





دانشکده مهندسی برق و رباتیک

رشته مهندسی رباتیک

پایان نامه کارشناسی ارشد

شناسایی، طبقه بندی و یادگیری اشکال و تصاویر در بستر وب

توسط یک موتور جستجو

نگارنده : مهرداد صالحی

استاد راهنما

دکتر علیرضا احمدی فرد

شهریور ۱۳۹۶

## دانشگاه صنعتی شاهرود

### دانشکده مهندسی برق و رباتیک

#### گروه مهندسی رباتیک

پایان نامه کارشناسی ارشد آقای مهرداد صالحی

تحت عنوان: شناسایی، طبقه بندی و یادگیری اشکال و تصاویر در بستر وب توسط یک موتور جستجو

در تاریخ ۱۳۹۶/۶/۱۴ توسط کمیته تخصصی زیر جهت اخذ مدرک کارشناسی ارشد  
مورد ارزیابی و با درجه ..... مورد پذیرش قرار گرفت.

امضاء	اساتید مشاور	امضاء	اساتید راهنما
	نام و نام خانوادگی :		نام و نام خانوادگی : علیرضا احمدی فرد
	نام و نام خانوادگی :		نام و نام خانوادگی :

امضاء	نماینده تحصیلات تکمیلی	امضاء	اساتید داور
	نام و نام خانوادگی : مسعود مهدی زاده رخی		نام و نام خانوادگی : سید علی سلیمانی ایوری
			نام و نام خانوادگی : امیر حسین نایبی آستانه
			نام و نام خانوادگی :
			نام و نام خانوادگی :

تقدیم به :

پدر و مادر عزیز و مهربانم

که در سختی‌ها و دشواری‌های زندگی همواره یآوری دلسوز و فداکار و

پشتیبانی محکم و مطمئن برایم بوده‌اند.

تشکر و قدردانی :

جناب آقای دکتر علیرضا احمدی فرد که در تمام مراحل این پایان نامه  
بنده را مورد راهنمایی قرار دادند.

## تعهد نامه

اینجانب **مهرداد صالحی** دانشجوی دوره کارشناسی ارشد رشته مهندسی رباتیک دانشکده مهندسی برق و رباتیک دانشگاه صنعتی شاهرود نویسنده پایان نامه طبقه بندی و یادگیری اشکال و تصاویر در بستر وب توسط یک موتور جستجو تحت راهنمایی دکتر **علیرضا احمدی** فرد متعهد می شوم .

- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است .
- در استفاده از نتایج پژوهشهای محققان دیگر به مرجع مورد استفاده استناد شده است .
- مطالب مندرج در پایان نامه تاکنون توسط خود یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است.
- کلیه حقوق معنوی این اثر متعلق به دانشگاه صنعتی شاهرود 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 تأثیرگذار بوده اند در مقالات مستخرج از پایان نامه رعایت می گردد.
- در کلیه مراحل انجام این پایان نامه ، در مواردی که از موجود زنده ( یا بافتهای آنها ) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است .
- در کلیه مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است اصل رازداری ، ضوابط و اصول اخلاق انسانی رعایت شده است

### تاریخ

### امضای دانشجو

### مالکیت نتایج و حق نشر

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

چکیده :

با افزایش روزافزون تولید داده در بستر وب نیاز به یافتن اطلاعات مورد نیاز با صرف کمترین زمان و هزینه، اهمیت موتورهای جستجو را روز به روز پررنگتر می‌کند. البته یافتن لیستی بلند بالا از تصاویر و نتایج نامرتبط با موضوع مورد جستجو همواره باعث ناراحتی کاربران شده و باید زمان و وقت بسیاری را صرف کنند تا بتوانند موضوع مورد نیاز خود را بیابند. یکی از مهمترین کاربردهای موتور جستجو بررسی صفحات وب می‌باشد، برنامه‌ی هوشمند که بعد از جستجوی و استخراج اطلاعات موجود در بستر وب به پردازش و گروه بندی اطلاعات می‌پردازد و نتایج را به ترتیب شباهت و اهمیت با موضوع مورد جستجو با یک اولویت خاص به کاربر نمایش می‌دهد.

راه حل پیشنهادی این پروژه شامل استخراج متن (متون نزدیک به تصویر) و تصویر از صفحات وب با استفاده از یک موتور جستجوی ساخته شده محلی و یادگیری تصاویر با یک شبکه عصبی عمیق می‌باشد در این جا چالش های مانده ایجاد یک پایگاه داده اولیه از بستر وب و یادگیری شبکه بدون ناظر و گروه بندی متون و تصاویر وجود دارد. در این پروژه به معرفی چالش های یک موتور جستجوی وب با محوریت تصویر و روش هایی برای حل این مشکلات به همراه نوآوری های انجام گرفته در آن می‌پردازیم و در انتها با استفاده از یادگیری عمیق به استخراج ویژگی از تصاویر و متون داخل صفحات وب پرداخته و به طبقه بندی این نتایج جهت دسترسی سریعتر به اطلاعات و گروه بندی اشیا می‌پردازیم.

کلمات کلیدی :

موتور جستجو ، یادگیری عمیق، طبقه بندی تصاویر، خزنده وب، پردازش موازی، موتور جستجوی

تصویر.

## فهرست مطالب

عنوان	صفحه
فصل اول : پیش گفتار .....	۱
۱-۱ مقدمه .....	۲
۱-۱-۱ پیش گفتار .....	۲
۱-۱-۲ صفحات وب .....	۲
۱-۱-۳ HTML .....	۳
۱-۲ اهداف پروژه .....	۴
فصل دوم : معرفی موتور جستجوگر وب .....	۷
۲-۱ معرفی .....	۸
۲-۲ خزنده وب چیست .....	۸
۲-۳ کاربردهای خزنده وب .....	۹
۲-۳-۱ موتور های جستجو .....	۹
۲-۳-۲ مدیریت فنی وب سایت .....	۹
۲-۳-۳ جمع آوری اطلاعات خاص .....	۱۰
۲-۴ نحوه کارکرد یک خزنده وب چیست .....	۱۰
۲-۴-۱ معماری .....	۱۱



- ۲-۵ فرایندهای پیش پردازش یک خزنده وب ..... ۱۱
- ۲-۵-۱ عادی کردن آدرس ..... ۱۱
- ۲-۵-۲ بررسی درس صفحات وب ..... ۱۵
- ۲-۶ قسمت‌های مهم یک خزنده وب ..... ۱۷
- ۲-۶-۱ شاخص‌گذاری ..... ۱۷
- ۲-۶-۲ پرس‌وجو ..... ۱۸
- ۲-۶-۳ رتبه‌بندی ..... ۱۸
- ۲-۶-۴ Frontier crawler ..... ۱۹
- ۲-۶-۵ دریافت‌کننده ..... ۲۰
- ۲-۶-۶ مخزن ..... ۲۰
- ۲-۷ انواع موتور جستجوی وب ..... ۲۰
- ۲-۷-۱ خزنده وب موازی ..... ۲۱
- ۲-۷-۲ خزنده وب توزیع شده ..... ۲۲
- ۲-۷-۳ خزنده وب افزایشی ..... ۲۴
- ۲-۷-۴ خزنده وب مخفی ..... ۲۵
- ۲-۸ معرفی چند موتور جستجوی مهم ..... ۲۷

۲۸	..... robots.txt	۲-۹ معرفی فایل
۳۱	.....	فصل سوم: معرفی چالش ها، راه حل ها و نوآوری های موتور جستجو
۳۲	.....	۳-۱ مقدمه
۳۲	.....	۳-۲ دریافت و بررسی صفحات وب
۳۲	.....	۳-۲-۱ مقدمه
۳۲	.....	۳-۲-۲ معرفی غالب یک صفحه وب
۳۵	.....	۳-۲-۳ دریافت کننده صفحه وب
۳۷	.....	۳-۲-۴ تجزیه کننده صفحه وب
۳۷	.....	۳-۲-۴-۱ شناسایی و اولویت بندی تگ های بی استفاده
۴۲	.....	۳-۲-۴-۲ استخراج ویژگی از تگ
۴۴	.....	۳-۲-۴-۳ استخراج ویژگی از تصویر
۴۴	.....	۳-۲-۴-۳-۱ معرفی تصویر
۴۵	.....	۳-۲-۴-۳-۲ شبکه عصبی کانولوشن
۵۰	.....	۳-۲-۴-۴ استخراج ویژگی از متن
۵۷	.....	۳-۲-۵ صف اولویت آدرس ها
۵۸	.....	۳-۲-۵-۱ روش آقای گوپتا

۶۰	..... ۳-۲-۵-۲ روش موتور جستجوی گوگل
۶۱	..... page ranking developer روش ۳-۲-۵-۲
۶۲	..... ۳-۳ مدیریت پایگاه داده
۶۳	..... ۳-۴ مدیریت پرس و جو
۶۷	..... فصل چهارم: معرفی الگوریتم موتور جستجوی وب
۶۸	..... ۴-۱ مقدمه
۶۸	..... ۴-۲ دریافت لیست صفحات وب
۷۱	..... ۴-۲-۱ محاسبه امتیاز دامنه
۷۳	..... ۴-۳ مازول دریافت کننده صفحات وب
۷۴	..... ۴-۴ قسمت تجزیه کننده
۷۴	..... ۴-۴-۱ استخراج تصاویر و متن از صفحات وب
۷۵	..... ۴-۴-۲ بررسی تصاویر متن
۷۷	..... ۴-۴-۳ استخراج لینک های اتصالی
۷۸	..... ۴-۵ قسمت کنترل پایگاه داده
۸۰	..... ۴-۶ پایگاه داده موتور جستجو
۸۹	..... فصل پنجم : نتیجه گیری

۵-۱	مقدمه	۹۰
۵-۲	اهمیت مکانی متون در صفحات وب	۹۰
۵-۳	تاثیر نرمال سازی در کاهش پردازش صفحات	۹۳
۵-۴	تعیین امتیاز هر صفحه وب	۹۴
۵-۴-۱	دلایل اصلی کاهش امتیاز یک صفحه وب	۹۵
۵-۵	بررسی اهمیت متون نزدیک تصاویر در صفحه وب	۹۸
۵-۶	بررسی نتایج تست موتور جستجو	۱۰۰

## فهرست جدول ها

<u>صفحه</u>	<u>عنوان</u>
۲۸ .....	جدول ۱-۲
۳۴ .....	جدول ۱-۳
۳۶ .....	جدول ۲-۳
۵۶ .....	جدول ۳-۳
۸۱ .....	جدول ۱-۴
۸۲ .....	جدول ۲-۴
۸۳ .....	جدول ۳-۴
۸۵ .....	جدول ۴-۴
۸۵ .....	جدول ۵-۴
۸۶ .....	جدول ۶-۴
۸۸ .....	جدول ۷-۴
۹۶ .....	جدول ۱-۵
۱۰۰ .....	جدول ۲-۵

## فهرست اشکال

عنوان	صفحه
شکل ۱-۱ نمای کلی یک موتور جستجوی وب .....	۳
شکل ۱-۲ معماری سطح بالای خزنده وب.....	۱۱
شکل ۲-۲ معماری کامل یک خزنده وب .....	۱۹
شکل ۳-۲ الگوریتم کامل مراحل یک خزنده وب .....	۲۱
شکل ۴-۲ خزنده وب موازی .....	۲۲
شکل ۵-۲ خزنده وب توزیع شده .....	۲۴
شکل ۶-۲ خزنده وب افزایشی .....	۲۵
شکل ۷-۲ خزنده وب مخفی .....	۲۷
شکل ۱-۳ غالب کلی یک تگ html .....	۳۳
شکل ۲-۳ چند تگ مهم و پرکاربرد .....	۳۴
شکل ۳-۳ بلوک دیاگرام ماژول دریافت کننده .....	۳۷
شکل ۴-۳ مثالی از غالب html .....	۳۸
شکل ۵-۳ بررسی تگ ها در صفحات وب .....	۴۱
شکل ۶-۳ مکان تگ ها در درخت html .....	۴۲

- شکل ۳-۷ تصویر به تفکیک پیکسل ..... ۴۵
- شکل ۳-۸ شبکه عصبی CNN تصویر..... ۴۶
- شکل ۳-۹ شبکه عصبی چند لایه ..... ۴۸
- شکل ۳-۱۰ شبکه کانولوشن ۳ بعدی ..... ۴۹
- شکل ۳-۱۱ بردار ورودی شبکه کانولوشن متن ..... ۵۲
- شکل ۳-۱۲ بردار ورودی شبکه کانولوشن متن ..... ۵۲
- شکل ۳-۱۳ بلوک یاگرام یادگیری تصویر در کافی ..... ۵۳
- شکل ۳-۱۴ بلوک دیاگرام یادگیری متن در کافی ..... ۵۴
- شکل ۳-۱۵ نمونه‌ای از فایل prototxt در کافی..... ۵۴
- شکل ۳-۱۶ لینک‌های ورودی و خروجی صفحه وب ..... ۶۰
- شکل ۳-۱۷ شبکه CNN برای محاسبه طبقه بندی پرسوجوی جدید ..... ۶۵
- شکل ۴-۱ گروه بندی دامنه ها ..... ۶۹
- شکل ۴-۲ انتخاب دامنه درمورد یک پرسوجوی خاص..... ۷۲
- شکل ۴-۳ متون نزدیک تصویر ..... ۷۴
- شکل ۴-۴ شبکه کانولوشن تصویر ..... ۷۶
- شکل ۴-۵ شبکه کانولوشن متن ..... ۷۶

