را جستجو کند اما کاربری که تازه به جمع محققین این زمینه پیوسته در جستجوی مطالبیست که عبارت

«ملک های کارایی» را در خود داشته باشد. موتور جستجوی گوگل توسط Google AdWords ابزار ساده Keyword Tool¹ را جهت کشف کلمات کلیدی متفاوت و بررسی حجم جستجو تقریبی این کلمات در اختیار کاربران قرار داده است. همچنین می‌توان با استفاده از ابزارهای تهیه شده در top search queries² نشان می‌دهد که سایت شما بیشترین بازدید را توسط چه کلماتی به خود اختصاص داده است.

در راهنمای SEO برای کاربران گوگل آمده است هرچه متن یک سایت بر موضوع بیشتر متمرکز باشد امتیاز بیشتری توسط موتور جستجو می‌گیرد که البته با توجه به رابطه منطقی فاکتور TFIDF در بررسی آماری متن قابل پیش‌بینی نیز می‌باشد. همچنین گفته شده از کپی محض یا تکرار مطالب نزدیک به مطالب دیگران نیز تا حد امکان پرهیز کنید.

5-3-3 متن پیوندی

متن پیوندی³ به آن متن‌های اتلاق می‌شود که قابل کلیک بوده و در صورت کلیک کاربر به آدرسی دیگر منتقل می‌شود. این متن که با استفاده از تگ <a> ایجاد می‌شود دارای دوبخش آدرس (herf) و متن نمایش می‌باشد. متن نمایش در واقع برای کاربران و موتورهای جستجو محتوای آدرس مورد اشاره را در چند کلمه شرح می‌دهد. برای نتیجه گیری بهتر است از اشتباهات رایج زیر اجتناب کنید.

- استفاده از کلمات عمومی مثل «اینجا را کلیک کنید» و غیره
- استفاده از متن و کلماتی که ارتباطی با صفحه مورد اشاره ندارند
- استفاده از متن بلند

¹ <https://adwords.google.com/select/KeywordToolExternal>

² <http://www.google.com/webmasters/edu/quickstartguide/sub1guide5.html>

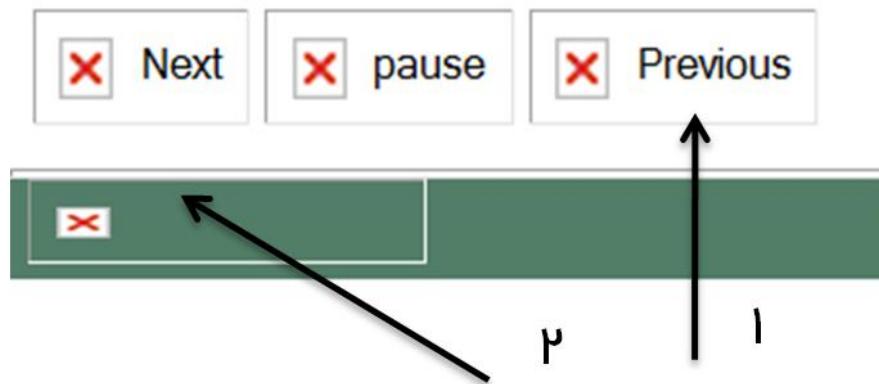
³ Anchor text

- استفاده از خود آدرس برای متن نمایش

- متن توضیح آدرس

6-3-3 تصاویر

تصاویر در HTML با استفاده از تگ img معرفی می شوند اما در میان صفت alt در این تگ نقش تعیین کننده‌ای برای موتورهای جستجو دارد. صفت alt به طراح سایت اجازه می‌دهد که یک متن جایگزین برای تصاویر تعریف کند تا در صورت عدم نمایش تصویر توسط مرورگر به هر علت کاربر بداند که در این مکان چه چیزی نمایش داده شده است. شکل 2-3 قسمتی از یک وب سایت را نشان می‌دهد که مرورگر نتوانسته است تصاویر آن را نمایش دهد. پیکان 1 به عکسی اشاره دارد که برای آن صفت alt برابر با previous قرار داده شده است و پیکان 2 به تصویری اشاره دارد که alt آن مقداردهی نشده است.



شکل 2-3 نمونه‌ای از یک عکس دارای alt (1) و عکس بدون alt (2)

مقدار صفت alt یا به تعبیری محتوای صفت alt در واقع شرح متنی یک تصویر است از این رو برای موتورهای جستجو که نمی‌تواند عکس‌ها را تفسیر کنند اهمیت زیادی دارد.

7-3-3 تگ های تیتر و زیر تیتر

تگ های تیتر¹ در صورتیکه به درستی استفاده شوند تاثیر زیادی در کسب رتبه مناسبی در نتایج موتورهای جستجو می‌شوند چرا که به طور طبیعی واژه‌های کلیدی بیشتری به علت ارائه درونمایه سایت در خود دارند. البته یادمان باشد که این تگ‌ها برای مدیریت مطالب سایت استفاده می‌شوند و برای برجسته کردن یک مطلب و ... از تگ‌های مخصوص اینکار مثل `` یا `` باید استفاده کرد. در کل توصیه می‌شود از اشتباهات زیر در رابطه با تگ‌های تیتر اجتناب کنید:

- قرار دادن جملات طولانی در یک تگ تیتر
- استفاده‌ای برای کاربردهایی غیر از ارائه ساختار سایت

¹ Headings

فصل چهارم:

پایگاه داده ها

1-4 مقدمه

همانگونه که در فصل قبل ذکر شد یکی از مشکلاتی که معمولاً تحقیقات مرتبط با پردازش زبان فارسی با آن روبرو می‌شود نبود منابع زبانی کافی و معتبر برای فارسی است. همانطور که در بخش کارهای گذشته ذکر شد در رابطه با این مسئله تلاش‌هایی صورت گرفته است اما تا زمان اجرای مراحل این تحقیق هیچ پایگاه داده‌ای که بصورت مشخص برای مسئله استخراج کلمات کلیدی باناظر قابل استفاده باشد (در ساده‌ترین حالت کلمات کلیدی آن از دید یک کارشناس استخراج شده باشد)، ارائه نشده است. به همین جهت برای پروژه استخراج کلمات کلیدی اقدام به ایجاد یک پایگاه داده کوچک نموده ایم که در این فصل به معرفی این پایگاه داده می‌پردازیم.

2-4 نقطه شروع

بسته به نوع تحقیق متون استفاده شده برای ایجاد پایگاه داده می‌توانند متون رسمی، خبری، اداری و حتی محاوره‌ای باشد. اما در صورتیکه بخواهیم پایگاه داده جامعیت پیدا کرده و برای تحقیقات مشابه نیز قابل استفاده باشد، به کارگیری متون رسمی دارای تنوع در موضوعات نسبت به سایر موارد ارجح است.

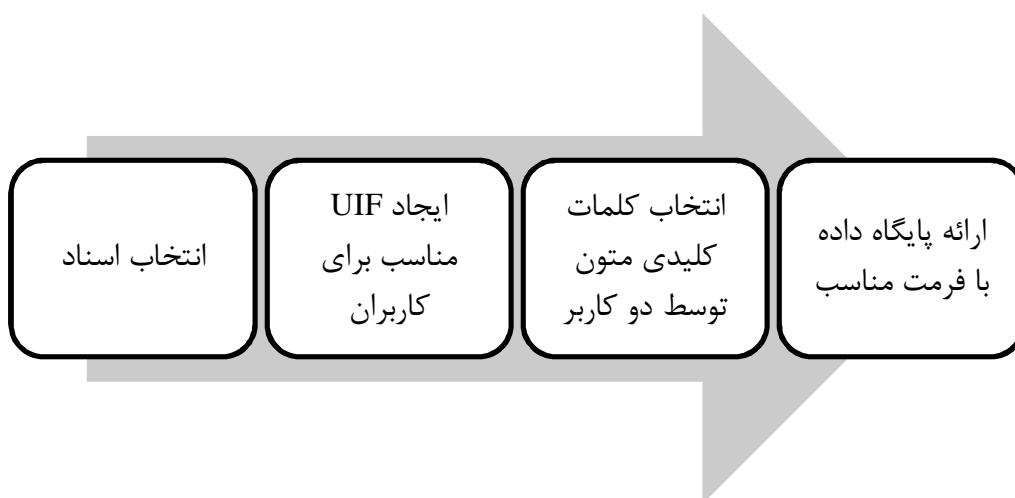
پایگاه داده پروژه‌های استخراج کلمات کلیدی می‌تواند تاثیر بسازایی در کارایی تحقیق داشته باشد این تاثیر با توجه به این مسئله که فاکتورهایی مثل IDF وابستگی مستقیم به پایگاه داده دارند بیشتر از سایر تحقیقات مرتبط با بازیابی اطلاعات خواهد بود.

به عنوان نمونه‌ای همانطور که در بحث کلمات پرتکرار مطرح شد یک کلمه می‌تواند در یک مجموعه متون مشخص فاقد ارزش تلقی شود در حالیکه در مجموعه دیگر کاملاً کلمه‌ای کلیدی محسوب می‌شود.

به عنوان مثال فرض کنید پایگاه داده یک تحقیق فقط شامل متون ورزشی باشد. آنگاه کلمه هایی مثل استادیوم، گل ، امتیاز، حریف و ... کلماتی بی ارزش تلقی می شوند. این در حالیست که اگر در یک مجموعه که فقط شامل متون معماری است یکی از کلماتی کلیدی متن «پژوهشی بر معماری استادیوم آزادی» کلمه «استادیوم» خواهد بود.

برای شروع کار از پایگاه داده همشهری استفاده کردیم. این پایگاه داده همانطور که در فصل دوم توضیح داده شد از جامعیت خوبی در تعداد سند و تعدد موضوعات برخوردار است. همچنین براحتی در دسترس بود و به علت آنکه مکرر در تحقیقات کنونی در حال استفاده شدن است می توان گفت اعتبار نسبتاً خوبی نیز دارد.

روند ایجاد پایگاه داده در شکل 1-4 نشان داده شده است.

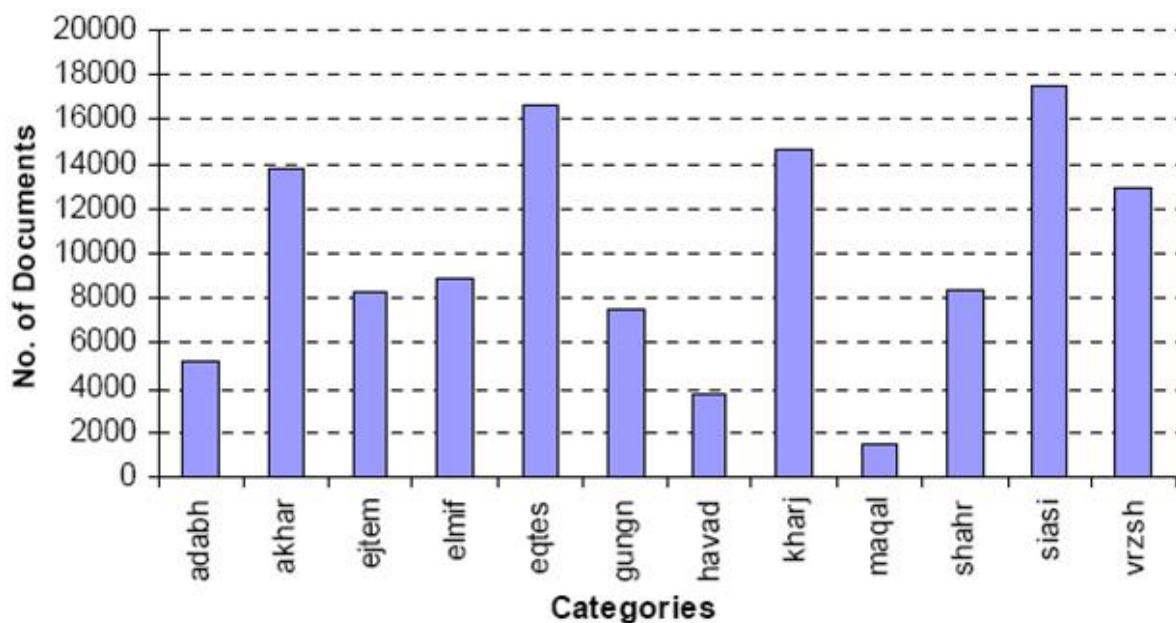


شکل 1-4 روند اتخاذ شده برای ایجاد پایگاه داده برای این پروژه

3-4 انتخاب اسناد

پایگاه داده همشهری از نظر توزیع اسناد دارای نموداری به صورت شکل 2-4 است. طول این اسناد متفاوت بوده و در برخی موارد به 8000 تا 9000 کلمه در یک سند نیز می رسد. اما از آنجا که برای انجام این تحقیق با کمبود منابع انسانی برای استخراج مواجه بوده ایم از اسنادی که کمتر از 1000 کلمه داشته‌اند استفاده کرده‌ایم. همچنین برای ایجاد جامعیت از 11 رسته مقالات پایگاه داده همشهری تعداد مناسبی از اسناد را انتخاب کرده ایم که در مجموعه به 1100 سند رسیده ایم. (از رسته مقالات به علت طولانی بودن هیچ سندی انتخاب نشده است.)

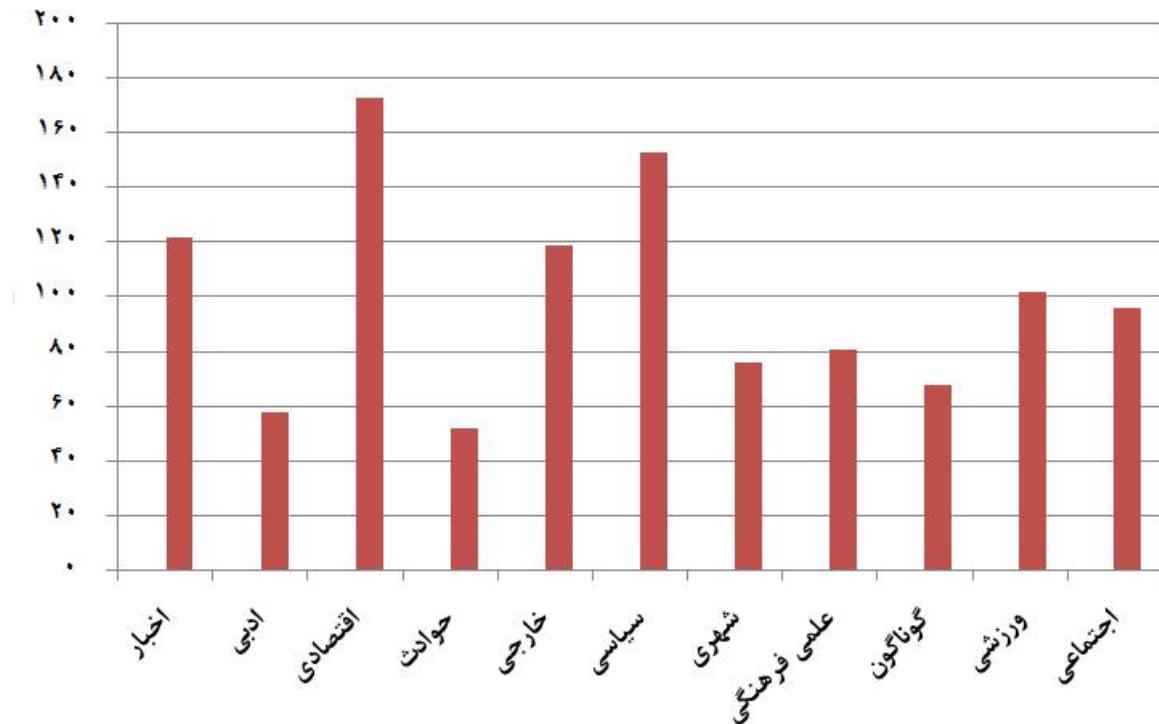
تعداد اسناد انتخاب شده از رسته های مختلف در شکل 4-3 نمایش داده شده است.