- شکل ۴-۶ بلوک دیاگرام ماژول کنترل پایگاه داده ..... ۷۹
- شکل ۴-۷ جدول دامنه ..... ۸۱
- شکل ۴-۸ جدول امتیاز دامنه ..... ۸۲
- شکل ۴-۹ جدول صفحه وب ..... ۸۳
- شکل ۴-۱۰ جدول تاریخچه ..... ۸۴
- شکل ۴-۱۱ جدول رابطه ..... ۸۴
- شکل ۴-۱۲ جدول تصویر ..... ۸۶
- شکل ۴-۱۳ جدول متن ..... ۸۷
- شکل ۴-۱۴ لیست روابط در پایگاه داده ..... ۸۸
- شکل ۵-۱ مقایسه مکان تصاویر دردانه ۱- ..... ۸۹
- شکل ۵-۲ مقایسه مکان تصاویر دردانه ۲- ..... ۹۱
- شکل ۵-۳ مقایسه مکان تصاویر دردانه ۳- ..... ۹۱
- شکل ۵-۴ مقایسه مکان تصاویر دردانه ۴- ..... ۹۲
- شکل ۵-۵ مقایسه مکان تصاویر دردانه ۵- ..... ۹۲
- شکل ۵-۶ مقایسه مکان تصاویر دردانه ۶- ..... ۹۳
- شکل ۵-۷ مقایسه نتیجه نرمال سازی صفحات وب ..... ۹۴



- شکل ۵-۸ متون نزدیک و مرتبط به تصویر در یک دامنه خاص ..... ۹۹
- شکل ۵-۹ بررسی ۱۰ صفحه درباره کلاغ زاغی ..... ۱۰۱
- شکل ۵-۱۰ بررسی ۱۰ صفحه درباره مارمولک فرشته ..... ۱۰۱
- شکل ۵-۱۱ بررسی ۱۰ صفحه درباره صندلی راحتی ..... ۱۰۲
- شکل ۵-۱۲ بررسی ۱۰ صفحه درباره کیسه خواب ..... ۱۰۲
- شکل ۵-۱۳ بررسی ۱۰ صفحه درباره دمنده دستی ..... ۱۰۳
- شکل ۵-۱۴ بررسی ۱۰ صفحه درباره طوطی سفید ..... ۱۰۳
- شکل ۵-۱۵ بررسی ۱۰ صفحه درباره دُمبِل ..... ۱۰۴
- شکل ۵-۱۶ بررسی ۱۰ صفحه درباره ماشین حمل زباله ..... ۱۰۴
- شکل ۵-۱۷ بررسی ۱۰ صفحه درباره ماشین ipod ..... ۱۰۵

## فهرست فرمول‌ها

<u>عنوان</u>	<u>صفحه</u>
فرمول ۱-۳.....	۴۹
فرمول ۲-۳.....	۵۷
فرمول ۳-۳.....	۶۰
فرمول ۱-۴.....	۶۹
فرمول ۲-۴.....	۷۲
فرمول ۳-۴.....	۷۲
فرمول ۴-۴.....	۷۲
فرمول ۱-۵.....	۹۷

فصل اول:

پیش گفتار

## ۱-۱ مقدمه

### ۱-۱-۱ پیش گفتار

با افزایش روز افزون تولید اطلاعات اینترنت<sup>۱</sup> طبقه بندی و دستیابی سریع به داده ها اهمیت بسیاری پیدا کرده است. در این حجم داده جدا سازی اطلاعات مفید از این حجم داده بسیار سخت می باشد. هر روزه در دنیا میلیون ها سایت<sup>۲</sup> با موضوعات و اطلاعات مختلف ایجاد می گردد که ما از آنها هیچ اطلاعاتی نداریم. پس چطور در این فضا اطلاعاتی را که مورد نظر ما است را کسب کنیم. برای این راهکار موتورهای جستجو<sup>۳</sup> متولد شدند. وظیفه اصلی یک موتور جستجو جمع آوری اطلاعات از بستر اینترنت، طبقه بندی داده های دریافتی در یکسری گروه مشخص و نمایش لیست نتایج مرتبط با موضوع مورد جستجوی به کاربر درخواست کننده می باشد. در حقیقت وظیفه اصلی یک موتور جستجو جمع آوری سیل عظیم این داده ها در غالب مشخص و سهولت دستیابی اطلاعات بستر وب می باشد.

### ۲-۱-۱ صفحات وب<sup>۴</sup>

منظور از صفحات وب همان مستندات و فایل های است که در بستر وب و توسط مرورگرهای وب<sup>۵</sup> قابل نمایش هستند. کنسرسیوم جهانی وب<sup>۶</sup> قالب استاندارد HTML را برای نمایش این مستندات در مرورگرهای وب معین کرده است. صفحات وب حاوی اطلاعات متنوعی میباشند از متن و تصویر گرفته تا فایل های تصویری. صفحات وب از نظر محتوا هم به چندین دسته تقسیم می شوند، مثلا بعضی صفحات وب مطالب کلی در باره هر موضوعی دارند و یا خاص منظوره اند به عنوان مثال صفحه ای درباره آخرین اخبار و اطلاعات ورزشی حاوی اطلاعات است و صفحه دیگری شامل جدیدترین اطلاعات سخت افزار کامپیوتر.

---

internet<sup>۱</sup>

web page<sup>۴</sup>

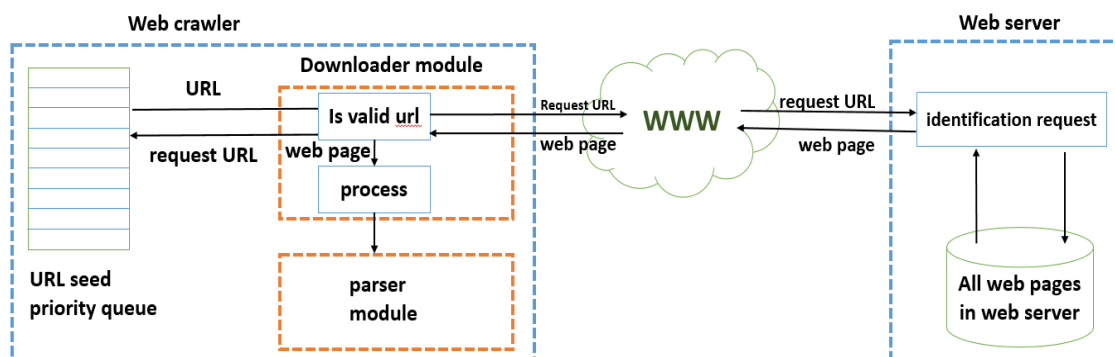
site<sup>۲</sup>

web browser<sup>۵</sup>

web crawler<sup>۳</sup>

world wide web (w3c)<sup>۶</sup>

یکی از وظایف موتورهای جستجو شناسایی موضوعیت کلی صفحه وب می‌باشد تا بهتر بتواند صفحه وب را مورد بررسی قرار دهد. به عنوان مثال اگر یک صفحه وب جدید ترین اطلاعات را درباره علوم ستاره شناسی نمایش می‌دهد پس شانس یافتن اطلاعات درباره جدیدترین روش هایی درمان سرطان در این صفحه وب کم و یافتن اطلاعات درباره یافتن جدیدترین کهکشان های کشف شده زیاد می‌باشد. در شکل ۱-۱ می‌توانید یک ساختار اولیه و کلی از موتور جستجو را ببینید.



شکل ۱-۱ نمای کلی یک موتور جستجوی وب

### HTML<sup>۱</sup> ۱-۱-۳

زبان نشانه گذاری ابر متن یا به اختصار **اچ‌دی‌ام‌ال** در کنار سی‌اس‌اس<sup>۲</sup> هسته فناوری ساخت صفحه‌های وب هستند. اچ‌دی‌ام‌ال زبان توصیف ساختار صفحه‌های وب است. زبانی است برای نشانه‌گذاری ابرمتن که برای تدوین قالب و طراحی صفحه‌های وب به کار می‌برند. دستورالعمل‌های این زبان، برچسب<sup>۳</sup> نام دارند که محتوای یک صفحه وب، با آن‌ها، نشانه‌گذاری شده و بدین ترتیب، نحوه نمایش آن صفحه برای مرورگرهای وب، توصیف می‌شود. آخرین نسخه این زبان در حال حاضر نسخه ۵ می‌باشد. در واقع زبانی که

<sup>۱</sup>hypertext markup language

<sup>۲</sup>cascading style sheets

<sup>۳</sup>tag

موتورهای جستجو آنرا می‌فهمند و از طریق آن به جمع‌آوری اطلاعات می‌پردازند همین زبان HTML می‌باشد. هر یک از برچسب‌های HTML، معنا و مفهوم و کاربرد خاصی دارند و تأثیر مشخصی بر محتوای صفحه می‌گذارند. به عنوان مثال برچسب‌های برای تغییرظاهری متن، نظیر درشت و ضخیم کردن یک کلمه یا برقراری پیوند به صفحات دیگر در HTML تعریف شده‌اند. و یا برچسب‌های برای نمایش فایل‌های صوتی تصویری مانند کلیپ‌های ویدیویی و حتی بعضی برچسب‌ها جهت تعریف متن و توضیحات در HTML تعریف شده‌اند.

یک سند HTML، یک پرونده مبتنی بر متن<sup>۱</sup> است که معمولاً با پسوند .htm یا .html نام‌گذاری شده و محتویات آن از برچسب‌های اچ‌تی‌ام‌ال تشکیل می‌شود. مرورگرهای وب، که قادر به درک و تفسیر برچسب‌های HTML هستند، تک‌تک آن‌ها را از داخل سند HTML خوانده و سپس محتوای آن صفحه را نمایان‌سازی<sup>۲</sup> می‌کنند.

HTML زبان برنامه‌نویسی نیست، بلکه زبانی برای نشانه‌گذاری ابرمتن است و اساساً برای ساختارمند کردن اطلاعات و جدایش اجزای منطقی یک نوشتار نظیر عناوین، تصاویر، فهرست‌ها، بندها و جدول‌ها به کار می‌رود. از سوی دیگر، HTML را نباید به عنوان زبانی برای صفحه‌آرایی یا نقاشی صفحات وب به کار برد. این وظیفه اکنون بر دوش فناوری‌های دیگری همچون CSS است. در موتورهای جستجو صفحات HTML توسط یک بخش مشخص به نام تجزیه‌کننده وب<sup>۳</sup> مورد بررسی قرار گرفته و اطلاعات و مطالب مهم از این صفحات استخراج می‌گردد. در فصل‌های بعد به اختصار درباره هر کدام از این بخش‌ها صحبت خواهد شد و همین‌طور به معرفی تگ‌ها و معرفی آنها می‌پردازیم.

---

text-based<sup>۱</sup>

render<sup>۲</sup>

webparser<sup>۳</sup>

## ۱-۲ اهداف پروژه

در این قسمت اهداف کلی این پروژه به طور مختصر عنوان می‌شود.

طراحی یک موتور جستجوی تحت وب :

۱- جستجو در بستر وب و جمع آوری اطلاعات از روی متن و تصاویر مرتبط با پرس‌وجو :

در این قسمت هدف دریافت صفحات وب و استخراج اطلاعات از داخل آن می‌باشد. و سپس

طبقه بندی اطلاعات و سرانجام نمایش نتایج.

۲- دادن امتیاز به صفحات بررسی شده جهت جستجوهای بعدی :

در اینجا به صفحات بازدید شده بر حسب ارزش و اطلاعاتی که به سوال مورد جستجوی ما

داده اند امتیازی داده می‌شود تا در جستجوهای بعدی درباره همین سوال یا سوالات و پرس‌وجوهای مشابه

(سوال های مشابه در یک طبقه) با سرعت و اطمینان بیشتری جستجو انجام شود.

۳- لیست اولویت :

این لیست حاوی پیوندهای<sup>۲</sup> مورد جستجو می‌باشد که بر حسب اولویت به موتور جستجو

جهت بررسی و دانلود داده می‌شوند.

۴- مبدل مساله<sup>۳</sup> :

استفاده از یک مبدل مساله جهت تولید سوالات مشابه با سوال پرسیده شده از موتور جستجو

و شاخص گذاری همه این سوالات در یک پایگاه داده جهت ساده سازی پرس‌وجوی بعدی.

---

<sup>۱</sup>webParser

<sup>۲</sup>link

<sup>۳</sup>problem generator

#### ۵- مقایسه نتایج جستجو :

مقایسه نتایج جستجو با نتایج و ویژگی های استخراج شده موجود در پایگاه داده و به روز رسانی ویژگی های مربوط به سوال پرسیده شده از موتور جستجو. بررسی پیوند های جدید بدست آمده از جستجو در صفحه وب جاری که بعضی از این پیوند ها مرتبط با موضوع مورد جستجوی ما هستند

و بعضی دیگر ارجاع به صفحات تبلیغات و فاقد محتوا هستند. تشخیص اهمیت و دادن اولویت به این صفحات از وظایف موتور جستجو است.

#### ۶- استخراج ویژگی :

کسب اطلاعات از داخل تصاویر و متون موجود در صفحه مورد بررسی (بعد از دانلود صفحه وب از لیست صفحات مورد بررسی، باید ویژگی های مورد نیاز را از داخل این صفحه استخراج کرده و همینطور پیوند های متصل به صفحات دیگر موجود در صفحه جاری را با توجه به اولویت به صف اولویت موتور جستجو اضافه کنیم).

- بررسی چالش ها، بررسی راهکارهای ارایه شده و معرفی نوآوری های پروژه :

در فصل سوم این پایان نامه به معرفی لیست چالش های یک موتور جستجوی تحت وب پرداخته سپس راهکار های ارایه شده تا کنون را مورد بررسی قرار می دهیم و در انتها به معرفی نوآوری های انجام گرفته در پروژه می پردازیم. در فصل چهارم هم به بررسی اجمالی از مراحل مختلف یک موتور جستجو می پردازیم.



## فصل دوم:

### معرفی موتور جستجوگر وب

## ۲-۱ معرفی

یک موتور جستجو در حقیقت نرم افزاری جهت بررسی و کسب اطلاعات درباره یک موضوع یا پرس و جوی خاص در صفحات وب، جمع آوری اطلاعات و نمایش نتایج می باشد. و از ابزار و قسمت های زیر تشکیل شده است.

- صف اولویت<sup>۱</sup>
- خزنده وب<sup>۲</sup>
- پایگاه داده : جهت ذخیره اطلاعات از قبیل صفحات وب و شاخص گذاری پرسش های انجام شده از موتور جستجو.

به طور خلاصه یک موتور جستجو از طریق یک نرم افزار به نام خزنده وب و با استفاده از یک صف شامل لیست آدرس صفحات وب به جستجو در صفحات وب پرداخته و نتایج بدست آمده را در یک پایگاه داده ذخیره می کند تا در جستجو های بعدی مطالب بدست آمده را در اختیار جستجو کننده قرار دهد.

## ۲-۲ خزنده وب چیست؟

خزنده وب به برنامه کامپیوتری گفته می شود که وب جهان گستر را به صورت مرتب و سلسله مراتبی بر حسب یکسری اولویت خاص بازدید کرده و اطلاعات آن را مورد پردازش قرار می دهد. از خزنده وب با عناوین دیگری مانند web robots و spiders ant نیز یاد می شود. البته نباید نامی ماننده خزنده وب باعث نگرانی از این موضوع شود که با یک ویروس کامپیوتری طرف هستیم.

---

<sup>۱</sup> priority queue

<sup>۲</sup> web crawler

## ۲-۳ کاربردهای خزنده وب

### ۲-۳-۱ موتورهای جستجو

در موتورهای جستجو، خزنده وب نقش بررسی کننده و کسب کننده اطلاعات را بازی می کند. خزنده وب سایت ها را به فرم دوره ای و بر حسب یک لیست اولویت بازدید کرده و نتایج جستجو را به موتور جستجو می فرستد. که بعد از پردازش و کسب اطلاعات و طبقه بندی داده ها و همینطور شاخص گذاری، نتایج را به کاربران نمایش می دهند.

### ۲-۳-۲ مدیریت فنی وب سایت

در یک وب سایت با هزاران صفحه و مطلب، و صدها کاربر که هر روز میلیون ها اطلاعات جدید را در سایت درج می کنند . مدیریت داده ها توسط نیروی انسانی کاری بسیار مشکل و حتی غیر ممکن به نظر می رسد. در اینجا هم وظیفه خزنده وب استخراج ویژگی ها از صفحات جهت سهولت مدیریت سایت می باشد. از طرف دیگر مدیریت فنی وب سایت شامل یافتن لینک های شکسته<sup>۱</sup>، اعتبار سنجی<sup>۲</sup> کدهای HTML، فایل های CSS و ... می باشد. در اینجا منظور از لینک شکسته در صفحه وب اشاره به لینک های است که در واقع بدون مقصد هستند یعنی وقتی که شما بر روی یک عبارت لینک شده کلیک می کنید، صفحه ی مربوط به آنها بارگذاری نمی شود و امکان دارد با خطای ۴۰۴ یا **Not Found** مواجه شوید. به این دلیل که یا صفحه به طور کامل از سایت حذف شده یا مشکلی برای صفحه رخ داده است و یا بحث مربوط به امنیت صفحه که مثلا توسط یک کاربر خاص قابل دیدن نیست و برای مشاهده اطلاعات داخلی این صفحه وب باید به اهراز هویت پرداخته شود.

---

<sup>۱</sup>broken links

<sup>۲</sup>validation

در حقیقت یکی از وظایف مهم خزنده وب تشخیص این لینک های شکسته است تا در بازدید های بعدی (به فرم دوره ای موتور جستجو هر چند وقت یکبار صفحات بررسی شده را دوباره دنبال می کند ) در لیست لینک های بررسی متوجه این لینک های معیوب شده و زمان و وقت خود را برای بررسی این لینک ها به هدر ندهد.

### ۳-۳-۲ جمع آوری اطلاعات خاص

کاربرد دیگر خزنده های وب جمع آوری اطلاعات خاصی مانند آدرس های ایمیل است. معمولا هدف از اینکار جلوگیری از ارسال هرزنامه می باشد. و یا حتی جمع آوری تصاویر از صفحات وب که موضوع این پروژه هم در این مورد می باشد.

### ۴-۲ نحوه کارکرد یک خزنده وب چیست؟

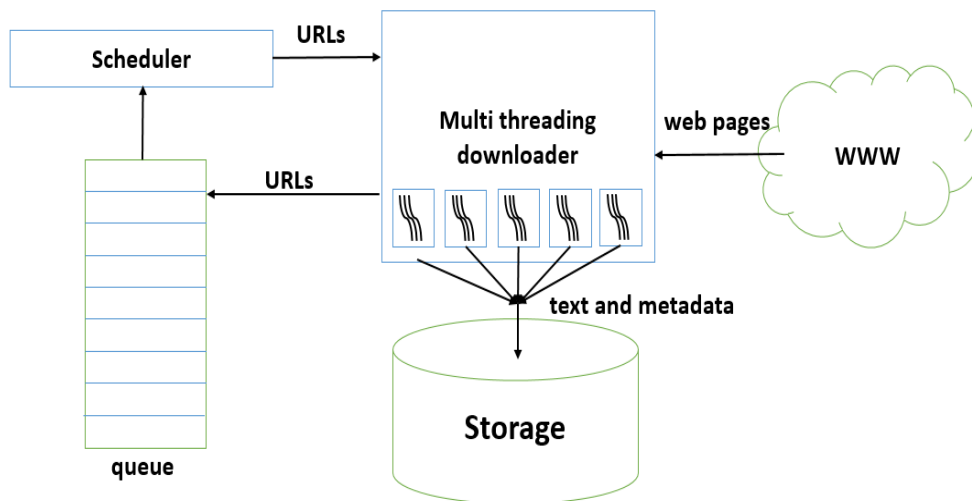
به صورت عمومی نحوه کار خزنده وب به این صورت است که ابتدا لیستی از URL<sup>۱</sup> ها (آدرس های وب) که به عنوان seed شناخته می شوند را برای بازدید پردازش می کند. البته این آدرس ها بر حسب یکسری امتیازها اولویت بندی می شوند. در مرحله بعدی صفحه مورد نظر دانلود شده و در مکانی مشخص ذخیره می شود. سپس باید اطلاعات داخل صفحه پردازش شوند هنگام پردازش آدرس ها، لیست لینک ها و آدرس های موجود در صفحات مورد جستجو را گردآوری کرده و به لیست ابتدایی اضافه می کنند. که اینجا هم صفحات اضافه شده به لیست باید بر حسب اولویت در صف قرار بگیرند بقیه اطلاعات را نیز با توجه به نیاز و هدف خود ذخیره و پردازش می کند.

---

<sup>۱</sup> uniform resource locator

## ۲-۴-۱ معماری

خزنده وب بخش اصلی و مرکزی هر موتور جستجو را تشکیل می‌دهد. به همین جهت الگوریتم و معماری آن‌ها به شدت مخفی نگه داشته می‌شود. با این وجود معماری سطح بالای آن در شکل ۱-۲ مشخص می‌شود.



شکل ۱-۲ معماری سطح بالای خزنده وب [1]

## ۲-۵ فرایندهای پیش پردازش یک خزنده وب

در زیر لیستی از فرایندها و پیش پردازش‌های یک خزنده وب را مشاهده می‌کنید.

### ۲-۵-۱ عادی کردن آدرس<sup>۲</sup>

منظور از عادی کردن آدرس، یکی کردن آدرس‌هایی می‌باشد که دارای خروجی یکسانی هستند. هدف از این کار جلوگیری از جمع آوری اطلاعات یکسان از چندین URL است. همانطور که قبلاً هم گفته شده ما

---

<sup>۱</sup> high level architecture

<sup>۲</sup> url normalize

لیست آدرس های صفحات وب را در یک صف اولویت قرار دادیم و اگر مدام آدرس های با مقادیر متفاوت ولی صفحه وب یکسان را در صف اضافه کنیم چیزی جز افزایش بارکاری برای موتور جستجو ندارد. بهتر است این آدرس ها یکی شوند. آدرس در وب برای دستیابی راحت تر کاربران به یک سایت یا محتوای خاص می باشد. هر صفحه وب دارای یک آدرس منحصر بفرد می باشد به عنوان مثال برای استفاده از سایت ویکی پدیا ما از آدرس منحصر بفرد <https://wikipedia.org> استفاده می کنیم. چند مرحله از عادی سازی یک آدرس در زیر نمایش داده شده است.

- تبدیل آدرس به حروف کوچک : هدف تبدیل تمام کاراکترهای بزرگ به کاراکترهای کوچک می باشد. به دلیل جلوگیری از تکرار آدرس به عنوان مثال تمام آدرس های زیر به یک صفحه وب اشاره می کنند

- <http://www.gamefa.com>
- <HTTP://WWW.GAMEFA.COM>
- <http://www.GAMEFA.com>

در زیر یک مثال از این فرایند را می بینید.

<HTTP://www.Example.com/> → <http://www.example.com/>

- افزودن / به انتهای آدرس در صورت نیاز : در زیر یک مثال از این فرایند را می بینید.

<http://www.example.com> → <http://www.example.com/>

- حذف آدرس ایندکس دایرکتوری : آدرس ایندکس فقط برای مدیریت صفحه وب می باشد و نتیجه آدرس بدون ایندکس با ایندکس یکی می باشد. پس بهتر است برای ایجاد آدرس های تکراری بالا ایندکس ها حذف گردد.

<http://www.example.com/default.asp> → <http://www.example.com/>

- بزرگ کردن حروف encode شده : بهتر است حروفِ انکد شده یا همان حروف بعد از علامت % حذف شوند.

<http://www.example.com/a%c2%b1b> → <http://www.example.com/a%C2%B1b>

حذف بخش زاید : بخش زاید بخشی است که معمولاً شامل بخش های فیلتری و محدود کننده می باشد. مثلاً در یک صفحه شامل ۱۰۰ تصویر اگر در انتهای آدرس عبارت `type=book` ? باشد شاید منظور این است که فقط تمام تصاویر با نوع کتاب را نمایش بدهد که با حذف این مورد ما به جای لیست فیلتر شده لیست کامل را ببینیم. و یا بخش های مانند اندیس که قسمتی از صفحه وب است که با علامت # شروع می شود و معمولاً کار خاصی انجام نمی دهد به عنوان مثال برای رفتن به بالای صفحه وب.

<http://www.example.com/bar.html#top> → <http://www.example.com/bar.html>

- حذف و تبدیل آی پی<sup>۱</sup> به دامنه : هر آدرس در بستر وب به یک آدرس آی پی (IP) تبدیل می شود این آدرس از یکسری عدد تشکیل شده است و کاملاً منحصر بفرد است. در حقیقت آدرس URL و آی پی هر دو یک آدرس را نمایش می دهند ولی آدرس URL عمومیت راحت تر و ساده تری دارد. پس بهتر است آدرس URL را بدست بیاوریم.

<http://208.77.188.166/> → <http://www.example.com/>

- اعمال محدودیت بر روی پروتکل ها مانند تبدیل https به http : دلیل این کار این است که معمولاً سایت ها در بستر اینترنت جهت جلوگیری از خرابکاری و هک<sup>۲</sup> از پروتکل های امنیتی مختلفی در

---

<sup>۱</sup> internet protocol (IP)

<sup>۲</sup> hack

سایت استفاده می کنند. مثلاً آدرس سایتی که با HTTPS شروع می شود غالباً از نظر امنیتی بسیار سخت گیر بوده و کار خزنده وب را برای جستجو در سایت دشوار می کند ولی آدرسی که با HTTP شروع می شود غالباً محدودیت های کمتری به خزنده وب اعمال می کند و مجوز جستجوی بیشتری را به خزنده وب می دهند.

`https://www.example.com/ → http://www.example.com/`

- حذف پورت پیش فرض (پورت ۸۰ به صورت پیش فرض برای http می باشد) : پورت در آدرس وب نمایشدهنده یک آدرس منحصر بفرد در سروری می باشد که سایت در آن قرار دارد (ممکن است خدمات دهنده تحت وب<sup>۱</sup> چند سایت را در خود داشته باشد که هر کدام پورت مخصوص به خود را دارند) ولی نیازی به وارد کردن آن در هنگام ورود به صفحه وب نیست پس جهت خلاصه سازی آدرس بهتر است مقدار پورت آدرس حذف گردد. از طرف دیگر ممکن است خدمات دهنده آدرس پورت را بنا به دلایلی عوض کند مثلاً از پورت ۱۲۳ به ۲۴۵ تغییر دهد پس خزنده وب دچار گمراهی می شود.

`http://www.example.com:80/bar.html → http://www.example.com/bar.html`

`http://www.example.com:80/bar.html → http://www.example.com/bar.html`

- حذف www از اول دامنه : این کلمه مشخص کننده دارا بودن لایسنس وب گستر جهانی سایت می باشد و بهتر است از آدرس حذف گردد.

`http://www.example.com/ → http://example.com/`

---

<sup>۱</sup>web server



- حذف متغیرهای اختیاری از رشته پرس‌وجو<sup>۱</sup>: در آدرس منظور از رشته پرس‌وجو یک پرس‌وجوی مشخص است که کار خاصی را روی صفحه وب انجام داده و محتوای سایت را تغییر می‌دهد بعضی از عناصر این پرس‌وجو تغییر خاصی در نتایج انجام نمی‌دهند مثلاً وجود پارامتر `&sort=ascending` یعنی نتایج را به ترتیب صعودی مرتب کنید که تاثیری در نتایج نمی‌دهد پس بهتر است از آدرس حذف شود. و یا مثلاً در صورتی که هیچ عبارت پرس‌وجوی در آدرس نیست بهتر است علامت ؟ از آدرس حذف گردد.