شکل 4-2 نمودار توزیع آماری اسناد در پایگاه داده همشهری از دید موضوع

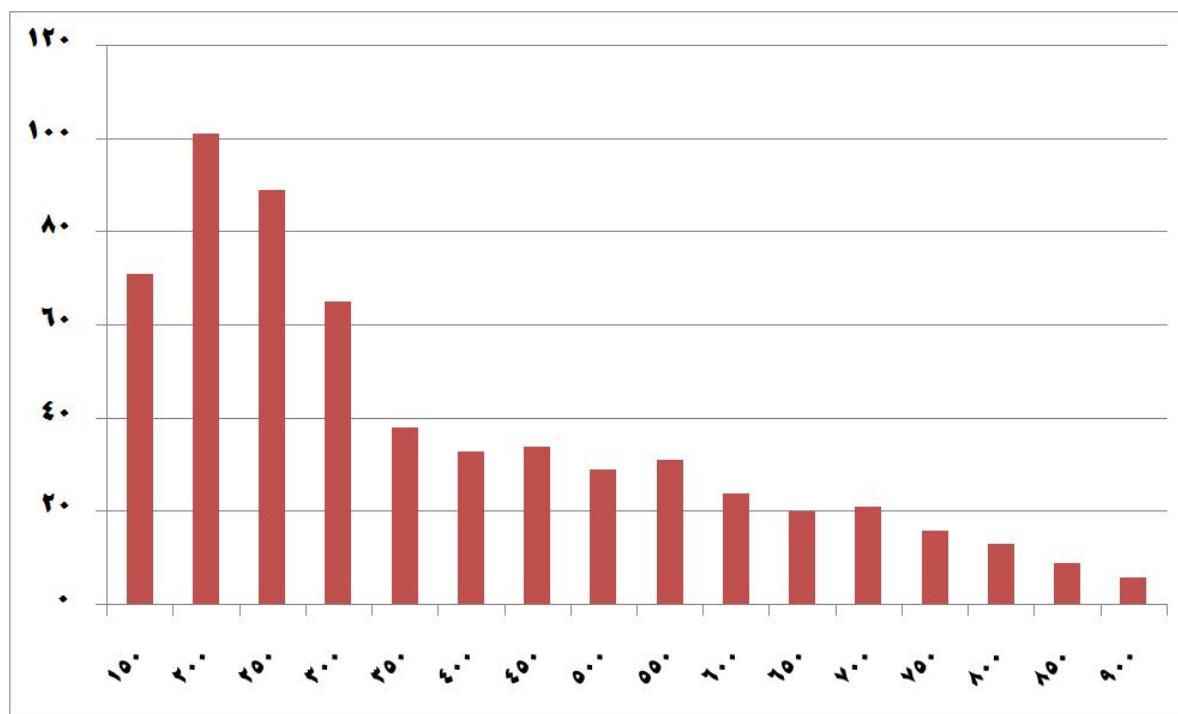
در واقع تعداد اسناد انتخابی بر اساس این موضوع که هریک از این رسته ها چند سند در پایگاه داده همشهری داشته‌اند انتخاب شده است.

اما برای آنکه جامعیت پژوهش محدود نشود پایگاه داده در دو حالت تعداد سند مساوی (50 سند از هر دسته) و تعداد سند متغیر در هر رسته در نتایج پایانی به کار برد شده و نتایج این مسئله نیز مورد بررسی قرار گرفته است.



شکل 4-3 توزیع موضوعی اسناد پایگاه داده ایجاد شده در این پژوهه

با زه تعداد کلمات موجود در اسناد انتخابی از 150 کلمه تا 900 کلمه است که بطور میانگین هر سند 4-4 372,5 کلمه را دارد. نوع توزیع اسناد انتخابی بر اساس تعداد کلمه موجود در اسناد در شکل آورده شده است.



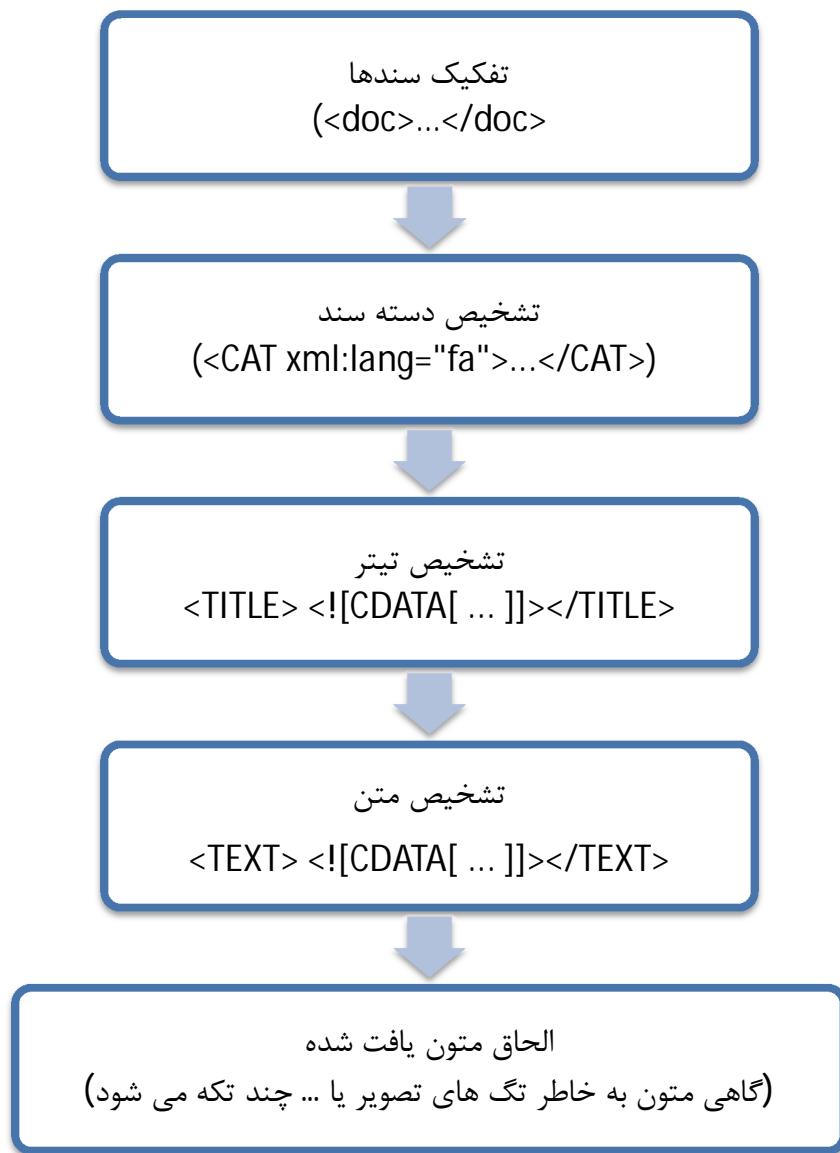
شکل 4-4-تعداد استناد پایگاه داده براساس طول بر حسب کلمه

پس از این مرحله پایگاه داده برای مرور و استخراج کلمات کلیدی با استفاده از یک واسط کاربری ساده به خوانندگان ارائه شده است.

4-4 پاک سازی تگ ها برای پایگاه داده

از آنجا که ساختار پایگاه داده همشهری بصورت XML طراحی شده است. حذف تگ های XML برای نمایش به کاربر و تطابق با پایگاه داده طراحی شده ضروری است. برای اینکار از روندی مطابق با شکل 5-4 استفاده کرده ایم.

پس از این مرحله پایگاه داده برای مرور و استخراج کلمات کلیدی با استفاده از یک واسط کاربری ساده به خوانندگان ارائه شده است.



شكل 4-5 روند پاک سازی یک سند پایگاه داده همشهری

4-5 انتخاب کلمات کلیدی

در این مرحله برای هر سند از دو کاربر خواسته ایم تا به انتخاب کلمات کلیدی بپردازد تا از کانالیزه شدن کلمات بر اساس تفکر یک کاربر جلوگیری شود.

واسط کاربری ایجاد شده به نحو بسیار ساده ای طراحی شده است. این واسط ابتدا متن هر سند را از پایگاه داده همشهری جدا کرده و پس از دریافت کلمات کلیدی کاربران آن را بصورت شکل 6-4 در یک رکورد پایگاه داده ذخیره می کند.

نام فیلد	ID	category	RTF doc address	Keyword1	Keyword2
توضیحات	یک شماره منحصر به فرد	دسته سند (سیاسی و...)	آدرس فایل ذخیره شده سند	کلید خارجی از جدول کلمات کاربران	کلید خارجی از جدول کلمات کاربران

شکل 6-4 شکل رکورد ذخیره سازی اطلاعات و کلمات کلیدی یک سند

همچنین یکی از نکاتی که منجر به بدست آوردن بازخوانی کمتر می شود معمولاً انتخاب کلماتی توسط کاربران است که در متن وجود ندارد. برای آن که تحلیلی در مورد تاثیرات این پدیده نیز داشته باشیم کلمات کلیدی را که رکورد های آن بصورت شکل 7-4 می باشد ذخیره کرده ایم.

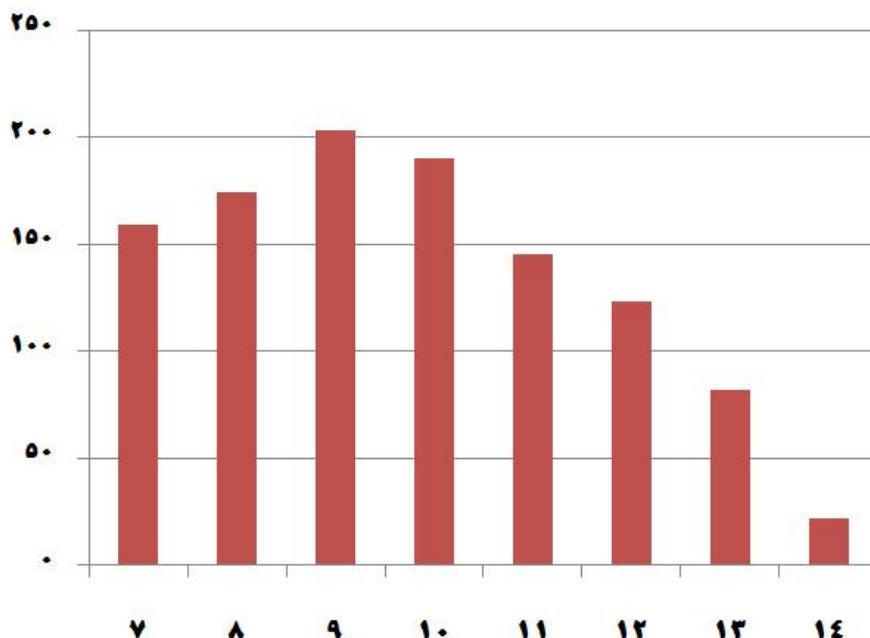
ID	Keyword	Does_exist
کلید جدول	کلمه کلیدی	یک فیلد باینری

شکل 7-4 ساختار رکورد جدول ذخیره کننده کلمات کلیدی

اگر کلمه کلیدی انتخابی کاربر دقیقا در متن وجود داشته باشد فیلد باینری Does_exist دارای مقداری برابر با 1 است و در غیر این صورت دارای مقدار صفر خواهد بود.

با توجه به این که از خوانندگان خواسته شده بود حداقل 10 و حداقل 4 کلمه کلیدی برای هر متن انتخاب کنند در مجموع بین 7 تا 14 کلمه به عنوان کلمه کلیدی به هر سند منصوب شده است. میانگین

تعداد کلمه انتخابی برای هر متن ۹,۷۲ می باشد. شکل ۴-۸ توزیع تعداد کلمات کلیدی را در متون نشان می دهد.



شکل ۴-۸- تعداد اسناد موجود در پایگاه داده بر اساس تعداد کلمات کلیدی انتخاب شده برای آن شکل ۴-۹ یک نمونه از اسناد پایگاه داده را نشان می دهد. این سند ۳۰۵ کلمه دارد و کلمات کلیدی منتخب کاربران برای آن به شرح زیر است:

کلمات کلیدی کاربر1: کوفی عنان - دبیرکل سازمان ملل - حقوق بشر - ایدز - تغییرات جوی -
جهانی شدن اقتصاد - خطرات بی مرز

کلمات کلیدی کاربر2: کوفی عنان - فلسفه وجودی سازمان ملل - همکاری های مشترک بین المللی -
تغییرات اقتصادی - آمریکا - بدون همکاری بین المللی - جهانی شدن اقتصاد و تجارت - حقوق بشر -
جامعه بشری

کوفی عنان: جهانیان! برای مقابله با
خطرات بی مرز متحد شوید!

سازمان ملل - خبرگزاری جمهوری اسلامی: ببیر کل سازمان ملل در پیامی به مناسبت آغاز سال جدید میلادی، خواستارگسترش همکاری های بین المللی برای مقابله با مشکلات جهانی شد.
کوفی عنان گفت: مردم جهان سرنوشت مشترکی دارند اما فقط زمانی می توانند بر این سرنوشت مشترک مسلط شوند که با یکدیگر، با آن روبرو شوند.

عنان با تکیه بر اینکه فلسفه وجودی سازمان ملل، همکاری های مشترک بین المللی است، گفت: قرنی که آغاز می شود امیدهای جدیدی را بر می انگیزد اما می تواند خطرهای جدیدی نیز به ارمغان آورد و یا خطرهای کهنه را در اشکالی جدید و هشداردهنده عرضه کند.
ببیر کل سازمان ملل از تغییرات اقتصادی، تعصب، خشونت، بیماری ها و صدمه به محیط زیست، به عنوان خطرهای عمدۀ ای که جهان را تهدید می کند نام برد.