`http://www.example.com/display?id=&sort=ascending` → `http://www.example.com/display`

## ۲-۵-۲ بررسی آدرس صفحات وب

- معرفی غالب دامنه<sup>۲</sup> و صفحات وب: دامنه یک آدرس منحصر بفرد می‌باشد که نشان دهنده یک سایت خاص است و هر وب سایت دارای یک دامنه می‌باشد. این دامنه در سرویس خدمات دهنده قرار دارد غالب یک دامنه را در زیر می‌بینید. در این جا قسمت‌های که داخل کروشه ([ ]) قرار دارند می‌توانند در یک دامنه باشند و یا وجود نداشته باشد.

`protocol://[subDomains.]domainName.netType[:portNumber]`

حال به معرفی قسمت‌های مختلف یک دامنه می‌پردازیم.

- پروتکل اینترنت<sup>۳</sup> (protocol): پروتکل اینترنتی دربردارنده یک قرار دادی است جهت مبادله اطلاعات در شبکه‌های اینترنتی و هر پروتکل مشخص کننده نوع اهمیت صفحه وب می‌باشد به عنوان مثال پروتکل http لایه امنیتی پایین‌تری نسبت به پروتکل https دارد.

<sup>۱</sup> query string

<sup>۲</sup> domain

<sup>۳</sup> internet protocol

- نام دامنه (domainName) : نام دامنه در حقیقت اسم وب سایت می باشد و به طور کامل منحصر بفرد است.
- زیر دامنه (subDomain) : ممکن است یک دامنه شامل چندین زیر دامنه هم باشد که در حقیقت وظیفه قسمت کردن و گروه بندی اطلاعات وبسایت را بر عهده دارند. به عنوان مثال اگر یک دامنه با نام `http://www.gamefa.com` داشته باشیم یکی از زیر دامنه های آن مثلا `http://www.cinema.gamefa.com` می باشد.
- نوع نت (netType) : این قسمت مشخص کننده استاندارد نت می باشد مثلا نوع نت `.com` نشان دهنده یک سایت عمومی و `.ir` نشان دهنده وجود دامنه در محدوده ایران می باشد.
- شماره پورت (portNumber) : هر خدمات دهنده وب شامل یک یا چند دامنه است که هر کدام را با یک شماره پورت منحصر بفرد مشخص می کند.
- حال به معرفی ساختار یک آدرس وب می پردازیم. در این جا قسمت های که داخل کروشه ([ ]) قرار دارند می توانند در یک صفحه وب باشند و یا نه.

`domain/path/[?query][#hash]`

- دامنه (domain) : این قسمت همان دامنه است که در قسمت قبل به طور کامل معرفی گردید.
- مسیر (path) : این قسمت مشخص کننده مکان مشخصی در دامنه است و شامل یک مسیر مشخص در دامنه است به عنوان مثال `https://a-z-animals.com/animals/dolphin/` مسیر مشخصی را در دامنه معرفی می کند که ما را به بخش حیوانات (`/animal`) و در نهایت صفحه وب معرفی دلفین (`/dolphin`) هدایت می کند.
- پرس و جو (query) : این قسمت جهت فیلتر کردن یک صفحه خاص مورد استفاده قرار می گیرد به عنوان مثال در یک صفحه شامل ۱۰۰ تصویر درباره کتاب های مختلف که با استفاده از پرس و جوی

مثل `?type=computer` فقط کتاب‌های با عنوان کامپیوتر را مشاهده می‌کنیم. بهتر است جهت جلوگیری از فیلتر و حذف قسمت‌های از صفحه وب این قسمت به طور کامل از آدرس حذف گردد.

- اندیس (hash) : این قسمت از بخش‌های اختیاری یک آدرس وب می‌باشد. این بخش فقط جهت ارجاع به یک قسمت مشخص از صفحه وب می‌باشد و کاربرد خاص دیگری ندارد به عنوان مثال وجود عبارت `#film` ما را به بخش معرفی فیلم در صفحه وب می‌برد. بهتر است این قسمت هم از صفحات وب حذف شود.

## ۶-۲ قسمت‌های مهم یک خزنده وب :

### ۱-۶-۲ شاخص گذاری<sup>۱</sup>

در این قسمت هدف نگاشت نتایج به دست آمده از صفحات وب به پرس‌وجوهای پرسیده شده از موتور جستجو می‌باشد. این کار به منظور سهولت در دستیابی به اطلاعات بدست آمده از صفحات وب می‌باشد. به عنوان مثال موتور جستجوی گوگل اطلاعات بیش از هزاران میلیارد سایت را در یکصد میلیون گیگابایت شاخص گذاری کرده است. البته بعضی از پرس‌وجوها دارای شباهت‌های از نظر معنای هستند که در یک گروه اندیس می‌شوند مانده دو پرس‌جوی "اطلاعاتی درباره انواع میزها" و "معرفی میزهای مختلف" که در حقیقت به یک موضوع مشابه اشاره دارند.

---

<sup>۱</sup> indexing

## ۲-۶-۲ پرس و جو<sup>۱</sup>

مفهوم پرس و جو در خزنده وب و بررسی موضوعات مشابه با شاخص ها (جهت سرعت در بررسی نتایج) است. در حقیقت پرس و جو یک نوع طبقه بندی نتایج صفحات وب است. زمانی که کاربر سوالی از موتور جستجو می پرسد. این جستجو در نهایت به یکی از طبقه بندی های موتور جستجو (یا چند تا با درجه اولویت مختلف) وصل شده و موتور جستجو بر حسب اولویت اطلاعات طبقه مورد بحث نتایج را نمایش می دهد. از طرف دیگر یکی از کاربردهای پرس و جو ایجاد یکسری سوال و درخواست توسط بخشی در موتور جستجو به نام مبدل مساله<sup>۲</sup> است. به این صورت که وقتی کاربر سوالی از موتور جستجو می پرسد مبدل مساله درخواست های مشابه با آن ایجاد می کند. این کار باعث گسترش و پیشرفت بهتر موتور جستجو می شود. در این پروژه همانطور که در فصل های جلوتر خواهیم دید ما از ۸۰۰ پرس و جوی مختلف برای تشخیص اشیاء استفاده می کنیم.

## ۲-۶-۳ رتبه بند<sup>۳</sup>

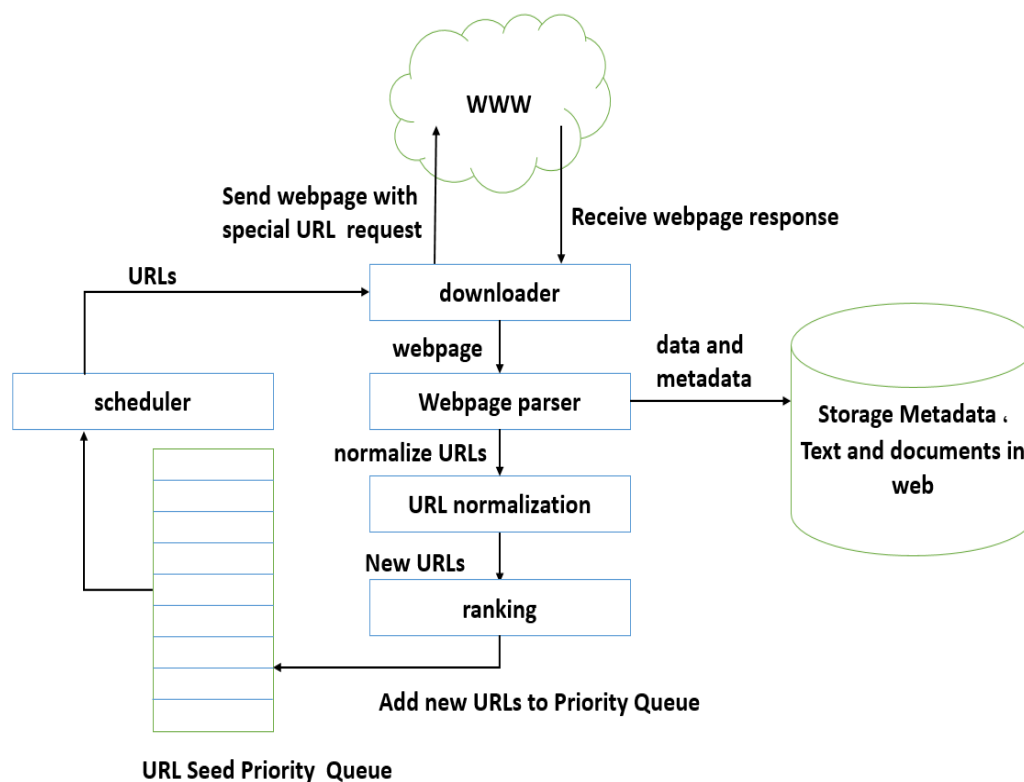
هدف از این قسمت دادن امتیاز به صفحات بررسی شده توسط موتور جستجو می باشد تا در بررسی های بعدی جهت به روز رسانی اطلاعات درباره یک موضوع یا موضوع های مشابه با آن بر حسب این امتیاز صفحات اولویت بندی شوند. آدرس صفحات وب بسیار زیاد می باشند و در بعضی مواقع هیچگونه شباهتی با موضوع مورد جستجو ندارند پس بهتر است به هر آدرس یک امتیاز خاص درباره یک موضوع خاص داده شود. در فصل های بعدی به معرفی راهکارهای برای دادن امتیاز به صفحات وب بر حسب یک روش پیشنهادی می پردازیم. در شکل ۲-۲ تصویری از معماری یک خزنده وب مشاهده می شود.

---

Query<sup>۱</sup>

problem generator<sup>۲</sup>

ranking<sup>۳</sup>



شکل ۲-۲ معماری کامل یک خزنده وب [1]

#### ۲-۶-۴ Frontier Crawler

این بخش شامل لیستی از آدرس صفحات وب معتبر<sup>۱</sup> می‌باشد. این لیست شامل مجموعه ای از آدرس صفحات وب با نام دانه<sup>۲</sup> می‌باشد. که یا توسط کاربران موتور جستجو معرفی شده اند و یا توسط خزنده وب در صفحات وب یافت شده‌اند. البته استخراج آدرس بر حسب اولویت می‌باشد. در پروژه ما ابتدا به معرفی یک لیستی از آدرس دامنه های مشخص پرداخته و خود موتور جستجو با استخراج آدرس های جدید از صفحات بررسی شده می‌پردازد.

<sup>1</sup> valid URL

<sup>2</sup> seed

## ۵-۶-۲ دریافت کننده<sup>۱</sup>

مهمترین کار این بخش همانطور که از اسمش مشخص است دریافت صفحه وب مشخصی از لیست Frontier می باشد (بر حسب اولویت). بنابراین این قسمت از خزنده یک درخواست HTTP برای صفحه مورد نظر به سروری که شامل این صفحه وب است فرستاده و نتیجه ارسال شده از سرور را که همان صفحه وب است را دریافت می کند اطلاعاتی که این بخش دریافت می کند همان صفحه وب با استاندارد HTML می باشد. و سپس به استخراج اطلاعاتی مانده متن، تصویر و آدرس صفحات دیگر داخل صفحه دریافتی می پردازیم.

## ۶-۶-۲ مخزن<sup>۲</sup>

این قسمت جهت ذخیره و مدیریت لیست بزرگی از داده ها استفاده می شود. البته مخزن تنها فقط صفحات وب استاندارد (html) را شامل می شوند به همراه تاریخچه ای از آنها. در اینجا منظور از تاریخچه اشاره به دوره های مختلفی است که موتور جستجو یک صفحه وب را دریافت می کند زیرا اطلاعات و ساختار صفحات وب دچار تغییرات دوره ای می شود که بسته به نوع صفحه وب کاملاً متغیر می باشد به عنوان مثال صفحه اصلی یک سایت

خبری تقریباً هر ۱۰ دقیقه یکبار تغییر می کند ولی یک صفحه آموزشی ممکن است هفته ها دچار تغییر نشود پس باید در دوره های مختلف به دریافت این صفحات وب پرداخت.

## ۷-۲ انواع موتورهای جستجوی وب :

در اینجا به معرفی انواع موتورهای جستجو می پردازیم. شکل ۲-۳ الگوریتم کامل یک خزنده وب را نمایش می دهد.

---

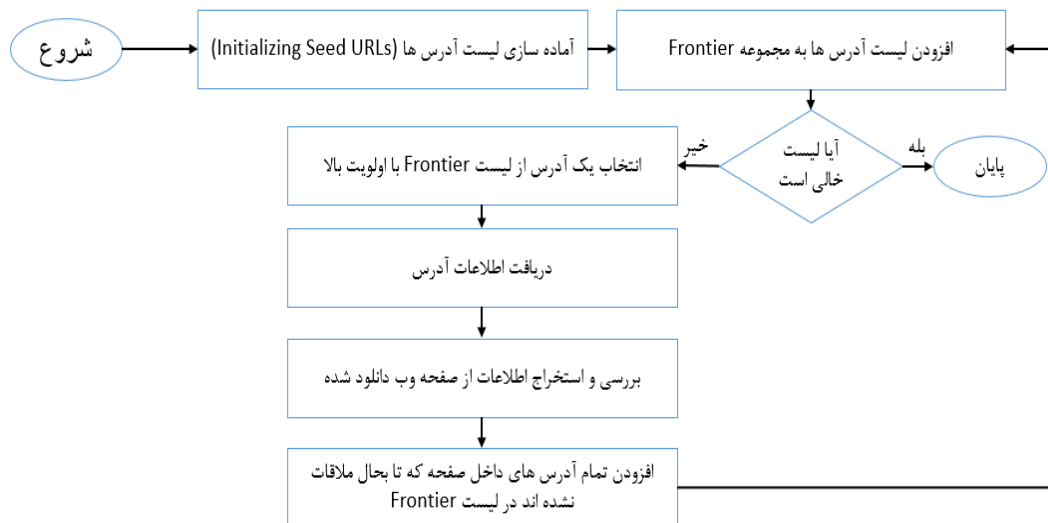
<sup>۱</sup> Downloader

<sup>۲</sup> repository

## ۱-۷-۲ خزنده وب موازی<sup>۱</sup>:

این نوع از خزنده وب شامل چندین پردازنده جهت بررسی صفحات وب به فرم موازی است (دریافت صفحات وب و ذخیره آنها به فرم محلی و استخراج لینک و اطلاعات از صفحات وب به فرم موازی توسط چند خزنده وب انجام

می‌گیرد). در این نوع موتورهای جستجو باید هریک از فرایندها با بقیه ارتباط داشته باشند جهت به روز رسانی موتور جستجو و یکی شدن اطلاعات هر کدام از فرایندها.



شکل ۲-۳ الگوریتم کامل مراحل یک خزنده وب [1]

• مزایا

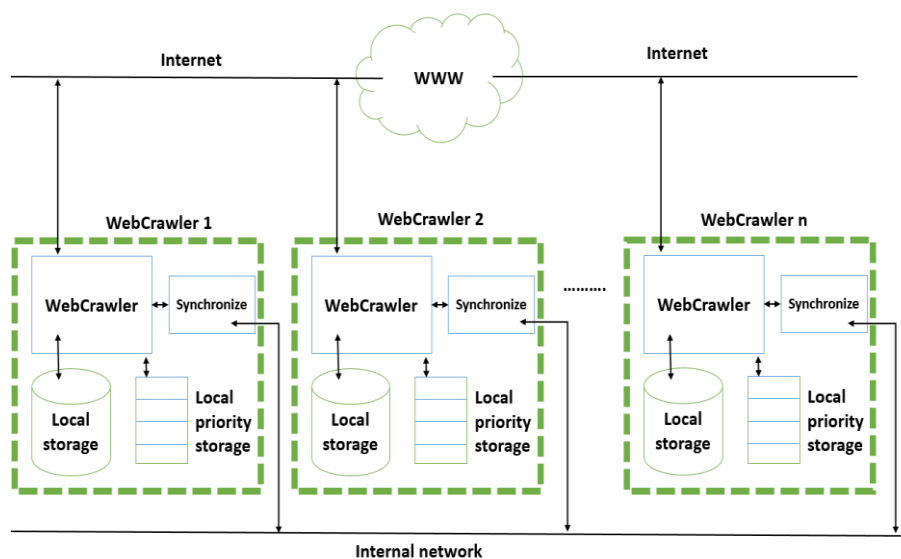
- پردازش موازی و سریع تر داده ها

<sup>۱</sup>Parallel

- مقیاس پذیری<sup>۱</sup>: تعداد صفحات وب به فرم نمایی در حال رشد است پس برای یک خزنده وب تنها غیر ممکن است بدست آوردن و بررسی تمام صفحات وب در یک زمان معقول پس به فرم همزمان از چند موتور جستجو استفاده می شود.
- قابلیت مقایسه نتایج جستجو و اشتراک گذاری
- معایب

- داده های تکراری: به علت عملکرد موازی ممکن است چند خزنده وب یک صفحه را به فرم جدا بررسی کنند و نتایج مشابه زیادی در موتور جستجو و مخزن آن ذخیره شود.

شکل ۲-۴ یک نمایش بلوک دیاگرامی از موتور جستجوی موازی را نمایش می دهد.



شکل ۲-۴ خزنده وب موازی [1]

۲-۷-۲ خزنده وب توزیع شده<sup>۲</sup>

<sup>۱</sup> Scalability

<sup>۲</sup> distributed



این خزنده وب صفحات وبی را که با هم رابط دارند دریافت کرده. این صفحات با موضوع مورد جستجوی خاص موتور جستجو هم ارتباط دارند. به فرم تخصصی این نوع خزنده از نمونه‌ها و کاوش در وب از طریق ارتباطات هدایت شده یاد می‌گیرد و شامل ۳ مولفه زیر است.

۱- کلاسه بند<sup>۱</sup>: این قسمت جهت طبقه‌بندی مستندات وب به چند گروه مختلف استفاده می‌شود

و اطلاعات به دست آمده از صفحه‌وب را در گروه‌های مختلف ادغام می‌کند.

۲- تقطیرکننده<sup>۲</sup>: به فرم خلاصه وظیفه این قسمت بررسی ارزش صفحه با موضوع مورد بررسی

می‌باشد. البته این ارزش روی صفحات و لینک‌های مرتبط با صفحه انتخابی هم تاثیر دارد، و

ممکن است ارزش این صفحات را تغییر دهد.

۳- خزنده: جهت بررسی صفحه از این قسمت استفاده می‌شود. و وظیفه آن دریافت صفحه وب

است.

- مزایا

- یادگیری هوشمند و گسترده تر داده ها: به خاطر طبقه بندی اطلاعات.

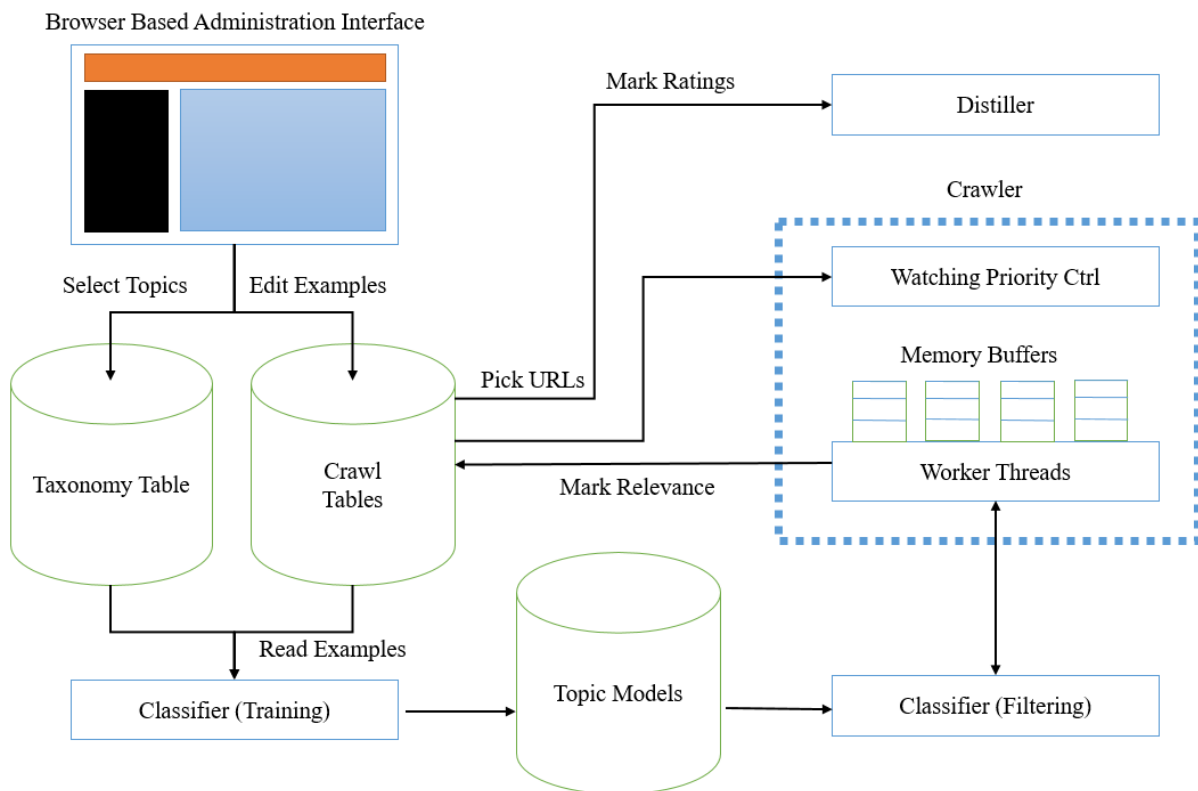
- سهولت دادن امتیاز به صفحات وب به خاطر استفاده از سیستم تقطیر کننده

- معایب

- عدم توانایی بررسی صفحات مخفی: منظور از صفحات مخفی صفحاتی است که جهت ورود به آنها

نیازمند تشخیص هویت می‌باشد.

شکل ۲-۵ یک نمایش بلوک دیاگرامی از موتور جستجوی توزیع شده را نمایش می‌دهد.



شکل ۲-۵ خزنده وب توزیع شده [1]

### ۳-۷-۲ خزنده وب افزایشی<sup>۱</sup>

در این نوع از خزنده وب صفحات وب با یک دوره زمانی مشخص بررسی می‌شوند. در این روش حتی امکان پاکسازی بعضی صفحات، جهت آزاد سازی فضا برای دریافت صفحات جدید است. در این نوع از خزنده وب اهداف زیر تحقق می‌یابد [1].

- تازه نگهداشتن اطلاعات محلی : با بررسی دوره‌ای صفحات امکان تازه نگهداشتن اطلاعات و کسب اطلاعات بهتر بیشتر می‌شود. به دلیل تغییر مداوم اطلاعات در بستر وب بسیار مهم است که موتور جستجو به سرعت از تغییرات صفحات وب آگاه شود.

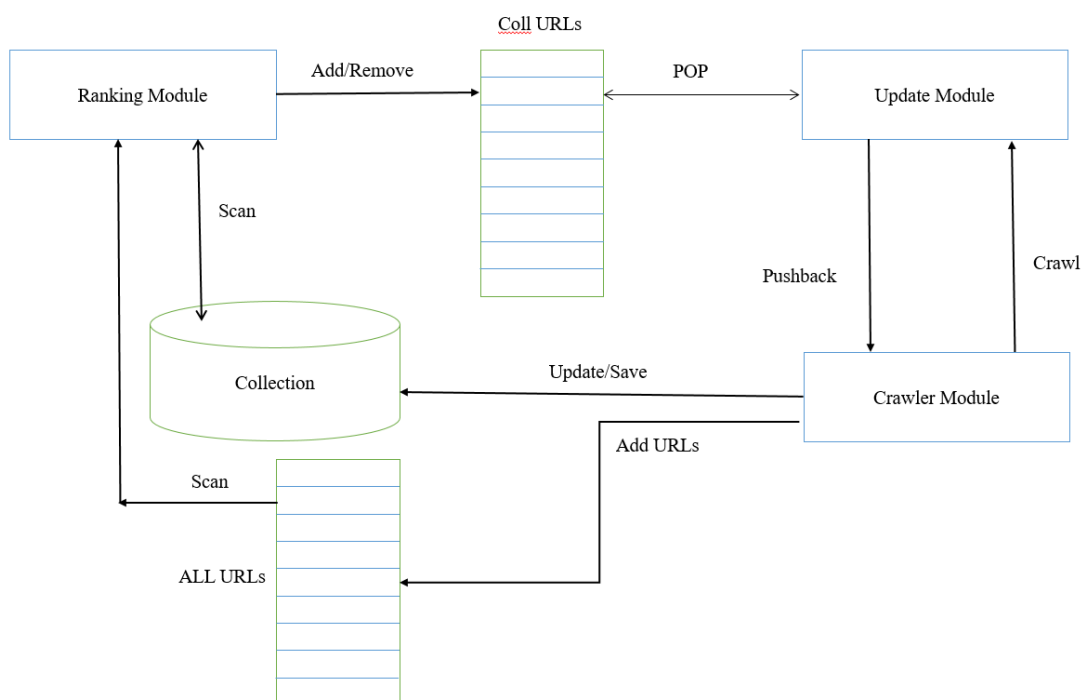
<sup>۱</sup>Incremental

- بهبود کیفیت مجموعه محلی : با بررسی دوره‌ای صفحات وب با صفحاتی که دارای اطلاعات با ارزش تری هستند آشنا می‌شویم. مثلاً اگر صفحه‌ای مدام ارزش اتصالش در بررسی‌ها پایین برود به دلایل مختلف مثل کاهش بازدید و یا شباهت کمتر با موضوعات پرس‌وجوی ما (در این پروژه ۸۰۰ گروه مختلف داریم) پس در دوره‌های بعدی ارزش کمتری در لیست اولویت می‌گیرد.

#### • معایب

- نتایج ممکن است خیلی با موضوع بررسی خزنده وب هم‌خوانی نداشته باشند
- بار شبکه : به علت بررسی مداوم و تکراری صفحات یک بار اضافه به شبکه تحمیل می‌شود.
- عدم توانایی بررسی صفحات مخفی

شکل ۲-۶ یک نمایش بلوک دیاگرامی از موتور جستجوی افزایشی را نمایش می‌دهد.



شکل ۲-۶ خزنده وب افزایشی [1]

## 4-7-2 خزنده وب مخفی

به فرم معمول موتور جستجو قادر به بررسی صفحاتی که نیازمند تشخیص هویت‌اند نیست از این جهت از خزنده های وب مخفی استفاده می‌شود [3] که با استفاده از یک جدول به نام انتساب دهنده ارزش<sup>۱</sup> که به اختصار LVS گفته می‌شود و ۴ مولفه زیر اقدام به تشخیص هویت و ورود به صفحات وب مخفی می‌کند.

- تجزیه و تحلیل فرم<sup>۲</sup>: بررسی اطلاعات داخل فرم ثبت نام درون وب سایت. مثلاً فرم حاوی یک المنت ورودی ایمیل است و یک فیلد تعیین جنسیت و ...

- انتساب دهنده ارزش<sup>۳</sup>: در این مرحله ما از جدول LVS که نوعی نگاشت دهنده از جدول به المنت های داخل فرم است استفاده کرده و فیلدهای داخل فرم را پر می‌کنیم. هر المنت داخل جدول حاوی لیستی از مقادیر متفاوتی است به عنوان مثال فیلد ایمیل شامل لیستی از ایمیل های معتبر می‌باشد که انتساب دهنده ارزش یکی از آنها را در فرم قرار می‌دهد.