عنان تأکید کرد: هیچکس از میزان شدت این خطرات دقیقاً مطلع نیست اما واضح است که این خطرات مرزهای کشورها را در می نوردد.

ببیر کل سازمان ملل با اشاره تلویحی به آمریکا گفت: حتی قدرتمندترین کشورها نیز نمی توانند به تنهاپی و بدون همکاری بین المللی، اتباع خود را در مقابل این خطرات حفظ کند.

عنان افزود: از طریق سازمان ملل و با تلاش مشترک می توان با ایدز و دیگر بیماری های واگیردار مقابله کرد، تغییرات جوی را کنترل نمود و آب و هوای پاکیزه برای همه فراهم ساخت.

وی خاطرنشان کرد: به علاوه سازمان ملل تلاش می کند جهانی شدن اقتصاد و تجارت به نفع همه باشد و به فقر اجازه دهد خود را از فقر خارج سازند.

عنان تأکید کرد: باید از طریق سازمان ملل، حقوق بشر را برای همه به یک واقعیت تبدیل کرد و به کل جامعه بشری این فرصت را داد در تصمیماتی که بر زنگی آنها تاثیر می گذارد، حق اظهار نظر واقعی داشته باشند.

شکل 9-4 نمونه ای از یک سند پایگاه داده ایجاد شده

فصل پنجم:

فرایند های پیش پردازش

1-5 مقدمه

در این فصل به معرفی فرایند پیش پردازش ایجاد شده برای این پروژه می‌پردازیم. در ابتدا با استفاده از تجربیات سایر تحقیقات به اعمال یکسری از قوانین ساده یکسان‌سازی می‌پردازیم. سپس در بخش بعد با نگاهی تکنیکی‌تر نسبت به آنچه در فصل دوم درباره کلمات پر تکرار ارائه شد به حذف آن‌ها می‌پردازیم. راه حل ارائه شده برای حذف کلمات پر تکرار در این پروژه استفاده از یک لیست ثابت است اما تغییرات کوچکی برای بهبود فرایند حذف کلمات بی ارزش ارائه شده است. بخش انتهایی این فصل به مسئله توکنیزه کردن کلمات، و جملات می‌پردازد. توکنیزه کردن جملات از آن جهت انجام می‌شود که بخشی از خصوصیات استفاده در فصل آینده به جملات و وضعیت آن‌ها بستگی دارد. این پروژه از ریشه یابی خاصی استفاده نکرده است چرا که الگوریتم‌های ریشه یابی موجود هنوز در ابتدای راه هستند و حتی برای برخی قواعد دستوری زبان فارسی در بعضی موارد هنوز تعریف رسمی یکسانی وجود ندارد.

2-5 یکسان‌سازی

در بعد یکسان‌سازی تکنیکی با توجه به غیر جامع بودن این تحقیق از جنبه اسناد ورودی چه از منظر نوع ساختارهای ورودی (در بخش پایگاه داده ذکر شد که اسناد به فرمت RTF برای استفاده ذخیره شده اند) و چه از منظر نوع رمزگذاری (با توجه به انتخاب اسناد از پایگاه داده همشهری که دارای رمزگذاری یکسان در سرتاسر اسناد است). یکسان‌سازی خاصی نیاز نداریم. اما در بخش قواعد زبانی با استفاده از مجموعه‌ای کوچک از قواعد که در مقاله [17] ذکر شده است به یکسان‌سازی نسبتاً خوبی رسیده ایم.

- تبدیل نویسه‌های «و» به «و»، «ی» به «ی» و «ا» به «ا».

- تبدیل نویسه‌های «ه» و «هء» به «ه» در آخر واژه‌ها.
- حذف «ی» از آخر واژه‌هایی مانند «خانه‌ی».
- حذف فتحه، کسره و ضمه (نویسه‌های ـ، ــ، ـــ) از واژه‌ها.
- حذف تنوین (نویسه‌های ـــ، ــــ، ـــــ) از انتهای واژه‌ها.
- حذف تشدید (شناسه‌ی ــ) از واژه‌ها
- حذف شناسه‌ی «ء» در آخر بعضی واژه‌ها مانند «شهداء»
- چسباندن پیشوندهای «می»، «درمی»، «برمی»، «نمی» و «بی» به ابتدای واژه‌ها
- چسباندن پیشوند «هم» به واژه‌های مانند «هم چنین» به ابتدای واژه‌ها
- چسباندن پسوندهای «ها»، «های»، «هایم»، «هایت»، «هایش»، «هایمان»، «هایتان» و «هایشان» به انتهای واژه‌ها
- چسباندن پسوندهای «تر» و «ترین» به آخر واژه‌ها
- حذف فاصله بعد از پیشوند «بر» در واژه‌هایی مانند «بر می‌گردد».
- حذف نویسه‌ی «ـ» که برای کشش نویسه‌های چسبان جهت تراز شدن طول خطها مورد استفاده قرار می‌گیرد. مانند تبدیل «بر» و «بر» به «بر».

شیوه یکسان سازی به این صورت بوده است که عبارات اشاره شده در قوانین در تمامی متون مورد بررسی جستجو شده و در صورت انطباق قانون مورد نظر اعمال شده است.

به عنوان مثال برای اعمال یکسان سازی منطبق بر آخرین قانون ابتداء نویسه «ـ» مورد جستجو قرار گرفته سپس در صورت یافتن شدن ادامه رشته تا رسیدن به کاراکتری غیر از «ـ» برای پیدا کردن تعداد

بیشتر این نویسه مورد بررسی قرار گرفته است. سپس نویسه «_» حذف شده و دو زیر رشته^۱ به جا مانده از حذف به یکدیگر الحق می شوند.

۵-۳ حذف کلمات پر تکرار

در این پژوهه برای حذف کلمات پر تکرار از یک لیست ثابت استفاده کرده ایم. اگرچه استفاده از روش های پویا دقت بالاتری داشته و مشکلات لیست ثابت را که در بخش قبل به آن اشاره شد در پی ندارد اما پیچیدگی های زمانی روش پویا بالاتر است. در واقع در توازن بین انتخاب بهتر کلمات پر تکرار و هزینه زمانی کمتر، برای روش ما با توجه به ارزیابی هایی که بعد از این نیز وجود دارد، وجود یک کلمه کم ارزش در لیست کلمات کلیدی به مرتب قابل تحمل تر از هزینه ای زمانی ساخت لیست بصورت پویاست.

لیست ثابت به این نحو تولید شده است که لیست 927 کلمه ای مرجع[20] را با لیست دیگری که طی پژوههای دیگر با استفاده از مرور 45000 وبلاگ فارسی بدست آمده بود ادغام کرده ایم.

اما آنچه که به عنوان ابتکار در این بخش صورت گرفته است توجه بیشتر به توالی کلمات پر تکرار است. تجربه نشان می دهد که یک توالی طولانی از کلمات پر تکرار خود یک عبارت با ارزش است. مثلا «بودن یا نبودن مسئله اینست» یا «در این درگاه که گه که که شود ناگه». در این پژوهها توالی های بیش از 4 کلمه را به عنوان یک عبارت کاندید ستاره دار در نظر می گیریم. و با عبارات کاندید ستاره دار مانند یک کلمه کلیدی استخراج شده توسط تابع امتیازدهی برخورد می کنیم.

روش حذف کلمات پر تکرار به این نحو است که در صورت انطباق کلمه ای با یکی از کلمات لیست ، یک توالی به طول یک تشکیل شده و ادامه رشته کلمات برای انطباق بیشتر پردازش می شود و هر انطباق به

^۱ زیر رشته اول از نویسه فاصله تا قبل از اولین نویسه «_» و زیر رشته دوم بعد از آخرین نویسه «_» تا رسیدن به اولین فاصله در نظر گرفته شده است

این توالی اضافه شده و طول افزایش می یابد با اولین عدم انطباق طول توالی بررسی خواهد شد. در صورتیکه طول توالی بیش از ۳ کلمه باشد توالی به عنوان یک عبارت کاندید ستاره دار ثبت می شود و در غیر اینصورت به عنوان یک یا چند عبارت پر تکرار حذف خواهد شد.

4-5 تعیین مرز جملات

در این قسمت روال ساده‌ای را، که با توجه به خصوصیات زبان فارسی برای حل مشکل نقطه جهت تعیین مرز جملات طراحی کردہ‌ایم، شرح می‌دهیم.

1. برای حل مشکل نقطه «.» در علائم اختصار یا اختصارات اسمی، هر نقطه‌ای که به یک تک حرف چسبیده باشد را نقطه پایان جمله حساب نمی‌کنیم. و در صورتیکه نقطه از هر دو طرف به یک تک حرف محدود شده باشد آن را یک علامت اختصار حساب می‌کنیم.
2. برای حل مشکل «.» در آدرس‌های وب، هرگاه نقطه‌ای از دو طرف به دنباله‌هایی چسبیده باشد که شامل اختصارات کلیدی وب از قبیل «.com»، «.www» و ... باشد، در این صورت آن دنباله را یک آدرس وب تلقی می‌کنیم.
3. در صورتیکه هریک از رخدادهای بالا اتفاق بیفتد دنباله بوجود آمده تفکیک نشده و به عنوان یک واژه رها می‌شود در غیر اینصورت رشته نویسه‌های سمت چپ نقطه تا رسیدن به آخرین نقطه دیده شده به عنوان یک جمله حساب می‌شود.

5-5 تعیین مرز کلمات

برای توکنیزه کردن کلمات به علت حل تعداد زیادی از مشکلات در فرایند یکسانسازی از هیچ الگوریتم خاصی استفاده نکرده ایم و به تفکیک با استفاده از نویسه های رایج بسته کرده ایم. تنها در چند مورد فرایندهای ساده ای را اجرا کرده ایم.

1. در زبان فارسی کلمه تک حرفی با معنی غیر از «و» وجود ندارد به همین علت در صورت برخورد

با کلمات تک حرفی در صورتیکه کلمه پیش از آن ها به حرف از راست پیوند ناپذیر¹ ختم شده باشد کلمه تک حرفی به عنوان حرف آخر کلمه قبل در نظر گرفته شده است و در غیر اینصورت به عنوان یک اشتباه حذف شده است.

2. برای نشانه «/» کلمات جانشین به عنوان دو کلمه هم اتفاق² در نظر گرفته شده اند و تفکیک

نشده اند این مسئله باعث می شود که در آدرس ها و تاریخ ها نیز تفکیکی صورت نگیرد. در واقع نویسه «/» به عنوان یک نویسه معمولی در نظر گرفته شده است.

3. برای نشانه «،» در صورتیکه در هر دو طرف نشانه عدد باشد از حذف نشانه چشم پوشی کرده و

کلمات و نشانه بینشان را به عنوان یک کلمه یکسان اعلام می کنیم.

باید توجه داشت که با توجه به استفاده از قالب های 1 تا 4 کلمه ای و نیز نوع مسئله استخراج کلمات کلیدی که بدون توجه به ساختارهای نحوی نیز با کیفیت خوبی قابل انجام است، از هر نوع تشخیص هم آئی کلمات ، افعال مرکب و ... چشم پوشی کرده ایم. این چشم پوشی تاثیر چندان زیادی بر نتیجه کلی نخواهد داشت.

¹ الف - د - ذ - ر - ز - ڙ - و

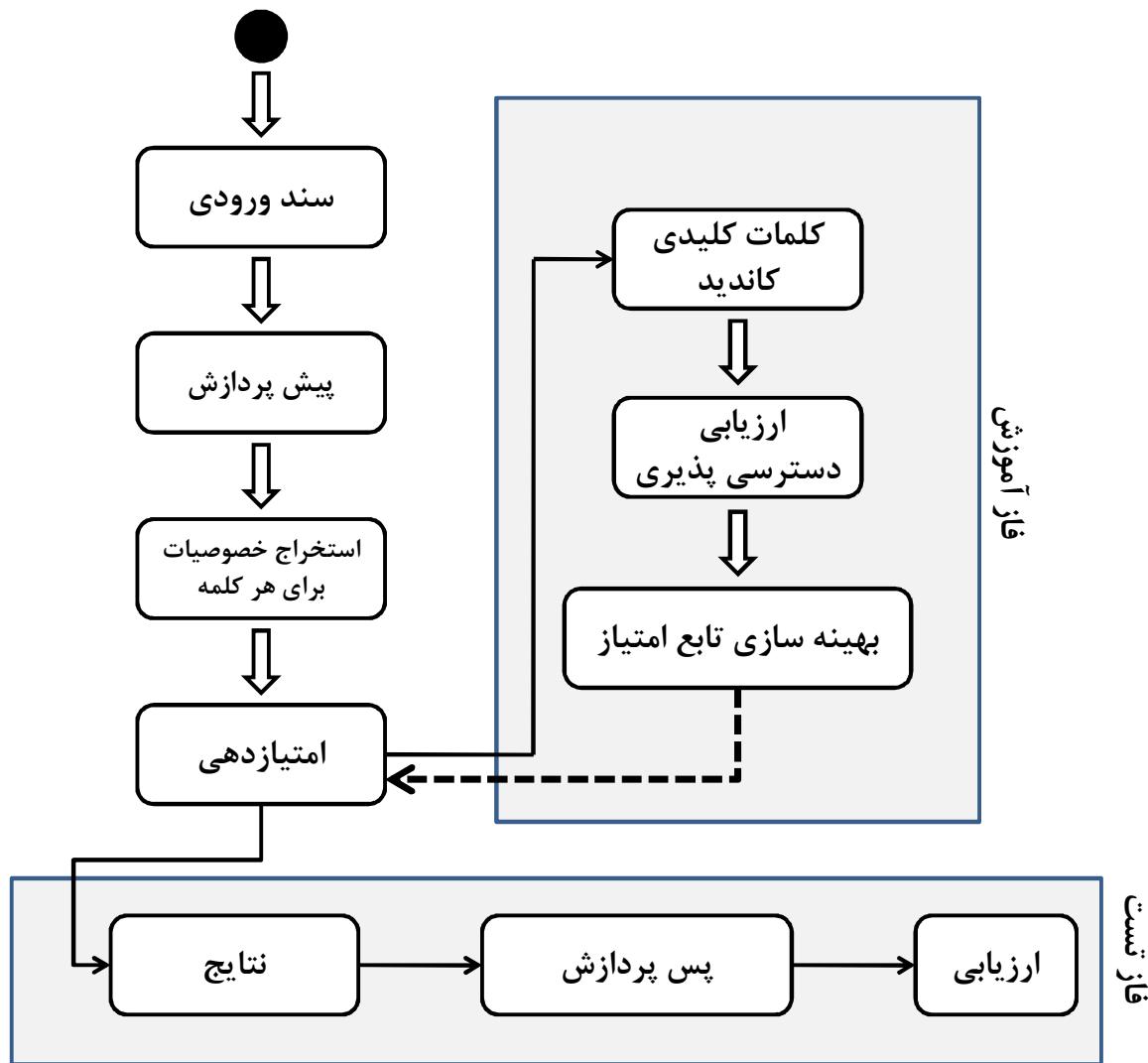
² Co-occurrence

فصل ششم:

استخراج کلمات کلیدی

1-6 مقدمه

شکل 1-6 شمای کلی روش ارائه شده در این پژوهه برای استخراج کلمات کلیدی را نشان می دهد فاز استخراج کلمات کلیدی شامل دو بخش می باشد: بخش آموزش و بخش تست



شکل 1-6 شمای کلی روش ارائه شده در این پژوهه برای استخراج کلمات کلیدی

روش ارائه شده در این پژوهه از دو مجموعه خصوصیات و دو تابع ارزیابی استفاده می‌کند. تابع اول که به نام تابع امتیاز دهی نامگذاری شده با استفاده از مجموعه خصوصیات¹ به هر یک از توکن‌های استخراج شده در بخش قبل امتیاز می‌دهد. بعضی از مجموعه خصوصیات¹ خصوصیات جمعی هستند و به همین علت این مجموعه خصوصیات وابستگی مستقیم به پایگاه داده دارند.

تابع دوم که به نام تابع ارزیابی نتایج موتورهای جستجو¹ (SERE) معرفی شده است تابعی است که بر اساس مجموعه خصوصیات² که از یک سند منفرد استخراج می‌شوند براساس نتایج موتورهای جستجو به کلمات امتیاز می‌دهد. در واقع بازخورد اهمیت کلمه انتخاب شده توسط تابع امتیاز دهی در تابع SERE نمایان می‌شود. تابع امتیاز دهی به این نحو طراحی شده طراحی شده است که ضرایب آن توسط الگوریتم ژنتیک بهینه سازی می‌شود. تابع ارزیابی الگوریتم ژنتیک مورد استفاده برای بهینه سازی برپایه SERE بنا نهاده شده و جمعیت را ضرایب استفاده شده در تابع امتیاز دهی معرفی کرده‌ایم. فرایند استخراج کلمات کلیدی طوری پیش می‌رود که ضرایب تابع امتیاز دهی در طی انجام پروسه‌ی آموزش به نحوی تنظیم شوند، که ماکزیمم امتیاز در تابع SERE توسط کلماتی که ماکزیمم امتیاز را در تابع امتیاز دهی بدست می‌آورند، کسب شود. در ادامه این فصل ابتدا نگاهی داریم به بخش‌های آموزش و تست سپس به مرور کامل این دو تابع، مجموعه‌های خصوصیات و سایر توابع و فرایندهای مورد استفاده در پژوهه می‌پردازیم.

2-6 بخش آموزش

در این بخش هر دو تابع امتیازدهی و SERE مورد استفاده قرار می‌گیرند.

¹ Search engines results evaluation (SERE)

در این پژوهه از پنجره‌های ۱-۴ کلمه‌ای استفاده کرده‌ایم به دلیل آنکه تحقیقات جامعی در مورد پیمانه‌های^۱ کلمات در فارسی وجود ندارد به همین علت از یک پنجره ساده استفاده کردایم. در مرحله اول این بخش برای هریک از پنجره‌های n تایی خصوصیات مجموعه ۱ برای استفاده در تابع امتیاز دهی استخراج می‌شود. خصوصیات مجموعه اول، گزیده‌ای از خصوصیات تعریف شده در تحقیقات دیگر هستند که تاثیرگذاری بیشتری در مسئله تشخیص کلمات کلیدی داشته‌اند. پس از مرحله امتیازدهی توکن‌ها بصورت صعودی از نظر امتیاز مرتب می‌شوند و با استفاده از رابطه ۱-۶ کلمات کандید استخراج می‌شوند.

$$\text{Score}(\text{word}_{i,d}) > \text{threshold} \quad 1-6$$

مقدار threshold یک مرز امتیازی است تا کلماتی که بیش از این مرز امتیاز می‌گیرند به عنوان کلمات کلیدی تلقی شوند. این مرز به نحوی تعیین می‌شود که حداقل ۶۰٪ از کلمات انتخاب شده کاربران در این لیست دیده شود.

پس از آنکه لیست کلمات کандید ساخته شد این لیست به مازول ارزیابی و بهینه سازی تحویل می‌شود. مازول ارزیابی و بهینه سازی ابتدا لیست را برای موتورهای جستجوی محبوب مانند گوگل، یاهو، MSN و... فرستاده و صفحه اول نتایج را به عنوان بازخورد دریافت می‌کند. صفحه اول نتایج با استفاده از SERE مورد تحلیل قرار می‌گیرد.

مازول ارزیابی و بهینه سازی با استفاده از الگوریتم ژنتیک به بهینه سازی ضرایب موجود در تابع امتیازدهی می‌پردازد و این فرایند آموزش تا رسیدن به بیشینه رسیدن تابع ارزیابی ادامه می‌یابد.

^۱ N-gram

1-2-6 تابع امتیاز دهی

تابع امتیاز دهی مورد استفاده در رابطه 6-2 نشان داده شده است. این تابع هم در بخش آموزش و هم در بخش تست مورد استفاده قرار می‌گیرد. اما در بخش آموزش با استفاده از الگوریتم ژنتیک بهینه سازی می‌شود.

$$\text{Scoring function (w)} = \text{NPW} * (\alpha * \text{NPLen} + \beta * \text{NSL} + \gamma * \text{NTFIDF}) / 3 \quad 2-6$$

در ابتدا تمامی ضرایب α , β و γ برابر با یک عدد تصادفی در بازه $[0.4, 0.6]$ قرار داده می‌شوند اما در طول فاز آموزش هریک از آن‌ها به سمت ماکزیمم کردن تابع ارزیابی پیش می‌روند که این منجر به بهتر شدن امتیاز کلماتی خواهد شد، که از نظر موتورهای جستجو مهمتر تلقی می‌شوند. همچنین به خاطر نوع تعریف مقدار threshold می‌توانیم تاحدود زیادی از بدست آوردن یک امتیاز بازخوانی خود مطمئن باشیم.

خصوصیات مورد استفاده به عنوان پارامترهای تابع امتیاز دهی که خصوصیات مجموعه اول را تشکیل می‌دهند به شرح زیر هستند:

1-1-2-6 NPW نرمالیزه شده تعداد کلمات

NPW^1 تعداد کلمات یک عبارت است که با تقسیم بر ماقزیمم تعداد کلمات پیمانه نرمالیزه شده است. از آنجا که حداکثر پیمانه مورد استفاده در این تحقیق 4 بوده پس مقادیری که این متغیر می‌تواند بدست آورد برابر با 1 برای عبارت مورد بررسی چهار کلمه‌ای، 0,75 برای عبارت سه کلمه‌ای، 0,5 برای یک عبارت دو کلمه‌ای و 0,25 برای عبارت یک کلمه‌ای است. انگیزه استفاده از این خصوصیت بصورت یک

¹ Normalized phrases words

ضریب برای سایر امتیازها، آنست که امتیاز بیشتری به عبارات چهار کلمه‌ای که با ارزش تراز عبارات سه کلمه‌ای هستند داده شود(همین طور این مسئله در مقایسه با یک کلمه‌ای و ...).

2-1-2-6 NPLen طول عبارت بصورت نرمالیز شده

NPLen¹ بصورت تعداد کلمات یک عبارت کاندید تقسیم بر تعداد کلمات جمله‌ای که شامل آن می‌باشد بدست می‌آید. این خصوصیت در صورتیکه کل جمله یک عبارت کاندید باشد برابر ۱ می‌شود. انگیزه استفاده از این خصوصیت آن است که، عبارت‌های موجود در عنوان و زیر عنوان‌ها به علت کوتاه بودن جملاتشان، امتیاز بیشتری بگیرند.

3-1-2-6 NSL امتیاز جایگاه جمله بصورت نرمالیزه شده

در بیشتر ارزیابی‌ها و روش‌های استخراج کلمات کلیدی ادعا شده است که کلماتی که در ابتدا یا انتهای یک متن وجود دارد به طور نسبی اطلاعاتی بیشتری در مورد اصل موضوع آن متن در خود دارند، به همین علت خصوصیت NSL² با استفاده از رابطه 3-6 طوری تعریف شده است که جملات ابتدایی و انتهایی در متن امتیاز بیشتری بگیرند.

$$NSL = \left(2 * \left(\frac{L}{m} \right) - 1 \right)^2 \quad 3-6$$

در رابطه 3-6 «L» نمایانگر رتبه جمله در متن و m به معنای تعداد جملات موجود در کل متن است. مقدار NSL برای اولین جمله (L=0) و آخرین جمله (L=m) برابر با ۱ و هرچه به مرکز متن نزدیک‌تر شویم کمتر از ۱ می‌شود.

¹ Normalized Phrase Length

² Normalized Sentence Location

NTFIDF 4-1-2-6

TFIDF یکی از مهمترین روش‌های وزن دهی به کلمات است. این ویژگی در بسیاری از روش‌های وزن-دهی مهم به عنوان یک پایه محسوب می‌شود [46]. برای مقید کردن این خصوصیت به بازه [0, 1] آن را با استفاده از رابطه 4-6 نرمالیزه کرده ایم.

$$NTFIDF(word) = \left(\frac{WF}{MF} \right) * \log_2 \left(\frac{N+1}{n+1} \right) \quad 4-6$$

WF تعداد تکرار کلمه در این سند و MF برابر با تعداد تکرار مکررترین کلمه سند است. N و n همانند آنچه در فصل قبل تعریف شد برابر با تعداد کل اسندها و تعداد اسناد شامل این کلمه می‌باشند.

6-2-2 تابع ارزیابی

این تابع در فاز آموزش به عنوان تابع ارزیابی الگوریتم ژنتیک مورد استفاده قرار می‌گیرد. از دید تکنیکی این تابع در واقع یک تابع وزن دهی به کلمات در یک سند منفرد به شمار می‌رود. از آنجا که بازخوردهای موتورهای جستجو صفحات وب هستند، استفاده از روش‌های آنالیز اهمیت عبارات با استفاده از برچسب‌های HTML برای کمک به وزن دهی بهتر در این تابع نیز امکان پذیر است.

تابع SERE در رابطه 5-6 نشان داده شده است.

$$SERE = PRF^*(TET + NLCT + NHCT + NIRT) \quad 5-6$$

پارامترهای تعریف شده در پرانتز در رابطه 5-6 بر اساس برچسب‌های HTML بدست آمده اند. با هم به مرور پارامترهای تابع SERE می‌پردازیم.

1-2-2-6 تعداد تکرار نسبی PRF

¹ PRF نشان دهنده تعداد تکرار کلمه در سند است که با تقسیم بر تعداد تکرار مکررترین کلمه نرمالیزه شده است. این پارامتر در واقع همان TF است که بالاترین مقدار یعنی 1 را در مکررترین کلمه گرفته و برای سایر کلمات ارزشی کمتر از 1 خواهد داشت.

2-2-2-6 وجود کلمه در عنوان TET

² TET برابر است با تعداد کلمات عبارت تقسیم بر تعداد کلمات موجود در عنوان سایت که با تگ <title> دیده می‌شود.

3-2-2-6 تعداد لینک‌های شامل عبارت NLCT

³ NLCT تعداد وقوع عبارت در برچسب <a> را نشان می‌دهد. برای نرمالیزه کردن این عدد را بر تعداد کل لینک‌های موجود تقسیم می‌کنیم.

4-2-2-6 تعداد تیترهای شامل عبارت NHCT

NCHT نشان دهنده تعداد وقوع عبارت در برچسب <h1> در نظر گرفته شده است. برای نرمال کردن این عدد آن را بر تعداد کل برچسب‌های <h1> تقسیم می‌کنیم.

5-2-2-6 تعداد تصاویر مرتبط با عبارت NIRT

⁴ NIRT تعداد تصاویر مرتبط با عبارت مورد ارزیابی را نشان می‌دهد. در تگ Img که نشان دهنده تصاویر است در HTML ⁵ است دو صفت مهم وجود دارد:

¹ The Phrase Relative Frequency

² Term Exist in Title

³ Number of Links Containing the Term

⁴ Number of Images Related with the Term

⁵ Attribute

Src: که آدرس تصویر را در خود دارد. (از این صفت به علت آنکه آدرس و نام فیزیکی تصاویر در بیشتر

موارد به زبان انگلیسی است استفاده نکرده ایم.)

.Alt: این صفت همانطور که در فصل سوم توضیح داده شد به عنوان توضیحی بر تصویر به کار می رود.

مقدار NIRT برابر با تعداد وقوع عبارت در تگ Img تقسیم بر تعداد کل تگ های Img موجود است.

در این تحقیق به همین چهار برچسب HTML که مهمترند بسنده کرده ایم اما می توان برای دقت بیشتر

با برچسب های دیگر نیز به این مسئله پرداخت. جدول 6-1 تعدادی از برچسب های HTML به همراه

توضیح مختصری برای آن ها و همچنین مثالی از کاربرد را نشان می دهد.