- ارسال کننده فرم<sup>۴</sup>: این قسمت مسئول ارسال فرم پر شده به سرور وب سایت مشخص شده می‌باشد

- تجزیه و تحلیل پاسخ<sup>۵</sup>: این قسمت مسئول بررسی پاسخ دریافتی از سرور می‌باشد و رفع مشکلات احتمالی (مثلاً اگر ایمیل قبلاً در سرور ثبت شده دوباره از جدول LVS یک ایمیل جدید در فیلد مربوطه قرار دهد و دوباره فرم را به سرور ارسال کند).

---

<sup>۱</sup> label value set

<sup>۲</sup> form analysis

<sup>۳</sup> value assignment

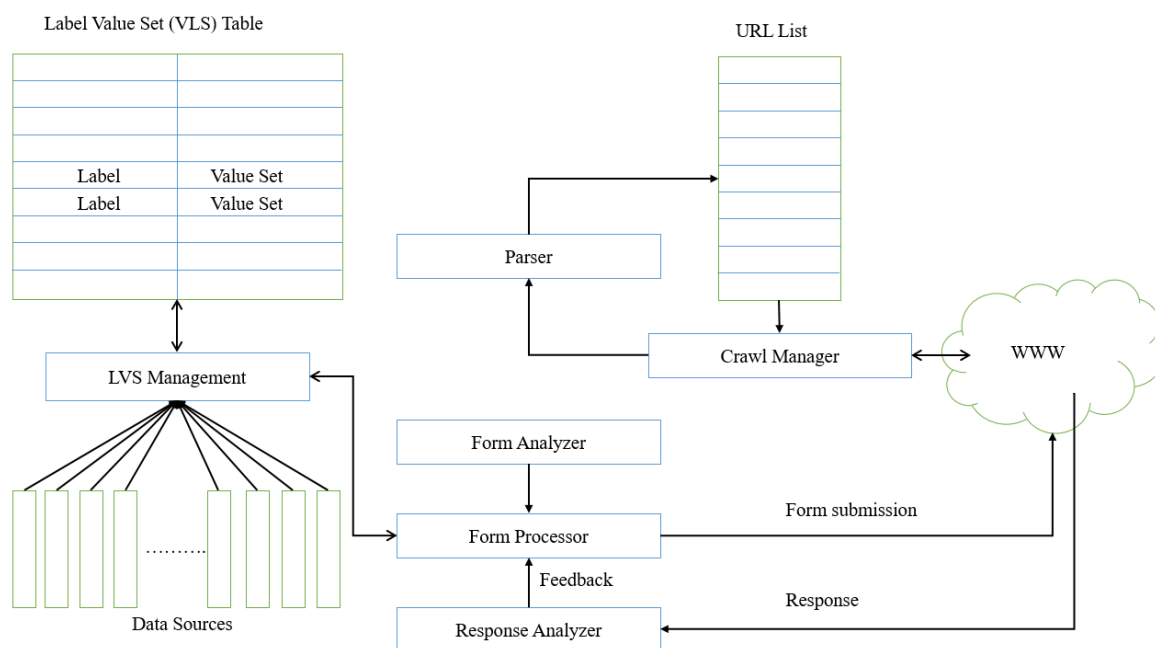
<sup>\*</sup> form submission

<sup>د</sup> response analysis

• معایب

- بار شبکه : به علت ارسال و بررسی مرتب فرم با جدول LVS نوعی بار اضافه در شبکه ایجاد می شود.

شکل ۷-۲ یک نمایش بلوک دیاگرامی از موتور جستجوی مخفی را نمایش می دهد.



شکل ۷-۲ خزنده وب مخفی [1]

## ۲-۸ معرفی چند موتور جستجوی مهم

بعضی موتورهای جستجو کاربرد خاص منظوره دارند به عنوان مثال فقط در زمینه تشخیص هرزنامه ها و یا تشخیص تصاویر در صفحات وب استفاده می شوند ولی بعضی از موتورهای جستجو کاربردهای عمومی تر دارند و در همه موارد استفاده می شوند. جدول ۱-۲ لیست از ۵ موتور جستجوی محبوب سال ۲۰۱۶ را نمایش می دهد البته بعضی موتورهای جستجو حالت کاربردی در علوم مختلف را هم ایفا می کنند به عنوان

مثال یافتن اطلاعات درباره یک موضوع خاص یا حتی جمع آوری داده برای یک موضوع مشخص مانده دریافت لیست تصاویر درباره صورت انسان جهت پردازش تصویر و یا ایجاد یک پایگاه داده اولیه برای یک فرایند یادگیری. در این پروژه هم با استفاده از استخراج متن و تصاویر به گروه بندی و تشخیص اشیاء می پردازیم.

جدول ۱-۲، ۵ موتور جستجوی محبوب وب

ردیف	نام	لوگوی موتور جستجو
۱	Google	
۲	Yahoo	
۳	Mahalo	
۴	webopedia	
۵	Bing	

## ۹-۲ معرفی فایل robots.txt

این فایل برای دادن اطلاعات اولیه در زمینه وب سایت مورد پردازش به موتور جستجو استفاده می گردد. به عنوان مثال با این فایل خدمات دهنده صفحه وب می تواند میزان دسترسی موتور جستجو به بعضی زیر شاخه ها را مشخص کند. به عنوان مثال دستورات زیر که داخل فایل robots.txt قرار دارد از دسترسی موتور جستجو

به مسیر <sup>۱</sup> /tmp/ جلوگیری می‌کند:

User-agent: \*

Disallow: /tmp/

فایل robots.txt یک استاندارد می‌باشد. به همین جهت موتور جستجو (معمولا موتورهای جستجو مخرب) می‌تواند آن را نادیده بگیرد. ولی این فایل بیشتر جهت راهنمایی یک موتور جستجو برای دسترسی سریع‌تر و با اطمینان بالاتری به اطلاعات کلی صفحات وب یک دامنه خاص است.

---

Path<sup>۱</sup>





فصل سوم:

معرفی چالش ها،

راه حل ها و نوآوری های

موتور جستجو

## ۳-۱ مقدمه

در این فصل ابتدا به معرفی ساختار داخلی صفحات وب پرداخته و در ادامه به بررسی چند مشکل و راه کارهای برای آنها در موتور جستجو می پردازیم.

## ۳-۲ دریافت و بررسی صفحات وب

### ۳-۲-۱ مقدمه

همانطور که گفته شد صفحات وب در کنسرسیوم جهانی وب قالب استاندارد HTML می باشد. اما غالب یک HTML به چه صورت می باشد. در فصل اول و به طور مختصر با مفهوم زبان HTML آشنا شدیم در اینجا بهتر است قبل از جمع آوری اطلاعات از یک صفحه وب با غالب صفحه وب و کلیاتی از آن آشنا شویم.

### ۳-۲-۲ معرفی غالب یک صفحه وب

یک صفحه ابر متن از یکسری المنت با نام تگ<sup>۱</sup> تشکیل می شود غالب کلی یک تگ را می توانید در شکل ۳-۱ ببینید. همانطور که در تصویر مشخص است هر تگ دارای یک نام منحصر به فرد است و یکسری ویژگی که هر کدام یک نام منحصر به فرد و مقدار دارند. هر تگ مخصوص به کار خاصی است مثل یکسری تگ برای درج متن، بعضی برای نمایش فایل های صوتی و تصویری و غیره استفاده می شوند. بعضی ویژگی ها نیز در تگ ها برای ظاهر و غالب تگ تعریف می شوند و بعضی ها مقادیر مشخصی به تگ می دهند به عنوان مثال ویژگی کلاس<sup>۲</sup> در هر تگ برای تعریف یک غالب و ظاهرو یا ویژگی value به مقدار تگ اشاره می کند. در جدول ۳-۱ لیستی از چند تگ مهم در صفحات وب معرفی شده است. شناسایی این تگ ها در صفحات

---

<sup>۱</sup> Tag

<sup>۲</sup> class

وب از اهمیت بالای برخوردار است به این دلیل که با شناسایی و کاربرد هر تگ در صفحات وب می‌شود از ماهیت و کاربردهای یک صفحه وب خاص

```
< tag-name  
  attribute1 = "value1"  
  attribute2 = "value2"  
  ⋮  
  attributen = "valuen"  
>  
  inside tag values  
< / tag-name>
```

شکل ۱-۳ غالب کلی یک تگ HTML

مطلع گردید و اهمیت صفحه درباره یک موضوع خاص مورد پرس‌وجو پررنگتر شود. مثلاً اگر در یک صفحه خاص تگ شکل زیادی موجود باشد می‌شود در بررسی تصاویر بیشتر به این صفحه مراجعه کرد و یا اگر صفحه ای دارای لینک پیوندی زیاد است به احتمال زیاد این صفحه یک صفحه خبری و یا حتی تبلیغاتی می‌باشد. البته باید به کمک تکنیک‌های هوشمند و روش‌هایی امتیاز دهی دستی ( رای کاربران در بازدید سایت بعد از درخواست یک سوال و انتخاب لینک نتیجه ) در گرفتن امتیاز صفحه در یک مورد خاص اهمیت داد. در شکل ۲-۳ لیستی از چند تگ پرکاربرد صفحات وب را می‌توان دید همانطور که در این شکل مشخص است تگ‌های قسمت (a) حاوی اطلاعات مفیدی نیستند و بیشتر به منظور تغییر ظاهر صفحه وب استفاده می‌شوند. از طرفی شانس مشاهده متن در تگ‌های قسمت (b) در شکل بالا است و همچنین احتمال مشاهده یک تصویر در تگ‌های قسمت (c) بالا است.

<br>      <p>      <img>  
 <script>   <span>   <figure>  
 <hr>      <strong>   <figcaption>  
 <style>   <ul>, <li>  
 a                      b                      c

شکل ۲-۳ چند تگ مهم و کاربرد

جدول ۱-۳ لیست تگ‌های پر کاربرد صفحات وب

| ردیف | نام تگ    | توضیحات   | نام چند ویژگی مهم   |
|------|-----------|---|---|
| ۱    | <img>     | شامل یک تصویر در صفحه وب و فاقد فرزندی می‌باشد  | src = آدرس تصویر<br>width = عرض تصویر<br>height = ارتفاع تصویر<br>alt = توضیحات تصویر |
| ۲    | < button> | تعریف یک دکمه قابل کلیک در صفحه وب  | name = نام کلیک   |
| ۳    | <figure>  | شبیه تگ img فقط دارای یک بخش توضیح هم است.  | figcaption = توضیحات تصویر<br>src = آدرس تصویر  |
| ۴    | <span>    | نمایش توضیحات و متن در صفحات وب   | value = مقدار   |
| ۵    | <menu>    | نمایش یک لیست انتخاب شوند در وب (غالباً در صفحات وب برای نمایش منو ها و اتصال به صفحات دیگر از این تگ استفاده می‌شود)   | menuitem = مقدار یک عنصر<br>value = مقدار   |
| ۶    | <p>       | نمایش یک پاراگراف در صفحات وب ( این تگ بیشتر برای نمایش توضیحات و پاراگراف ها در وب استفاده می‌شود)   |   |
| ۷    | <meta>    | این تگ جهت توضیحات صفحه وب بکار می‌رود و به عنوان مثال ویژگی keywords شامل لیستی از کلمات کلیدی صفحه وب فوق می‌باشد که برای راهنمایی موتور جستجو و یا فریب آن استفاده می‌شود. |   |

---

|   |     |   |                             |
|---|-----|---|-----------------------------|
| ۸ | <a> | این تگ جهت ارتباط با دیگر صفحات وب استفاده می‌شود و بیشتر در صفحاتی با غالب خبری و تبلیغاتی استفاده می‌شود. | آدرس صفحه وب = href<br>دیگر |
|---|-----|---|-----------------------------|

---

### ۳-۲-۳ دریافت کننده<sup>۱</sup> صفحه وب : در شکل ۳-۳ دیاگرامی از وظایف این بخش نمایش داده شده

است. حالا به معرفی قسمت‌های مختلف این قسمت از موتور جستجو می‌پردازیم.

۱- درخواست آدرس صفحه وب با بالاترین اولویت از صف seed و بررسی آن : منظور از بررسی آدرس

اطمینان از عدم مشکل یا هرزنامه<sup>۲</sup> بودن آدرس می‌باشد. این کار به صورت تجربی و بررسی ارزش

صفحه وب بدست می‌آید.

۲- ارسال درخواست صفحه وب انتخابی به خدمات دهنده<sup>۳</sup> آن و انتظار دریافت نتیجه از خدمات

دهنده

۳- بررسی دریافت موفق صفحه وب : آیا صفحه وب با موفقیت از خدمات دهنده آن دریافت گردیده

است (در صورت دریافت صحیح قسمت دریافت کننده پیغام ۲۰۰ را که نماینده جواب صحیح است

دریافت می‌کند). در جدول ۲-۳ لیست نتایج برگردانده شده از خدمات دهنده دارای آدرس صفحه

وب را نمایش می‌دهد. این نتایج در تعیین امتیاز یک صفحه وب کمک زیادی می‌کنند به عنوان

مثال با ایجاد یک خطا در دریافت صفحه وب می‌شود امتیاز آن را کاهش داد.

---

<sup>۱</sup> web page downloader

<sup>۲</sup> spam

<sup>۳</sup> web server

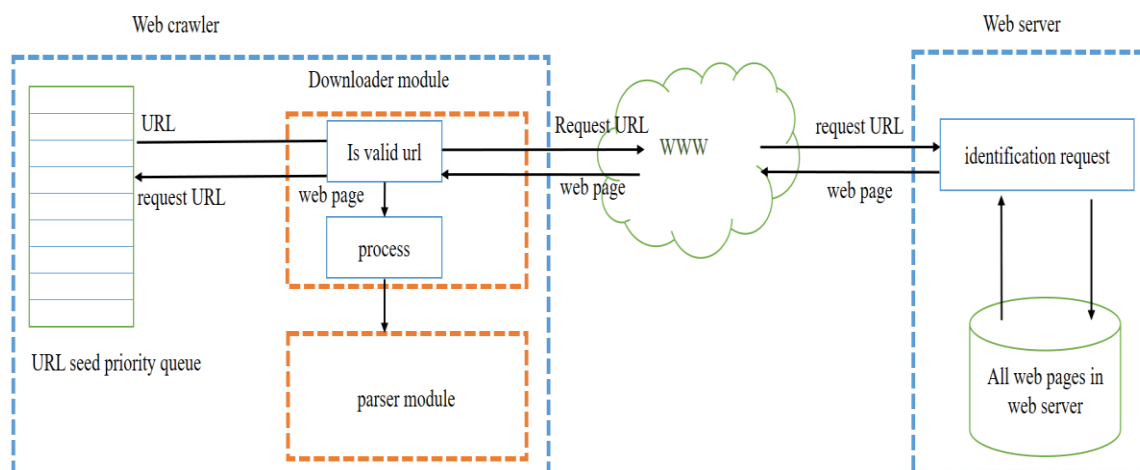
<sup>۴</sup> redirection

جدول ۲-۳ لیست نتایج درخواست یک صفحه وب

| ردیف | نام تگ | توضیحات  |
|------|--------|--|
| ۱    | ۲۰۰    | صفحه وب با موفقیت از خدمات دهنده دریافت شد.  |
| ۲    | ۴۰۴    | صفحه وب مورد نظر یافت نشد. (به دلایل مختلفی مانده حذف موقت صفحه وب به منظور به روز رسانی اطلاعات آن یا حذف دائمی آن و یا حتی تغییر آدرس صفحه وب <sup>۴</sup> ) |
| ۳    | ۵۰۱    | جهت مشاهده اطلاعات صفحه وب نیاز به اعمال اهراز هویت می باشد  |
| ۴    | ۴۰۶    | خدمات دهنده به فرم موقت قادر به ارسال پاسخ (صفحه وب) نیست به دلایلی مانده حجم بالای درخواست های ارسالی به خدمات دهنده  |

پیش پردازش روی صفحه وب دریافتی : قبل از ارسال صفحه وب دریافتی به قسمت تجزیه کننده باید یکسری پیش پردازش روی صفحه وب انجام بگیرد :

- بررسی عدم خرابی تگ ها : در html باید تمام تگ ها بسته شوند در غیر این صورت در هنگام تجزیه در مازول تجزیه کننده ممکن است مشکلی ایجاد گردد. منظور از بسته بودن تگ داشتن انتها برای آن است (</tagname>) به دلیل غالب درختی یک صفحه وب.
- بررسی صحیح غالب صفحه وب : صفحه وب دریافتی از نظر غالب و سلسله مراحل تگ ها درست است. به عنوان مثال تگ <body> باید بعد از تگ <head> در صفحه وب قرار داشته باشد زیرا تگ head یکسری اطلاعات مهم را قبل از بارگذاری تگ body در صفحه وب دریافت می کند.



شکل ۳-۳ بلوک دیاگرام ماژول دریافت کننده

### ۳-۲-۴ تجزیه کننده<sup>۱</sup> صفحه وب

وظیفه مهم یک تجزیه کننده بعد از دریافت یک صفحه وب تجزیه کردن و استخراج اطلاعات مهم از داخل صفحه وب می‌باشد. حالا به معرفی چند مشکل در این ماژول و راهکارهای برای آنها می‌پردازیم.

#### ۳-۲-۴-۱ شناسایی و اولویت بندی و حذف تگهای بی استفاده

یکسری از تگ ها در صفحات وب به هیچ عنوان جنبه اطلاعات نداشته و فقط برای کارهای مانند زیباسازی صفحات وب، کارهای پیش زمینه<sup>۱</sup> (مثلا وقتی روی یک دکمه کلیک می‌کنید یک کار خاص انجام بگیرد) و ... استفاده می‌شوند. یکی از وظایف تجزیه کننده شناسایی این تگ ها در صفحه وب و حذف آنها است مطابق با شکل ۳-۲ قسمت (a). و یا حتی باید به تگ ها اولویت داده شود مثلا در یک صفحه خاص بعضی تگ ها دارای اهمیت و اولویت بالای هستند به عنوان یک مثال فرض کنید وارد یک صفحه وب شدید که کاملا در زمینه معرفی گل‌ها است پس محتوای تگ های مانده <img> و <figure> در این صفحه باید

<sup>۱</sup>Parser

اولویت بالای را بگیرند ( البته اهمیت درباره پرس وجوهای با موضوعیت گل هم دارای اهمیت می باشد ).  
چه راهکاری برای این امر باید انجام شود و تجزیه کننده چگونه باید تگ های بدون استفاده را از صفحه جدا کرده و برای هر تگ یک اولویت را در نظر بگیرد.

**راه حل :** بعضی از تگ ها فاقد محتوای مفید هستند مانند تگ style همانطور که در شکل ۳-۴ نمایش داده شده این تگ فقط جهت بهبود قالب صفحات وب می باشد و هیچگونه اطلاعات مفیدی در باره صفحه وب به ما نمی دهد اما تگ span نسبتا مفید است و حاوی اطلاعاتی درباره صفحه وب.

```
<style>
    .view-input{
        border-left:none;
        border-right:none;
        border-top:none;
        border-bottom:1px solid blue;
    }
</style>
<span
    type="text"
    style="color:#34ED29"
    class="view-input"
>
    محتوای این سایت درباره هوش مصنوعی است
</span>
```

شکل ۳-۴ مثالی از HTML

## نوآوری

راهکار دیگر که در این پروژه استفاده شده است حذف تگ های می باشد که به مرور زمان ارزش و اهمیت خود را درباره یک پرس وجوی خاص از دست می دهند. به عنوان مثال اگر موضوع مورد بررسی موتور جستجو درباره اخبار باشد طبق بررسی های انجام گرفته معمولا تگ های های مانند <p> و <h1> اولویت و داده های بیشتری نسبت به تگ <span> دارند. به عنوان مثال دیگر در صفحات وبی مانند pinterest.com اهمیت تگ <img> که برای نمایش تصاویر می باشد بسیار بالا است و این صفحه غالبا از تصاویر تشکیل گردیده است.

ولی در صفحه وبی مانند www.cnn.com ما می بینیم که تگ <a> که یک تگ اتصالی بین صفحات وب



است بسیار زیاد می‌باشد که غالباً هم همینطور است به دلیل ماهیت خبری این صفحه باید لینک اتصال به صفحات خبری دیگر هم زیاد باشد البته راهکار فقط روی زیاد بودن یک تگ در صفحه متمرکز نیست بلکه با استفاده از تکنیک یادگیری عمیق<sup>۱</sup> و شبکه های عصبی<sup>۲</sup> به مرور زمان متوجه می‌گردد که در صفحه وب مشخص مثلاً تگ <h1> چه اهمیتی دارد. با بالا رفتن اولویت یک تگ در یک دامنه ارزش اطلاعاتی آن تگ در آن دامنه در یک پرس‌وجوی خاص هم بالا می‌رود.

پس تجزیه کننده برای هر صفحه یک جدول از لیست تگ های مهم با اولویت را مشخص می‌کند. و این لیست مدام در حال به روز رسانی می‌باشد. بعد از بررسی صفحات و استخراج ویژگی ها و با استفاده از پاداشی که از کاربران بازدید کننده می‌گیرد و آموزش های که می‌بیند این اطلاعات به روز می‌شود. پس این تگ در این صفحه به مرور زمان اهمیت بالای می‌گیرد.

در حقیقت ما در این پروژه و در این قسمت از یک تکنیک یادگیری عمیق برای امتیاز دهی هر صفحه استفاده می‌کنیم. البته بعضی از تگ ها و ویژگی های آنها هم در حذف و یا کاهش اولویت آنها اهمیت دارد که یا به فرم تجربی قابل استفاده است و یا با روش های یادگیری ماشین و شبکه عصبی. به عنوان مثال فرض کنید یک تگ شکل دارای ویژگی طول ۲ پیکسل می‌باشد بدیهی است این اندازه کوچک کمکی به موتور جستجو خصوصاً در زمینه پردازش تصویر نمی‌کند و یا مثلاً یک تگ پاراگراف که

فاقد محتوا است هم کمکی به پردازش نمی‌کند. و یا در زمینه یادگیری فرض کنید ۱۰۰ بار موتور جستجو صفحه مورد نظر را بررسی کرده است و به مرور متوجه شده ارزش تگی مانده <li> (این تگ برای لیست بندی صفحات بسیار مهم است) در صفحه بسیار مورد توجه کاربران قرار گرفته و یا تگی به نام <meta>

---

<sup>۱</sup> deep learning

<sup>۲</sup> neural network

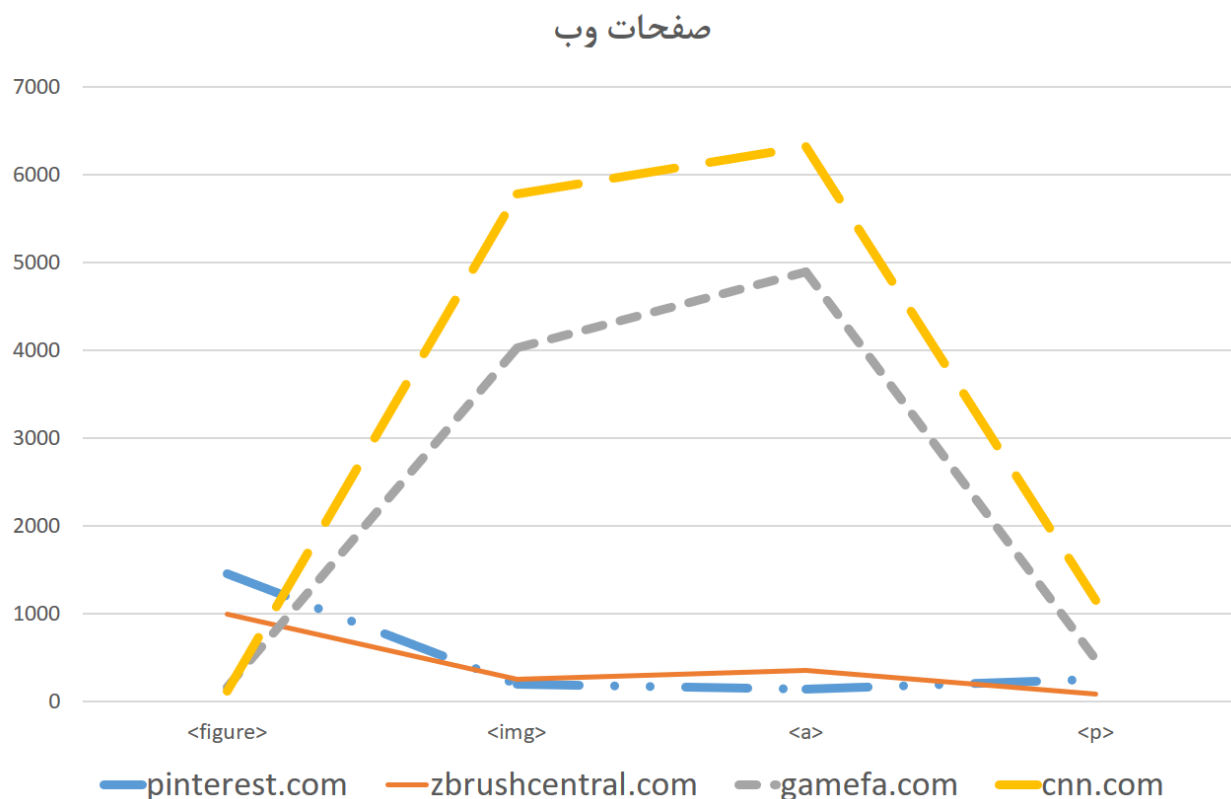
با ویژگی به نام keyword که در صفحه وب مشخص کننده کلمات کلیدی سایت است در صفحه این کلمات با موضوع صفحه ارتباطی ندارند پس به مرور زمان ارزش این تگ در این صفحه کم رنگ تر می گردد.

پس موتور جستجو به طور مداوم ارزش هر تگ را به فرم دوره ای در هر صفحه وب تغییر می دهد.

در این پروژه ما بعد از استخراج متن از نزدیک تصاویر صفحات وب به ایجاد ۸۰۰ گروه مختلف (پرس و جو) با استفاده از یادگیری عمیق و متن می پردازیم. و شناسایی و گروه بندی اشیاء در صفحه وب با کمک این ۸۰۰ گروه مشخص می شود.

در شکل ۳-۵ به دلخواه ۴ دامنه انتخاب شده اند ۲ دامنه با محوریت اخبار (gamefa.com، cnn.com) و ۲ دامنه دیگر با محوریت تصویر (pinterest.com، zbrushcentral.com).

همچنین ۴ تگ خاص از نظر تعداد تکرار و شباهت با پرس و جوهای مشخص شده در پروژه در این دامنه ها بررسی شده اند. (تگ تصویری <img> و <figure>، تگ ارتباطی <a>، تگ پاراگراف <p>) البته این ۴ دامنه شامل بیش از هزاران صفحه وب داخلی اند و اطلاعات مفید در زمینه های مختلف را در این ۴ دامنه بررسی شده است، همانطور که می بینید اهمیت تگ های تصویری و اتصالی در صفحات مربوط به گالری تصاویر بسیار بالا بوده ولی در صفحات اخباری تگ های پاراگراف و اتصال از تکرار و اهمیت موضوعی بالای برخوردارند. در اینجا ما ۸۰۰ گروه مشخص شده را با تمام متن ها و تصاویر نزدیک آنها در هر صفحه وب دامنه ها مقایسه کرده و درجه شباهت آنها را بدست می آوریم.