جدول 6-1 نمونه ای از برچسب های HTML به همراه توضیح و مثال

برچسب	توضیح	مثال
<html>	برای نشان دادن شروع و پایان کل یک صفحه وب	<html>
	به کار می رود	All other content goes here...</html>
<head>	این برچسب برای نشان دادن عنوان صفحه وب به کار می رود و تقریباً مهمترین برچسب بشمار می	<head>
<title>Welcome to my site</title>	رود.	<title>
<body>	برای مشخص کردن قسمت اصلی متن یک صفحه استفاده می شود	<body>
		Page body goes here</body>
<a>	برای نشان دادن لینک به صفحات دیگر به کار می	href="http://www.hanuz.blogfa.com">Vis
	رود	it the my site

	<p>توجه کنید که این برچسب برچسب مهمی است که وجود کلمه ای چه در قسمت آدرس آن (<i>herf</i>) و چه در قسمت توضیح آن می تواند اهمیت زیادی داشته باشد</p>	
	برای درونی سازی یک تصویر به کار می رود	
This is bold text while <i>this text is in italic</i>	به ترتیب از بالا برای تیره سازی ¹ کج نویسی ² و نوشته های زیر خطدار به کار می روند.	 <i> <u>
<table> <tr> <td>This is a Cell in Column 1</td> <td>This is a Cell in Column 2</td> <tr/> <table/>	<p>این برچسب ها در ساخت جدول ها به کار می روند. سطرهای جدول و td ستون های آن را ایجاد می کنند.</p> <p>از آنجا که معمولاً اطلاعات ارائه شده در جداول اطلاعاتی خلاصه و مفید هستند. این برچسب ها در سیستمهای امتیاز دهی مبتنی بر صفحات وب می توانند با ارزش تلقی شوند</p>	<table> <tr> <td>
<h1>Hello World</h1>	<p>این برچسب که مشخص کننده تیتر است می تواند تا زیر تیترهای درجه 6 ادامه پیدا کند بسته به کاربرد می توان برای هر درجه از تیتر ارزش خاصی قائل شد.</p>	<h1>
 First item with a bullet Second item with a bullet 	<p>تگ هایی که برای ساخت لیست های شماره دار و بی شماره به کار می روند. از آنجا که لیست ها معمولاً شامل اطلاعات مرتبط به هم هستند این</p>	

¹ Bold

² Italic

<pre> First item with a number Second item with a number </pre>	برچسب ها را می توان به نوبه خود با اهمیت تلقی کرد.	
--	--	--

6-2-3 الگوریتم ژنتیک اجرا شده

الگوریتم ژنتیک یک روش جستجوی تصادفیست که بر پایه انتخاب طبیعت و ارزیابی بنا نهاده شده است. یک GA برای حل مسئله، مجموعه بزرگی از راه حل ها را تولید می کند و هر راه حل را با استفاده از یک "تابع تناسب¹" ارزیابی می کند. آنگاه تعدادی از بهترین راه حل ها با استفاده از اپراتور های ژنتیکی باعث تولید راه حل های جدیدتری می شوند این کار تا رسیدن به جواب مطلوب ادامه پیدا می کند. به این ترتیب فضای جستجو در جهتی تکامل پیدا می کند که به راه حل مطلوب برسد. به طور عمومی اپراتورهای ژنتیکی به در دو دسته آمیزش و جهش تقسیم می شوند و هریک دارای انواع مختلفی می باشند. می توان ادعا کرد در صورت انتخاب صحیح پارامترها الگوریتم ژنتیک می تواند بسیار موثر عمل کند. در واقع ژنتیک الگوریتم به جای جستجو راه حل ها را با استفاده از ترکیب راه حل های فعلی می سازد و در هر بار نسل جدید را جایگزین بخشی از نسل قدیم می کند.

شبه کد شکل 6-2 الگوریتم ژنتیک را بصورت مختصر شرح می دهد:

¹ Fitness function

```

Function genetic algorithm (population, fitness_fn,cross_over, mutation)
return an individual
Input: population a set of individual
Fitness_fn: a function that measure the fitness of an individual
Cross_over: is a percent of population that use for child produce (cross over)
Mutation: is a percent of population that use for mutation
Population<- InitPopulation();
Repeat
new_population <-empty set;
crossoverSize<-crossSize(Cross_over,population.size);
mutationSize<-mutationSize(Mutation,population.size);
new_population<- generateNewPopulation(population,
crossoverFunction(crossoverSize,population))
mutationFunction(new_population,mutationSize);
new_population<-bestchildof(population);
population<-new_population;
Until some individual is fit enough, or enough time has
elapsed return the best individual in population.

```

شکل 6-2 شبیه کد الگوریتم ژنتیک

1-3-2-6 ویژگی های الگوریتم ژنتیک

- الگوریتم های ژنتیک در مسائلی که فضای جستجوی بزرگی داشته باشند می تواند بکار گرفته شود. همچنین در مسائلی با فضای فرضیه پیچیده که تاثیر اجزا آن در فرضیه کلی ناشناخته باشند می توان از GA برای جستجو استفاده نمود.
- در بهینه سازی گسسته بسیار مورد استفاده قرار می گیرد.
- الگوریتم های ژنتیک را میتوان براحتی بصورت موازی اجرا نمود از اینرو می توان کامپیوترهای ارزان قیمت تری را بصورت موازی مورد استفاده قرار داد.
- امکان به تله افتادن این الگوریتم در مینیمم محلی کمتر از سایر روشهاست.

- از لحاظ محاسباتی پرهزینه هستند.

- تضمینی برای رسیدن به جواب بهینه وجود ندارد.

6-3-2-2 خصوصیات الگوریتم ژنتیک پیاده سازی شده

همانطور که در شکل 6-2 نشان داده شده است. تابع الگوریتم ژنتیک ارائه شده دارای چهار پارامتر است:

1. جمعیت: برای این پارامتر مجموع ضرایب را بصورت یک رشته 24 بیتی در نظر گرفته ایم و از آنجا که 3 ضریب بایستی مقدار دهی شوند هر ضریب یک بایت را اشغال می کند.

2. تابع تناسب: این تابع بر طبق تابع SERE به صورت رابطه 6-6 تعریف می شود

Fitness function = $Se_1.weight * SERE(word \text{ in } Se1.result)$

+ $Se_2.weight * SERE(word \text{ in } Se2.result)$

.

.

.

+ $SeN.weight * SERE(word \text{ in } SeN.result)$

6-6

که در این تابع برای هر موتور جستجوی دلخواه 1تا n وزن دلخواهی می توان در نظر گرفت.

استفاده از این وزن دلخواه باعث می شود که بتوانیم برای موتورهای گوناگون میزان اهمیت

متفاوت قائل شویم.

3. اپراتور ژنتیک استفاده شده برای آمیزش در این مسئله، یک تابع ساده آمیزش است که نیمی از

بایت های ضریبی را با نیم بایت همان ضریب در جمعیت دیگر ترکیب می کند.

4. اپراتور ژنتیک برای جهش: برای تاثیر گذاری منطقی اپراتور جهش روی بیت های 3تا 5 هر بایت

(برای هر ضریب) بصورت تصادفی بعد از هر 100 بار آمیزش عمل می کند.

از آنجا که تابع تناسب ارائه شده در این بخش ارزش یک کلمه کلیدی را در نتایج موتورهای جستجو ارائه می کند از این تابع تناسب به عنوان یک معیار دیگر سنجش کارایی استفاده کرده ایم و آن را **میزان دسترسی پذیری معرفی می کنیم.**

3-6 فاز تست

در این فاز دیگر از تابع SERE استفاده نمی کنیم و ضرایب هر آنچه که در مرحله قبل بدست آمده باشد ثابت می ماند.

1-3-6 استخراج نتایج

در این بخش یک سند ورودی پس از انجام مراحل پیش پرداز و امتیاز دهی کلمات توسط تابع با استفاده از مرز بدست آمده در مرحله قبل لیستی از کلمات را به عنوان کلمات کلیدی کاندید مطرح می کند. این لیست تعداد کلمات زیادی دارد و به خاطر فرایند استخراج و مرز تعیین شده از میزان بازخوانی خوبی برخوردار است اما میزان دقیقت این لیست در مقدار پایینی قرار دارد. در واقع اگر به دنبال کسب نتایج مناسب در یک موتور جستجو هستیم می توان کار را به این لیست یا تعداد مناسبی از کلمات بالای لیست بسنده کرد اما اگر به دنبال ملاک های علمی تر و نه عملی تر دقیقت و بازخوانی هستیم به طوریکه فرایند استخراج به رفتار انسانی شبیه تر باشد استفاده از فرایند پس پردازش بعد ضروری به نظر می رسد.

3-2-3 پس پردازش برای بهبود میزان دقت با استفاده از کلمات برجسته ساز

همانطور که قبلا هم ذکر شد انتخاب کلمات براساس رفتار انسانی باعث افزایش میزان دقت و بازخوانی در نتایج خواهد شد. در این مرحله با استفاده از کلمات برجسته ساز¹ روش جدیدی را برای گزینش کلمات از یک لیست کلمات کاندید جهت بهبود میزان دقت ارائه داده ایم.

کلمات برجسته ساز عباراتی از قبیل «در این مقاله» هستند که معمولا در همان جمله شامل این عبارت قبل یا بعد از آن، عباراتی که ایده اصلی یک متن را در خود دارند می آید. یکی از محاسن استفاده بیشتر از حوزه های نزدیک به این کلمات برای استخراج کلمات کلیدی می تواند کاهش بسیار زیاد محاسبات به علت محدود شدن ناحیه های جستجو و امتیازدهی باشد.

تعريف ما از کلمات برجسته ساز اینست:

کلماتی که احتمال حضور کلمات کلیدی منتخب کاربر انسانی در همسایگی آن ها بیشتر باشد.

در ادامه به شرح روش ابتکاری استخراج این کلمات و نیز نحوه به کارگیری آن ها می پردازیم اما برای شروع ابتدا بهتر است تعريف تعدادی اصطلاح که در ادامه متن زیاد به کار گرفته می شوند را تعريف کنیم.

همسایگی چپ یا **LN²**: به کلمات واقع شده در سمت چپ یک کلمه مشخص x (در اینجا یک کلمه برجسته ساز) همسایگی چپ آن کلمه می گوییم. رابطه همسایگی چپ در یک زبان از «راست به چپ نویس» مثل فارسی، از کلمه قبل از x شروع شده و تا پایان جمله قبل ادامه می یابد.

¹ Attention attraction strings

² Left neibor

همسایگی راست یا RN¹: کلمات واقع شده در سمت راست یک کلمه مشخص y (در اینجا یک کلمه برجسته ساز) را همسایگی راست آن کلمه می‌گوییم. رابطه همسایگی راست در یک زبان از «راست به چپ نویس» از کلمه بعد از y شروع شده و تا انتهای جمله ادامه می‌یابد.

کلمات هم رخداد²: کلماتی که در همسایگی با یک کلمه کلیدی خاص و نه همه کلیدی حضور دارد. به عنوان مثال کلمه اگر در مجموعه متون آموزش ما عبارت «فرشچیان» یک کلمه کلیدی منتخب انسانی باشد. کلمه «استاد» یک کلمه هم اتفاق محسوب می‌شود نه یک کلمه برجسته ساز، چرا که به رغم حضور در همسایگی کلمه کلیدی تنها به همسایگی یک کلمه کلیدی خاص محدود می‌شود.

6-3-2-1 استخراج کلمات برجسته ساز

شکل 6-3 روند استخراج کلمات برجسته ساز در این پروژه را نشان می‌دهد. ابتدا تمامی کلماتی که در مجموعه آموزش، یک رابطه RN یا LN با یک کلمه کلیدی برقرار می‌کنند را لیست کرده ایم و آن‌ها را به عنوان کلمات برجسته ساز احتمالی برچسب گذاری می‌کنیم. اصطلاح احتمالی به این خاطر است که این کلمات ممکن است با کلماتی که با یک کلمه کلیدی خاص هم‌رخداد باشند اشتباه گرفته شوند. برای جلوگیری از این تداخل، لیست کلمات برجسته ساز احتمالی را با استفاده از رابطه 6-7 لیست را امتیازدهی کرده ایم.

$$AAS-score(w) = \frac{D_w}{D_T} * \frac{NRN_w + NLN_w}{NW O_w} \quad 7-6$$

D_w: تعداد اسنادی که شامل کلمه w هستند. •

D_T: تعداد کل اسناد مورد تحقیق •

¹ Right niebor

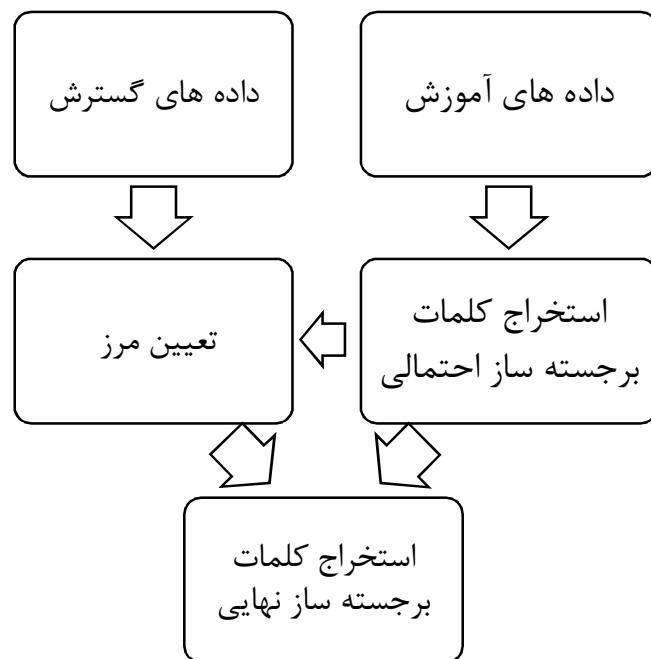
² Co-occurrence

NRN: تعداد دفعاتی که کلمه برچسب RN گرفته است. •

NLN: تعداد دفعاتی که کلمه برچسب LN گرفته است. •

NWO: تعداد دفعات وقوع کلمه •

از آنجا که کلمات برجسته ساز کلمات عمومی تری نسبت به کلمات هم‌رخداد محسوب می‌شوند قسمت اول رابطه 7-6 باعث می‌شود که این کلمات امتیاز بهتری نسبت به کلمات هم‌رخداد بگیرند و قسمت دوم مطابق با تعریف باعث می‌شود که کلماتی که در همسایگی یک کلمه کلیدی قرار می‌گیرند از امتیاز بالاتری برخوردار شوند.



شکل 6-3 روند استخراج کلمات برجسته ساز

بعد از این مرحله نیاز به تعریف یک مرز امتیازی برای تشخیص کلمات هم‌رخداد و کلمات برجسته‌ساز از یکدیگر است. اگر مرز تعریف شده در حد پایینی باشد کلمات هم‌رخداد نیز جزو کلمات برجسته‌ساز محسوب شده و نتایج را مخدوش می‌کنند و اگر این مرز در حد بالایی باشد ممکن است برخی از کلمات

برجسته ساز را از دست بدھیم. به همین خاطر به نوعی بهینه‌سازی برای تعیین این مرز نیازمندیم. برای

این منظور تابع مرز را بصورت رابطه ۸-۶ تعریف می‌کنیم.

$$KWNA(\text{threshold}) = \frac{\text{Number of keywords with no AAS neighbor}}{\text{Number of total keywords}} \quad 8-6$$

این تابع یک تابع گسسته است که دامنه آن را مرز انتخابی (threshold) در بازه $[0, 1]$ و برد آن را نسبت کلمات کلیدی که با این مرز، در مجموعه توسعه^۱، بدون همسایگی کلمات برجسته ساز احتمالی می‌ماند، به کل کلمات کلیدی، تعیین می‌کند.

مقدار مرزی که این تابع را کمینه کند یک مقدار بهینه محسوب می‌شود چرا که بیشترین میزان همسایگی بین کلمات کلیدی و کلمات برجسته ساز را در مجموعه توسعه ایجاد کرده است. برای پیدا کردن این کمینه از روش مطرح شده در مرجع [47] با استفاده از ابزار بهینه سازی² در نرم افزار Matlab استفاده کرده‌ایم.

۲-۳-۶ گزینش کلمات کلیدی

بعد از آنکه کلماتی که امتیازی بالاتر از مرز بهینه گرفته اند را به عنوان کلمات برجسته ساز انتخاب کرده ایم و برای هریک از این کلمات احتمال حضور یک کلمه کلیدی در همسایگی چپ و راست آن هارا با استفاده از روابط ۹-۶ بدست می‌آید.

¹ Develop database

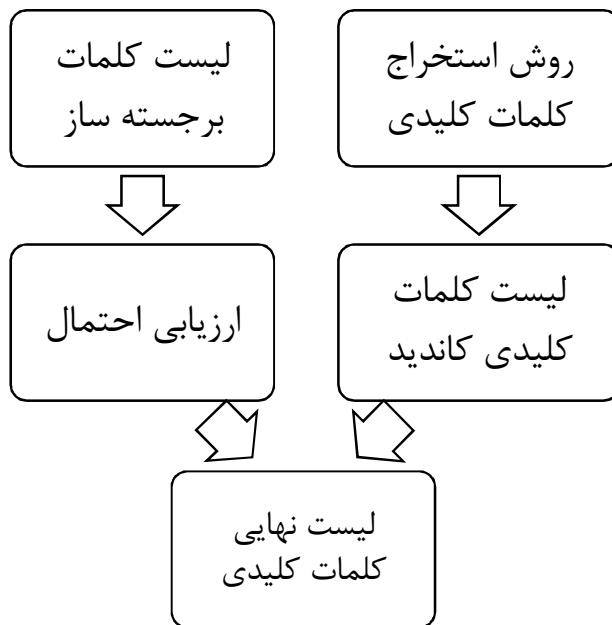
² Optimization tool

$$\text{EPKLN}(\text{AAS}_i) = \frac{\text{number of } \text{AAS}_i \text{ get a LN label}}{\text{total number of } \text{AAS}_i \text{ occurrence}}$$

9-6

$$\text{EPKRN}(\text{AAS}_i) = \frac{\text{number of } \text{AAS}_i \text{ get a RN label}}{\text{total number of } \text{AAS}_i \text{ occurrence}}$$

پس پردازش ارائه شده در این پژوهه مستقل از شیوه استخراج لیست کلمات کلیدی کاندید عمل می‌کند و برای افزایش میزان دقت هر روشی می‌توان از آن استفاده کرد. شکل ۶-۴ فرایند گزینش نهایی کلمات کلیدی را نشان می‌دهد.



شکل ۶-۴ فرایند گزینش نهایی کلمات کلیدی

براساس قاعده احتمال کل^۱ اگر Ω باشد برای هر واقعه دلخواه A از Ω داریم:

^۱ Total Probability Theorem

$$P(A) = \sum_{i=1}^m P(A | B_i) P(B_i) \quad 10-6$$

در صورتیکه همسایه های چپ و راست یک کلمه برجسته ساز را افزای حساب کنیم احتمال وقوع همسایه شدن یک کلمه کلیدی را می توان با استفاده از رابطه 11-6 بدست آورد.

$$PKEWN(W) = \frac{EPKLN_w + EPKRN_w}{2} \quad 11-6$$

حال با استفاده از 11-6 تابع امتیاز دهی 12-12 را بصورت زیر تعریف کنیم

$$PS(CKW) = \sum_{j=1}^n PKEWN(W_j) \quad 12-6$$

در رابطه 12-6 CKW به معنای کلمه کلیدی کاندید و n تعداد کلمات برجسته سازیست که با CKW رابطه RN یا LN دارند.

تابع امتیاز 12-6 یک تابع احتمال نیست اما هرچه خروجی این تابع بیشتر باشد احتمال اینکه کلمه ورودی، یک کلمه کلیدی منتخب کاربر باشد بیشتر است.

3-3-6 ارزیابی نهایی

بدست آوردن تمام ملاک‌های ارزیابی در حد ایده آل اگر ناممکن نباشد کاری دشوار است اما آنچه که در اینجا ارائه شده است حق انتخاب را به کاربر می‌دهد که تا ضمن کسب میزان قابل قبولی از هر ملاک به هر میزان که نیاز وجود دارد در یک ملاک خاص امتیاز بالاتری کسب کند آنچه که در فصل بعد می‌آید بررسی کارایی روش توضیح داده شده بر روی پایگاهداده‌ها و نتایج عملی بدست آمده از این روش است.

فصل هفتم:

نتائج عملی

1-7 مقدمه

در این فصل ابتدا به تحلیل میزان تاثیر گذاری اندازه پایگاه داده بر فرایند آموزش تابع امتیازدهی کلمات کلیدی پرداخته ایم. سپس با توجه به نتایجی که در تابع تناسب به دست آورده ایم میزان تاثیر گذاری هریک از پارامترهای در نظر گرفته شده برای تابع امتیازدهی را بدست می آوریم. همچنین میزان تاثیرگذاری روش پسپردازش ابتکاری نیز برروی روش ارائه شده و دو روش ساده دیگر مورد بررسی قرار گرفته است.

شایان ذکر است که در تمامی نتایجی بدست آمده، از 10 سایت اول، اولین صفحه نتایج چهار موتور جستجو¹, Google², yahoo³ و MSN⁴ با وزن های یکسان استفاده شده است. این مسئله می تواند در تحقیقاتی که بر روی یک موتور جستجوی خاص تمرکز بیشتری دارند بصورت استفاده صرف از این موتور جستجو یا استفاده بیشتر از نتایج ارائه شده توسط این موتور یا وزن دهنی سنگین تر برای موتور جستجوی مطلوب مطرح شود.

7-2 بررسی تاثیرگذاری اندازه پایگاه داده

همانطور که در فصل چهارم گفته شده جهت انجام این پروژه، برای 1100 سند از پایگاه داده همشهری کلمات کلیدی انتخاب کردہ ایم که در این میان از 600 سند برای مرحله استخراج کلمات کلیدی استفاده می کنیم.

¹ www.google.com

² www.yahoo.com

³ www.MSN.com

⁴ www.altavista.com

برای تعیین میزان تاثیر حجم پایگاه داده آموزش برای فرایند استخراج لیست کلمات کلیدی از 600 سندی که برای آنان کلمات کلیدی را انتخاب کرده ایم 100 سند را به عنوان مجموعه تست انتخاب کرده ایم و سپس با مجموعه های مختلف پایگاه داده به آموزش تابع امتیازدهی پرداخته ایم. جدول 7-1 نشان دهنده نتایج متفاوت با تعداد مختلف مجموعه آموزش است.

برای آنکه آزمایشات مستقل از هم باشند برای هر مجموعه آموزش بطور جداگانه کل سیستم را از ابتدا راهاندازی کرده‌ایم، همچنین در هر بار انتخاب مجموعه آموزش برای آنکه پاسخ‌های سیستم به مجموعه تست مستقل از اسناد باشد سندهای تست و آموزش بصورت تصادفی انتخاب شده‌اند.

همانطور که مشاهده می‌شود بهترین نتایج سیستم در حالتی بدست می‌آید که از 400 سند مجموعه پایگاه داده برای آموزش استفاده کرده‌ایم. این نتایج در مقایسه عددی با سایر تحقیقات استخراج کلمات کلیدی [35] و [44] و [48] نتایج قابل قبولی به شمار می‌رود اگرچه بایستی در نظر گرفت که شرایط و حتی زبان مورد تحقیق تفاوت دارد.

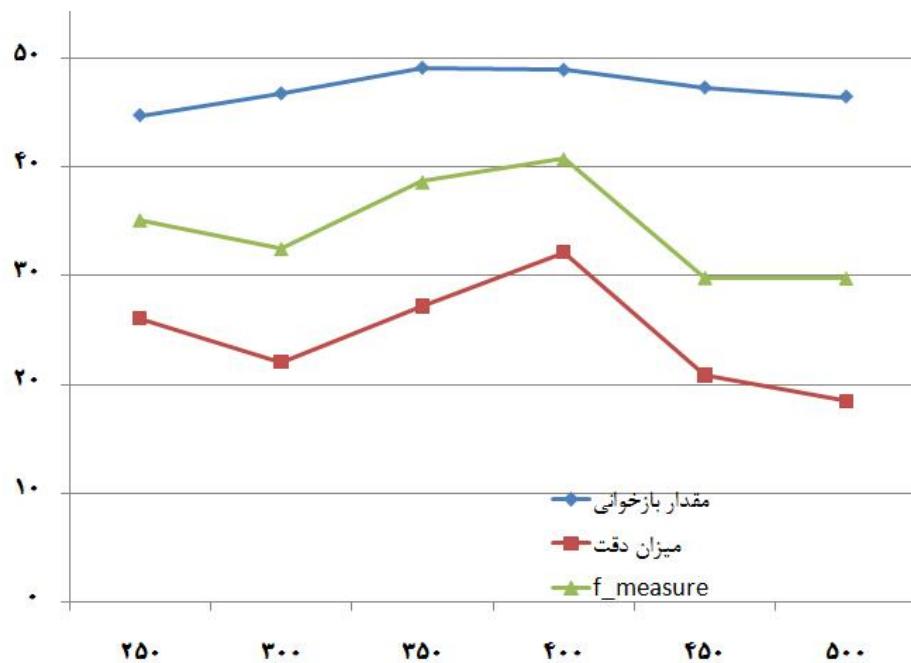
همانطور که در این نتایج می‌توان مشاهده کرد (شکل 7-1) آموزش بیش از حد باعث افزایش دنباله روی ضرایب برای تقلید از رفتار موتور جستجو شده و در نتیجه میزان دقت کاهش پیدا می‌کند. در واقع با وجود اینکه انتظار می‌رود افزایش اسناد مجموعه آموزش تاثیر مثبتی بر روند استخراج کلمات کلیدی بگذارد این اتفاق نیفتاده است چرا که الگوریتم موتور جستجو نیز کاملاً بر فرایند انتخاب انسانی منطبق نیست و در نتیجه روند استخراج از نوع رفتار انسانی فاصله گرفته است.

هرچند به طور کلی با توجه به اعوجاجی که رفتار میزان دقت و در نتیجه آن رفتار $F_measure$ نمی‌توان بطور قطع اظهار نظر کرد چراکه تعداد اسناد به کار گرفته شده در این مسئله آنقدر نیست که ثبات رفتاری توابع قابل اثبات باشد.

نکته ای دیگری که حائز اهمیت است ثابت ماندن تقریبی معیار بازخوانی با وجود ایجاد تغییرات در پایگاه داده است که البته به خاطر نوع تعریف مرزیست که ارائه کرده ایم.

جدول 7-1 تأثیرات اندازه پایگاه داده بر ملاک ارزیابی کارایی

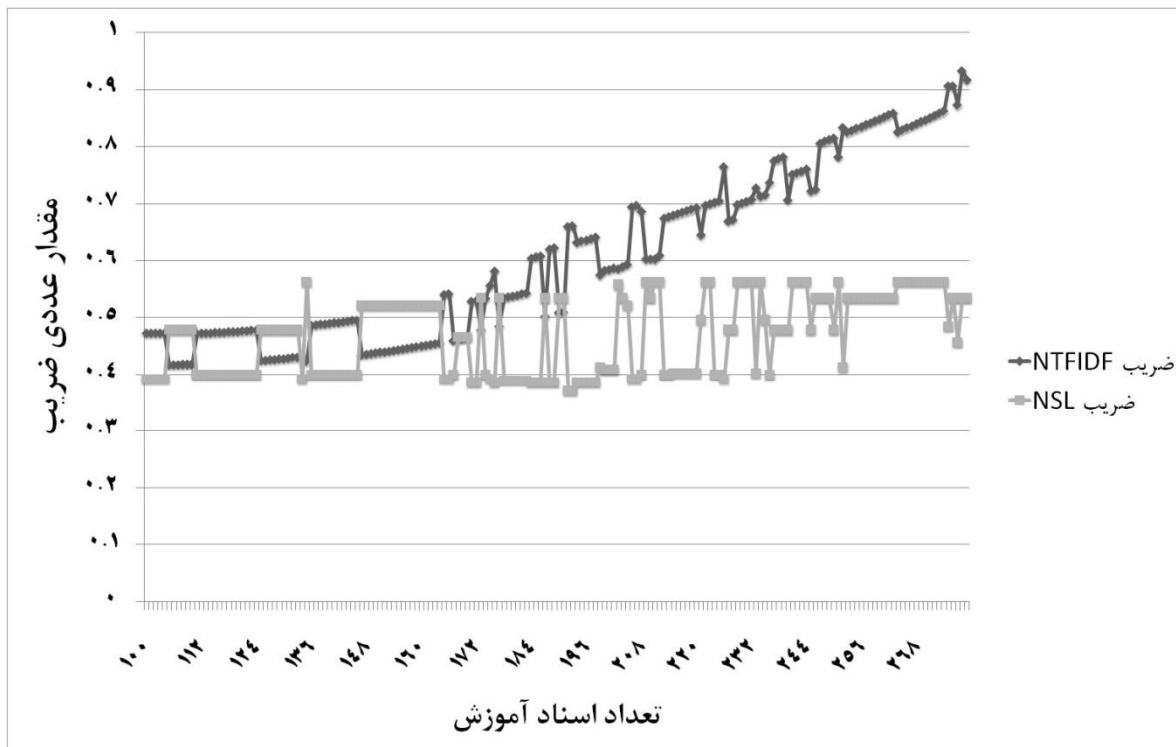
میانگین F	میانگین دقت	میانگین بازخوانی	تعداد اسناد مجموعه آموزش
33,18	26,14	44,77	250
34,54	22,09	46,78	300
38,51	27,26	49,14	350
40,82	32,21	48,95	400
29,86	20,9	47,32	450
29,8	18,51	46,49	500

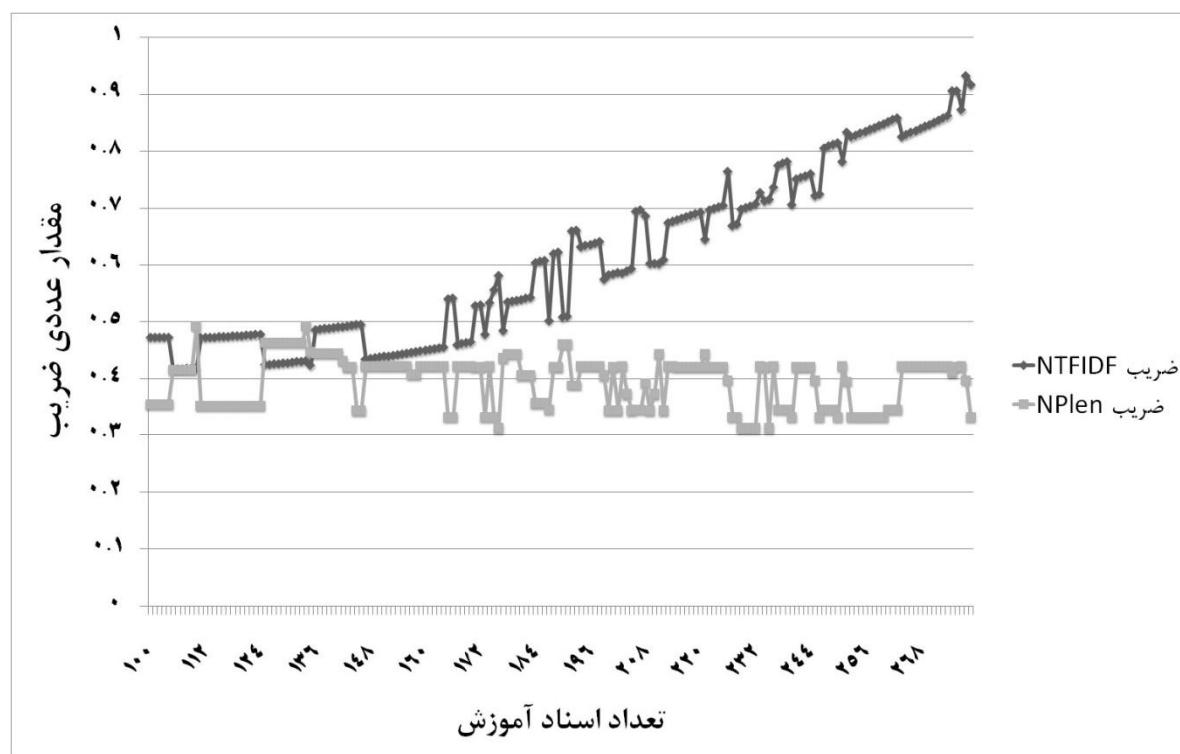


شکل 7-1 نحوه تغییرات پارامترهای کارایی در اثر تغییرات تعداد اسناد پایگاه داده آموزش

7-3 میزان تاثیر گذاری هریک از خصوصیات بر دسترسی پذیری

به نظر می‌رسد که می‌توان با استفاده از تطابق رشد ضرایب با توجه به صعودی مطلق بودن مقدار تابع تناسب در طول فرایند آموزش، میزان اهمیت تاثیر گذاری هریک از خصوصیات استخراج شده در این تابع را بررسی کرد (خصوصیات تعریف شده در صفحه 78 و 79). همانطور که در شکل 7-2 ملاحظه می‌کنید در طی فاز آموزش ۷ بازه بزرگتری نسبت به β می‌گیرد که این نشان از تاثیر گذاری بیشتر NTFIDF نسبت به NSL دارد.



شكل 7-3 مقایسه رشد γ (ضریب NTFIDF) با α (ضریب NPlen) در طول فرایند آموزش

7-4 نتایج فرایند پس پردازش با استفاده از کلمات بر جسته ساز

برای آنکه فرایند پس پردازش استقلال کامل خود را از فرایند استخراج کلمات کلیدی حفظ کند در این مرحله از 500 سندی که در حین فرایند استخراج نقشی نداشته اند استفاده کرده ایم، 450 سند را به عنوان پایگاه داده آموزش ، 50 سند را به عنوان مجموعه توسعه¹ در نظر گرفته ایم . توجه کنید از آنجا که این روش به عنوان پیشپردازش بعد از استخراج کلمات کلیدی به کار گرفته می شود نیازی به مجموعه تست جداگانه وجود ندارد. برای بررسی تاثیرات میزان حجم پایگاه داده بر عملکرد فرایند پس پردازش با حجم های مختلفی به ساخت پایگاه داده پرداخته ایم و مثل مرحله قبل جهت جلوگیری از تداخل نتایج آزمایشات قبل هر بار برای هر حجم پایگاه داده سیستم را مجددا بازسازی کرده ایم.

¹ Develop set

همچنین برای نشان دادن استقلال فرایند بهبود از روش استخراج، سه روش ساده دیگر براساس مراجع [39]، [40]، [44]¹ طراحی و به بهبود نتایج آن‌ها نیز پرداخته‌ایم. روش‌های مذکور با نام‌های ۱ و ۲ و ۳ در جدول آورده شده‌اند. برای هریک از روش‌های ۱ و ۲ و ۳ از مجموع ۲۵۰ سند برای آموزش استخراج کلمات کلیدی استفاده کرده‌ایم. از آنجا که فرایند بهبود به افزایش میزان دقت می‌انجامد سعی کرده‌ایم لیست تولید شده توسط روش‌های حاوی تعداد کلمات بیشتری باشد تا بازخوانی مناسبی ایجاد شود. همچنین نتایج بهبود برای روش ارائه شده در این پژوهه به نام روش ۴ آورده شده است برای روش ۴ از تعداد ۴۰۰ سند، که بهترین نتایج را به همراه داشته است، به عنوان مجموعه آموزش استفاده کرده‌ایم همچنین تعداد ۱۰۰ سند به عنوان داده تست برای تولید نتایج استفاده شده است.

لازم به ذکر است که این تحقیق قصد مقایسه نتایج بدست آمده از روش‌ها با یکدیگر را ندارد و تنها به بررسی میزان بهبود در هر روش به عنوان سندی بر درستی ادعای مستقل بودن پیش پردازش ارائه شده از روش استخراج کلمات کلیدی پرداخته است.

جدول 7-2 نتایج بدست آمده برای فرایند پس پردازش با حجم مختلف پایگاه داده طی فاز آموزش برای روش ۱

با استفاده از پسپردازش			عدم استفاده از پسپردازش			حجم پایگاه داده
F_measure	دقت	بازخوانی	F_measure	دقت	بازخوانی	
29,2	27,91	29,16	27,54	23,52	30,67	200
29,7	32,68	25,63	26,62	25,27	28,03	300
30,35	31,64	28,34	25,39	19,96	29,73	400

¹ نگارنده هرگز ادعای دوباره طراحی کردن روش‌های مذکور را ندارد و سه روشی که در نتایج از آن‌ها بهره‌گیری شده تنها براساس پایه‌های این تحقیقات به ویژه خصوصیات استفاده شده در آن‌ها برای استخراج کلمات کلیدی بنا نهاده شده است. در نتیجه اختلافات فاحش نتایج به همین دلیل عدم تطابق طراحی می‌باشد استفاده از این روش‌ها تنها برای نشان دادن میزان بهبود است.

جدول 7-3 نتایج بدست آمده برای فرایند پس پردازش با حجم مختلف پایگاه داده طی فاز آموزش برای

روش 2

با استفاده از پسپردازش			عدم استفاده از پسپردازش			حجم پایگاه داده
F_measure	دقت	بازخوانی	F_measure	دقت	بازخوانی	
27,42	25,4	29,67	27,11	20,87	31,89	200
25,45	25,28	25,58	22,88	18,09	27,14	300
29,33	30,21	27,27	24,32	19,7	29,52	400

جدول 7-4 نتایج بدست آمده برای فرایند پس پردازش با حجم مختلف پایگاه داده طی فاز آموزش برای

روش 3

با استفاده از پسپردازش			عدم استفاده از پسپردازش			حجم پایگاه داده
F_measure	دقت	بازخوانی	F_measure	دقت	بازخوانی	
27,42	25,4	29,67	27,11	20,87	31,89	200
26,88	29,59	24,19	25,77	20,64	28,16	300
29,77	31,05	28,34	23,03	17,46	32,02	400

جدول 7-5 نتایج بدست آمده برای فرایند پس پردازش با حجم مختلف پایگاه داده طی فاز آموزش برای روش

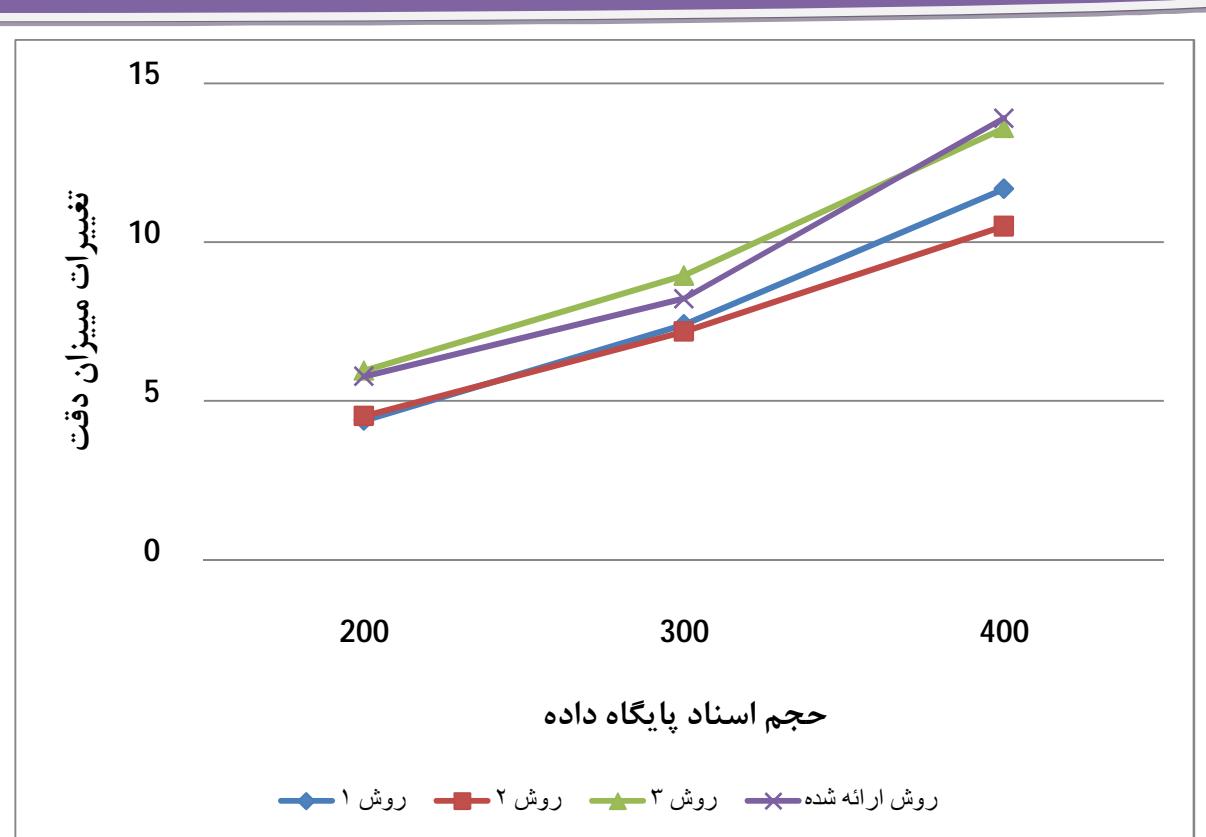
ارائه شده در این پژوهه

با استفاده از پسپردازش			عدم استفاده از پسپردازش			حجم پایگاه داده
F_measure	دقت	بازخوانی	F_measure	دقت	بازخوانی	
43,76	40,65	46,2	41,98	34,87	48,21	200
43,22	41,74	43,28	39,77	33,51	46,15	300
48,9	48,24	49,5	42,53	34,33	50,71	400

همانطور که در نتایج مشاهده می شود افزایش در میزان اسناد پایگاه داده در فرایند آموزش جهت استخراج کلمات برجسته ساز به طرز قابل توجهی موجب افزایش در معیار میزان دقت می شود. چراکه انتخاب هرچه دقیق تر این کلمات منجر به نزدیک شدن بیشتر گزینش کلمات به نقطه نظر انسانی می شود. این مطلب در نمودارهای صعودی افزایش میزان دقت بر حسب تعداد سند (شکل 4-7) و نیز در میزان تغییرات میانگین دقت جدول 7-6 قابل مشاهده است.

جدول 7-6 میزان تغییرات میانگین دقت در اندازه های مختلف پایگاه داده

400	300	200	روش
11,68	7,41	4,39	1
10,51	7,19	4,53	2
13,59	8,95	5,96	3
13,91	8,23	5,78	4



شکل 7-4-نمودار تغییرات میزان دقت بر حسب اندازه پایگاه داده

بهترین میزان بهبود در دقت زمانی رخ می دهد که از 400 سند در پایگاه داده استفاده می کنیم. البته در این مورد با کاهش اندکی (حداکثر 3,68%) در بازخوانی روبرو هستیم . در واقع پس پردازش معرفی، در اکثر موارد کلماتی را پالایش می کند که توسط کاربران به عنوان کلمه کلیدی شناخته نشده‌اند.

در پایان در جدول 7-7 میزان تغییرات در هر سه معیار کارایی را برای بهترین حالت، یعنی زمانی که استخراج کلمات برجسته ساز با استفاده از 400 سند صورت می گیرد ارائه شده است. توجه داشته باشید که میزان منفی ثبت شده برای میزان تغییرات معیار مورد بررسی بدست آمده است نه مثدار خود ملاک زیرا مقدار تمامی ملاک‌های مورد بررسی همواره صفر و یک می باشد.

جدول 7-7 میزان تغییرات در متوسط ملاک های کارایی در زمان استفاده از 400 سند

F_measure	دقت	بازخوانی	روش
4,96	11,68	-1,39	1
5,01	10,51	-2,25	2
6,74	13,59	-3,68	3
6,37	13,91	-1,21	4

فصل هشتم:

نتیجه گیری و

پیشنهاداتی برای آینده

1-8 مقدمه

این فصل شامل دو بخش می باشد:

1. نتیجه گیری که در آن نتایج بدست آمده از انجام این پژوهه بررسی و اعلام می شود
2. پیشنهاداتی برای آینده که به ارائه چند پیشنهاد برای تحقیقاتی می پردازد، که قصد ادامه این روش یا مقایسه روش معرفی شده با روش های دیگر را دارند.

2-8 نتیجه گیری

در این تحقیق روشی جهت استخراج خودکار کلمات کلیدی ارائه شده است که علاوه بر خصوصیات رایج، از بازخوردهای موتورهای جستجو نیز استفاده می کند. این روش از دوتابع امتیاز دهنی و یک الگوریتم ژنتیک بهره می گیرد. تابع امتیازدهی اول یک تابع انتخاب کلمات کلیدی با استفاده از فرایندهای آماریست که ضرایب بردار خصوصیات آن قابل تغییر تعریف شده است. تابع امتیاز دهنی دوم که در این تحقیق تابع ارزیابی نتایج موتورهای جستجو نامیده شده، یک تابع انتخاب کلمات کلیدی برپایه یک متن منحصر به فرد است که کلمات کلیدی بدست آمده از تابع اول را بر اساس نتایج موتورهای جستجو ارزیابی می کند. الگوریتم ژنتیک به کار گرفته شده در این روش، از یک تابع تناسب که برپایه تابع دوم بنانهاده شده، جهت بهینه سازی ضرایب تابع استخراج کلمات کلیدی کاندید می پردازد. نتایج بدست آمده نشان می دهد که روش ارائه شده علاوه بر انتخاب کلمات بالارزش از نظر موتورهای جستجو، در دو معیار میزان دقت و بازخوانی نیز امتیاز قابل قبولی کسب می کند.

همچنین در این پژوهه به تعریف کلمات برجسته ساز پرداخته ایم و برای چگونگی استخراج آن‌ها روش جدیدی به کار گرفته شده است. کلمات برجسته ساز طبق تعریف کلماتی هستند که احتمال وقوع یک کلمه کلیدی در همسایگی آن‌ها بیشتر است. پس از استخراج کلمات برجسته ساز با استفاده از قاعده احتمال کل، پس‌پردازش جدیدی جهت افزایش معیار دقت ارائه گردیده است. بررسی‌های عملی نشان می‌دهد که پس‌پردازش ارائه شده با کمی کاهش بازخوانی، تاثیر بسزایی بر معیار دقت دارد.

از دیگر تحقیقات انجام شده در این پژوهه بررسی میزان تاثیر حجم پایگاه داده آموزش، برای استخراج کلمات کلیدی و استخراج کلمات برجسته ساز می‌باشد. آزمایشات نشان می‌دهد که افزایش حجم پایگاه داده در روش ارائه شده برای استخراج کلمات برجسته ساز منجر به نتایج مطلوب‌تر می‌شود در حالیکه افزایش حجم پایگاه داده برای روش استخراج کلمات کلیدی باعث آموزش بیش از حد و پیروی رفتار تابع امتیاز دهی از رفتار موتورهای جستجو و در نتیجه کاهش دقت و بازخوانی می‌شود. هرچند نحوه رفتار تابع امتیازدهی آنقدر پایدار نیست که بتوان با قطعیت رای بر صحت نتیجه بدست آمده داد.

همچنین در بخش‌های پیش‌پردازش که به زبان مورد تحقیق وابستگی وجود دارد، ایده‌های تازه در فرایندهای توکنیزه سازی و حذف کلمات پر تکرار ارائه شده است.

3-8 پیشنهاداتی برای آینده

موارد زیر به منظور ادامه و تکمیل این پژوهه پیشنهاد می‌شود

3-8-1 ایجاد پایگاه داده بزرگتر برای اطمینان بخشی به نتایج بدست آمده

در بخش استخراج کلمات کلیدی اگرچه به نظر می‌رسد افزایش حجم پایگاه داده بیش از یک میزان بهینه، باعث ضعف در عملکرد تابع می‌شود. اما استفاده از پایگاه داده منجر به حرکت نرمتر و با ثبات‌تر

نمودار ملاک‌های کارایی روش بر حسب تعداد اسناد (شکل 1-7) شده و اطمینان بیشتری از نتایج بدست

آمده حاصل می‌شود.

همچنین در بخش استخراج کلمات برجسته ساز افزایش حجم پایگاه داده در طول فرایند آموزش باعث جامعیت لیست کلمات استخراج شده و بهبود نتایج خواهد شد.

همچنین پیشنهاد می‌شود این پایگاه داده‌ها استفاده از ساختارهای توزیع شده و تحت وب با استفاده از تعداد بیشتری از کاربران ایجاد شود. چرا که این امر منجر به جامعیت هرچه بیشتر پایگاه داده مذکور خواهد شد.

2-3-8 استفاده از یک روش ریشه‌یابی

اگرچه روش‌های ریشه‌یابی در زبان فارسی هنوز در ابتدای راه خود هستند اما استفاده از این یک روش ریشه‌یابی می‌تواند منجر به بهبودهای احتمالی در نتایج بدست آمده شود.

3-3-8 بررسی تاثیرات ایجاد یک مرز دینامیک امتیاز

از آنجا که در این پژوهه مقدار مرزی امتیاز جهت انتخاب کلمات کلیدی بر اساس ملاک بازیابی صورت گرفت پیشنهاد می‌شود روشی برای تعیین این مرز بر اساس طول و نوع ساختار هر سند برای تعیین کلمات کلیدی طراحی شود. پیش‌بینی می‌شود این نوع انتخاب به افزایش میزان دقت و برعکس کاهش ملاک کارایی منجر شود که در حالت کلی می‌توان آن را به عنوان یک روش تاثیر گذار بر دقت بررسی کرد.

مراجع

- [1] James Sanger and Ronen Feldman, "**The Text Mining Handbook advanced approaches in analyzing unstructured data**", Cambridge university press, 2007.
- [2] J. Kaur and V. Gupta, "**Effective Approaches For Extraction of Keywords,**" Journal of Computer Science, vol. 7, 2010, pp. 144-148.
- [3] C. Zhang, "**Automatic Keyword Extraction from Documents Using Conditional Random Fields,**" Journal of Computational Information Systems, vol. 3, 2008.
- [4] M.S. K. Taghva, R. Beckley, "**A Stemming Algorithm for the Persian Language,**" International Conference on Information Technology: Coding and Computing (ITCC 2005), 2005.
- [5] شمس فرد م, "پردازش متون فارسی: دستاوردهای گذشته، چالش های پیش رو" دومین کارگاه پژوهشی زبان فارسی و رایانه، صفحه 189-172، 1385
- [6] K. Taghva, J. Coombs, R. Pareda, and T. Nartker, "**Language Model-Based Retrieval for Farsi Documents,**" Internationa Conference on Information Technology: Coding and Computing (ITCC 2004), 2004.
- [7] A.N. Taghiyareh F., Darrudi E., Oroumchian F., "**Compression of Persian Text for Web-Based Applications, Without Explicit Decompression.**"
- [8] J.W. Amstrup and N.M.S.U.C.R. Laboratory, "**Persian-English Machine Translation : An Overview of the Shiraz Project**", Computing Research Laboratory, New Mexico State University, 2000.
- [9] S.M. Assi, "**Farsi Linguistic Database (FLDB),**" international Journal of Lexicography, vol. 10, 1997, pp. 5-10.

[10] K.S. Esmaili, H. Abolhassani, and M. Neshati, “**Mahak : A Test Collection for Evaluation of Farsi Information Retrieval Systems**,” Cities, 2007, pp. 639-644.

[11] A. AleAhmad, H. Amiri, E. Darrudi, M. Rahgozar, and F. Oroumchian, “**Hamshahri: A standard Persian text collection**,” Knowledge-Based Systems, vol. 22, 2009, p. 382–387.

[12] حافظی م م ، ثامتی ح ، منصوری ن ، منتظری ن ، بحرانی م ، موثق ح، ”**یک مدل دستوری برای بهبود دقیق سیستم‌های بازشناسی گفتار پیوسته‌ی فارسی** ”، دومین کارگاه پژوهشی زبان فارسی و رایانه، صفحه 91-80 ، 1385

[13] D.J.& J. H.Martin, “**Speech and Language Processing: An introduction to speech recognition, natural language processing, and computational linguistics**,” 2005, pp. 29-32.

[14] مدرس خیابانی ش ، قیومی م، ”**نقش پیکره‌های زبانی در باهم‌آیی: رویکردی مقایسه‌ای** ”، دومین کارگاه پژوهشی زبان فارسی و رایانه، صفحه 55-55-65، تهران 1385

[15] داداش میری پ. ”**تشخیص انتهای کلمات فارسی و ایجاد فاصله میان کلمات** ”، پایان نامه کارشناسی، دانشگاه علم و صنعت 1380.

[16] بیجن خان م، ”**تشخیص کسره اضافه** ”، طرح تحقیقاتی پژوهشگاه فرهنگ و هنر دانشگاه تهران 1384

[17] مرتضی آنالویی ، محسن مشکی، ”**یک روش آماری مبتنی بر پیکره برای جداسازی واژه‌های به هم چسبیده** ”، کنفرانس فازی و هوشمند. 2008 ،

- [18] P. Schäuble, **Multimedia information retrieval : content-based information retrieval from large text and audio databases**, Boston: MA:kluwer Academic publishers, 1997.
- [19] K. Taghva, R. Beckley, and M. Sadeh, **A list of Farsi stopwords**, las vegas: Citeseer, 2003.
- [20] M.R. Davarpanah, M. Sanji, and M. Aramideh, “**Farsi lexical analysis and stop word list**,” emerald group publishing, vol. 27, 2009, p. 435–449.
- [21] C.S. alahati Qadimi Fumani, M. R.; Ramachandra, “**The concept of stop words in Persian chemistry articles: A discussion in automatic indexing.**,” 2008.
- [22] اوری، ح. احمد گیوی .. دستور زبان فارسی 2، تهران: انتشارات فاطمی چاپ 21، 1380
- [23] ”**تحلیل سیستم یافتن خودکار کلمات کلیدی متون زبان فارسی**“،دانشگاه علم و صنعت، ویرایش اول ، 1388
- [24] M.F. Porter, “**an algorithm for suffix stripping**,” 1983.
- [25] B. Lovins, “**Development of a Stemming Algorithm**,” Mechanical transaltion and computational linguistics, vol. 11, 1968, pp. 22-31.
- [26] N.F. and M.M. Maristella Agosti, Michela Bacchin, “**Improving the Automatic Retrieval of Text Documents**,” Advances in Cross-Language Information Retrieval, vol. 2785, 2003.
- [27] M.M. Michela Bacchin, Nicola Ferro, “**A probabilistic model for stemmer generation**,” Information Processing & Management, vol. 41, 2005, pp. 121-137.
- [28] M.meybodi, “**Bon: First Persian Stemmer**,” First Eurasia Conference on information2, 2003.

[29] B. Esfahbod, “**FarsiTEX and the Iranian TEX Community**,” East Asia, vol. 23, 2002, pp. 41-45.

[30] مجتبی محمدی نصیری، کیومرث شیخ اسماعیلی، حسن ابوالحسنی .. ”یک ریشه یاب آماری برای زبان فارسی ”، مجموعه مقالات یازدهمین کنفرانس بین المللی کامپیوتر. (csicc), 1384.

[31] J.D. Cohen, ‘**Highlights : Language- and Domain-Independent Indexing Terms for Abstracting Automatic**,” journal of the american society for information science, vol. 46, 1995, pp. 162-174.

[32] Y. Matsuo and M. Ishizuka, “**Keyword extraction from a single document using word co-occurrence statistical information**,” International Journal on Artificial Intelligence Tools, vol. 13, 2004, p. 157–170.

[33] L.F. Chien, “**PAT-tree-based keyword extraction for Chinese information retrieval**,” ACM SIGIR Forum, ACM, 1997, p. 50–58.

[34] G. Ercan and I. Cicekli, “**Using lexical chains for keyword extraction**,” Information Processing & Management, vol. 43, Nov. 2007, p. 1705–1714.

[35] A. Hulth, “**Improved automatic keyword extraction given more linguistic knowledge**,” Proceedings of the 2003 conference on Empirical methods in natural language processing -, 2003, pp. 216-223.

[36] M.J.M.G. G. Salton, “**Automatic Text Structuring and Retrieval –Experiments in Automatic Encyclopaedia Searching**,” Fourteenth SIGIR Conference, 1991, pp. 21-30.

[37] “**Domain-Specific Keyphrase Extraction**,” 16th International Joint Conference on Arifical Intelliegence, Stockholm, Sweden: 1999, pp. 668-673.

[38] K. Zhang, H. Xu, J. Tang, and J. Li, “**Keyword extraction using support vector machine**,” Advances in Web-Age Information Management, 2006, p. 85–96.

- [39] I.H. Witten, G.W. Paynter, E. Frank, C. Gutwin, and C.G. Nevill-Manning, “**KEA: Practical automatic keyphrase extraction,**” Proceedings of the fourth ACM conference on Digital libraries, ACM, 1999, p. 254–255.
- [40] P.D. Turney, “**Learning algorithms for keyphrase extraction,**” NRC Technical Report ERB-1057, 2000.
- [41] J.B.K. Humphreys, “**Phraserate: An html keyphrase extractor,**” Riverside: University of California, Riverside, 2002, pp. 1-16.
- [42] تشكري.م، ”ساخت يك نمایه ساز خودکار برای متون فارسی،“ دانشگاه صنعتی اميرکبير، .1380
- [43] El-Shishtawy, T., Al-Sammak, A. “**Arabic Keyphrase Extraction using Linguistic knowledge and Machine Learning Techniques**”, Proceedings of the Second International Conference on Arabic Language Resources and Tools, 2009
- [44] J. Wang, H. Peng, and J.S. Hu, “**Automatic keyphrases extraction from document using Neural Network,**” Machine Learning and Cybernetics, 2005. Proceedings of 2005 International Conference on, IEEE, 2005, p. 3770–3774.
- [45] **Search Engine Optimization Starter Guide**, Google Inc, 2010.
- [46] C.T.Y. s 1 G. Salton, C. S. Yang, “**A Theory of Term Importance in Automatic Text Analysis,**” Journal of the American society for Information Science, vol. 26, 1975, pp. 33-44.
- [47] D.C. sorensen jorge J.more, “**Computing a Trust Region Step,**” SIAM Journal on Scientific and Statistical Computing, vol. 3, 1983, pp. 553-572.
- [48] R. Mihalcea and P. Tarau, “**textRank: bringing order into texts,**” Proceedings of EMNLP 2004, 2004, pp. 404-411.

Abstract

Generally, automatic keyword extraction methods attempt to get better results in criterions such as recall and precision, which show the ability of the method to perform the desired action as an author.

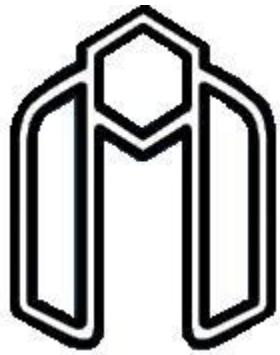
According to the undeniable role of search engines in today's world, it seems that keywords selection methods should pay the same attention to web content accessibility by search engines and the information retrieval criterions.

In this thesis, a new Automatic keyword extraction method is presented which acquires as good recall/precision as other methods and at the same time increases the contents accessibility by search engines. To extract keywords two score functions are used. The first function is a statistical keyword selector and the second one is a search engine results evaluator. The keyword selector function attempts to get a better recall/precision, while the search engine result evaluator function optimizes the first one, using genetic algorithm and search engine results during training phase.

The introduced keyword extraction method also uses a new post-processing phase which employs attention attraction strings to improve precision without losing noticeable recall.

Experimental results show that using proper documents in training phase, reasonable balance of the recall, precision and web content accessibility can be obtained.

Keywords: Automatic keywords extraction, Search engine optimization, SEO, Farsi stopwords, Genetic algorithm, Farsi attention attraction strings, Precision improvement, Accessibility improvement.



Shahrood University of Technology

Faculty of Computer and Information Technology

Optimized Keyword Selection for General Search

Engines

Hamze Hodhokian

Supervisor(s):

Dr.Morteza Zahedi

Dr.Hamid Hassanpur

July 2011