شکل ۳-۵ بررسی تگ ها در صفحات وب

### اهمیت مکان تگ ها

آیا قرار گرفتن تگ در مکان خاص بر اهمیت آن می افزاید و یا تاثیر در آن ندارد. به عنوان مثال اگر یک تگ در ابتدای صفحه بیاید از نظر ارزش و اهمیت با تگی که در انتهای تصویر می آید تفاوتی با هم دارند.

البته بیشتر اهمیت یک تگ روی محتوای داده آن است ولی مکان قرار داشتن تگ هم تاثیری هر چند ناچیز در اهمیت تگ ایفا می کند. به عنوان مثال تگ <h1> که معمولاً در ابتدای صفحات وب می آید و بیشتر برای نمایش سرتیتر اطلاعات است از اهمیت بسیار بیشتری به نسبت تگ مشابهی که در انتهای صفحه می آید دارد. زیرا در ابتدای صفحه جنبه یک عنوان و در انتهای صفحه بیشتر جنبه تبلیغ و بی ارزشتری را دارد. در اینجا منظور از مکان تگ فاصله تگ تا ابتدای صفحه وب می باشد. در شکل ۳-۶ مشاهده می کنید

که به عنوان مثال تگ <p>

در فاصله ۵ و تگ <span> در فاصله ۳ تا ریشه صفحه وب (تگ <body>) قرار گرفته اند.

```
<body contenteditable="false">
  <div class="level1">
    <div class="level2">
      <span> level 3 </span>
      <div class="level3">
        <p> i love my mother </p>
      </div>
    </div>
  </div>
  <div class="level2">
    <li>1</li>
  </div>
</body>
```

سطح چهار نسبت به ریشه

سطح پنج نسبت به ریشه

سطح سه نسبت به ریشه

شکل ۳-۶ مکان تگ در درخت html

**نواوری :** در این پروژه برای هر تگ در هر صفحه یک مکان قرار داده می شود (منظور افزودن یک ویژگی است که مکان این تگ را در صفحه وب معین می کند) و بعد از بررسی تگ ها و داده ها به فرم یادگیری عمیق و به

مرور زمان متوجه اهمیت مکان تگ در یک صفحه خاص می شویم. یعنی بعد از استخراج متن و تصویر از تگ ها و بررسی درصد شباهت آنها با پرس وجوهای داخل موتور جستجو به اهمیت و شباهت آنها با یک پرس وجوی خاص می رسیم.

## ۲-۴-۲-۳ استخراج ویژگی از تگ

بعد از تجزیه صفحه وب و حذف تگ های غیرضروری (پیش پردازش) و دادن امتیاز به هر تگ باید تگ ها مورد پردازش قرار بگیرند و بردار ویژگی آنها برای مراحل بعدی استخراج گردد. در اینجا هدف یافتن تصاویر و متون مفید و مرتبط با پرس وجوهای مشخص شده می باشد. البته همانطور که گفته شد هر تگ حاوی یکسری اطلاعات و ویژگی می باشد که در صورت نیاز این اطلاعات از تگ مربوطه استخراج می شود.

## نوآوری :

در این پروژه ما یک ماشین یادگیری بدون ناظر داریم (تصاویر فاقد برچسب و گروه مشخص هستند و فقط از صفحات وب گرفته شده اند. همینطور ۸۰۰ پرسوجوی ایجاد شده نیز از متون نزدیک صفحات وب استخراج شده اند) پس باید ماشین به فرم تجربی بیاموزد و قابلیت شناسایی تصاویر مفید را داشته باشد. در اینجا به جز گرفتن تصاویر و استخراج بردار ویژگی برای آنها ما از تکنیک استخراج ویژگی از متون نزدیک تصویر هم استفاده کرده ایم. متون نزدیک به تصویر را با استفاده از یک تکنیک شبکه عصبی چند لایه و پردازش متن مورد پردازش قرار داده و بعد از یادگیری می شود به هر متن نزدیک تصویر یک درجه اهمیت داد بدین صورت که چقدر این تصویر بر اساس این متن با یک پرسوجوی داخل موتور جستجو شباهت دارد. البته این قسمت فقط درصدی از اهمیت تصویر را معین می کند. خود تصویر و استخراج ویژگی های آن و درصد شباهت آن با دیگر تصاویر یک پرسوجوی خاص در این فرایند تاثیر دارد.

به عنوان مثال در سایت [pinterest.com](http://pinterest.com) بعد از بررسی مشخص شد ویژگی `alt` که در تگ `<img>` قرار دارد و خلاصه ای از موضوع تصویر مورد جستجو را به ما می دهد کمک شایانی در تشخیص گروه تصویر می کند ولی در سایت دیگری مانده <http://www.flowerpicturegallery.com> که نمایش دهند تصاویری در زمینه گل می باشد. تگ پدر تصویر از نوع `<span>` دارای توضیحات خوبی درباره نوع و ارزش تصویر می باشد و تگ `alt` خود تصویر مطلب خوبی درباره تصویر نمی دهد (بیشتر به معرفی خود دامنه پرداخته است). پس باید برای هر صفحه وب بر حسب اهمیت هر تگ به محتوای متنی آن اهمیت و ارزشی بدهیم. دقت شود که برای هر تصویر ما چند لایه بالاتر آن را هم مورد بررسی قرار می دهیم و متون داخل آنها را با ارزش و اهمیت های مختلف جمع آوری می کنیم.

بعد از استخراج متون از این تگ ها به استخراج ویژگی از خود تصویر می پردازیم.

### ۳-۲-۴-۳ استخراج ویژگی از تصویر

قبل از استخراج ویژگی‌های یک تصویر لازم است درباره ماهیت تصویر اطلاعاتی کسب کنیم.

#### ۱-۳-۲-۴-۳ معرفی تصویر

تصویر دیجیتال، از تعداد بسیار زیادی از مربع‌های کوچک به نام پیکسل<sup>۱</sup> تشکیل شده‌است. هر پیکسل دارای یک شماره رقمی<sup>۲</sup> است که بیان‌گر مقدار روشنایی آن پیکسل است. مقدار عددی پیکسل هم از صفر (سیاه) تا ۲۵۵ (سفید) متغیر است. به این نوع از تصاویر، تصاویر رستری<sup>۳</sup> هم می‌گویند. هر تصویر رستری از تعدادی سطر و ستون تشکیل شده است. در زیر چند نکته در مورد تصویر به صورت مختصر آورده شده است [17].

**معرفی پیکسل:** برای مشخص کردن رنگ یک پیکسل، روش‌های مختلفی استفاده می‌شود. آنچه که متداول‌تر است RGB<sup>۴</sup> (آبی قرمز و سبز ۳ رنگ اصلی) است، که ۳ کانال مختلف برای ۳ رنگ قرمز، سبز و آبی در نظر می‌گیرند. اما در پردازش تصویر از فضاهاى رنگی دیگر استفاده بیشتری می‌شود. برای مثال فضای رنگ HSV. هر کانال دارای ۲۵۶ حالت خواهد بود در نتیجه هر پیکسل می‌تواند ۱۶۷۷۷۲۱۶ رنگ مختلف را نشان دهد.

**تفکیک‌پذیری تصویر:** تفکیک‌پذیری تصویر به تعداد پیکسل‌ها در طول و عرض تصویر بستگی دارد. مانده شکل 3-7 که همانطور که می‌بینید تعداد پیکسل‌ها نمایش داده شده است [17].

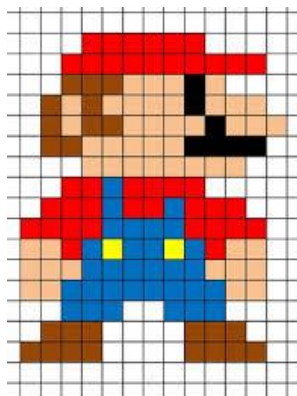
---

<sup>۱</sup> Pixel

<sup>۲</sup> digital number

<sup>۳</sup> raster

<sup>۴</sup> red-green-blue



شکل ۳-۷ تصویر به تفکیک پیکسل [17]

**مشکل در پردازش تصویر :** در موضوعاتی مانده پردازش تصویر تشخیص چهره یا هر موضوع خاص دیگری استخراج ویژگی توسط یک الگوریتم مشخص انجام می‌گیرد مانده SIFT<sup>۱</sup> یا PCA<sup>۲</sup> و ... ولی در اینجا ما با تصاویر با فرمت ها و اندازه هایی مختلف مواجه هستیم و از طرفی خروجی ما نیز درصد شباهت تصویر با لیست پرس وجوهای مختلف می‌باشد که از متون نزدیک تصاویر در صفحات وب بدست می‌آید. برای استخراج ویژگی و پردازش تصویر و متون [10] در این پروژه ما از شبکه‌های عصبی کانولوشن<sup>۳</sup> و یادگیری عمیق<sup>۴</sup> استفاده کرده‌ایم.

## ۲-۳-۴-۲ شبکه های عصبی کانولوشن

شبکه های عصبی کانولوشن شبیه شبکه های عصبی مصنوعی معمولی می‌باشند. این نوع شبکه ها متشکل

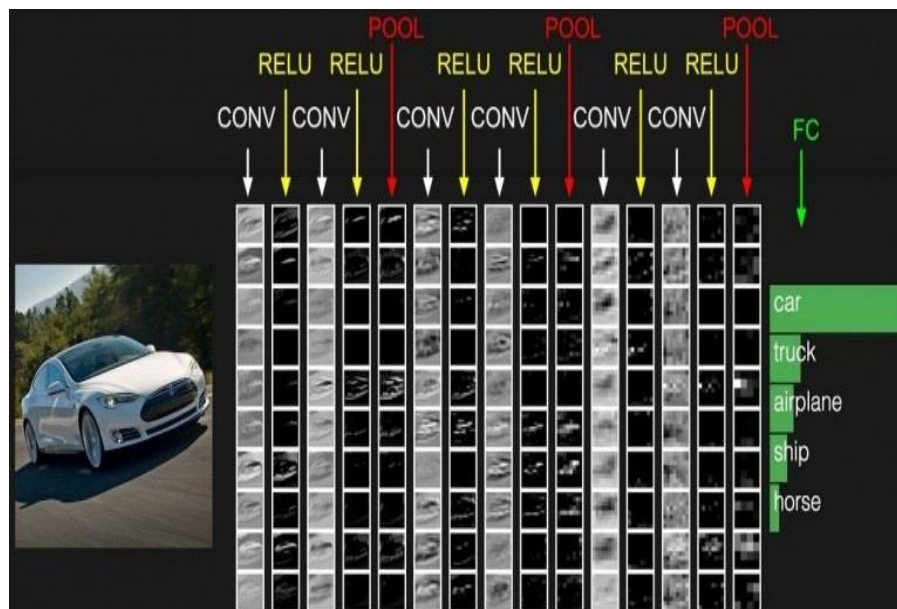
از نورونهایی<sup>۵</sup> با وزنها و بایاسهای<sup>۶</sup> قابل یادگیری هستند. [20] هر نورون تعدادی ورودی دریافت کرده و

<sup>۱</sup>Scale-invariant feature transform  
<sup>۲</sup>Principal component analysis  
<sup>۳</sup>convolution neural network (CNN)  
<sup>۴</sup>deep learning  
<sup>۵</sup>bios<sup>۶</sup>

<sup>۷</sup>neurons

<sup>۸</sup>bios

سپس حاصل ضرب وزن‌ها در ورودی‌ها را محاسبه کرده و در انتها با استفاده از یک تابع تبدیل (فعال سازی)<sup>۱</sup> غیرخطی درصد شباهت ورودی با بردار خروجی (لیست پرس‌وجو‌ها) را تعیین می‌کند. همین‌طور شبکه شامل یک تابع امتیاز مشتق پذیر<sup>۲</sup> را شامل می‌شود. که در یک طرف آن پیکسل‌های خام تصویر ورودی و یا متن نزدیک تصویر و در طرف دیگر آن امتیازات مربوط به هر پرس‌وجو قرار دارد. این نوع شبکه‌ها هنوز یک تابع هزینه مثل SVM یا Softmax (در لایه آخر) تمام مرتبط<sup>۳</sup> دارند و تمامی نکات مطرح در مورد شبکه‌های عصبی معمولی در اینجا هم صادق است. پس تفاوت این دو شبکه در چیست؟ معماری‌های شبکه‌های عصبی کانولوشن بصورت صریح فرض می‌کنند که ورودی‌های آنها تصاویر هستند . با این فرض ما می‌توانیم ویژگی‌های مشخصی را درون معماری تعبیه کنیم. با این عمل تابع پیشرو را میتوان بصورت بهینه‌تر پیاده سازی کرد و همین‌طور با این کار میزان پارامترهای شبکه نیز بشدت کاهش پیدا میکند . در شکل ۳-۸ یک شبکه کانولوشن تصویر نمایش داده شده است.



شکل ۳-۸ شبکه عصبی CNN تصویر [20]

<sup>۱</sup> activator function

<sup>۲</sup> differentiable

<sup>۳</sup> fully connected



## خلاصه معماری شبکه عصبی کانولوشن :

شبکه های عصبی یک ورودی دریافت می کنند (در قالب یک بردار) و سپس آنرا از تعدادی لایه مخفی<sup>۱</sup> عبور می دهند. و در نهایت یک خروجی که نتیجه پردازش لایه های مخفی است در لایه خروجی نمایش داده می شود که نمایش دهنده درصد شباهت ورودی با یکی از پرس وجوهای خروجی می باشد. هر لایه مخفی از تعدادی نورون تشکیل شده که این نورون ها به تمام نورون های لایه قبل از خود متصل می شوند . نورونهای هر لایه بصورت مستقل عمل کرده و هیچ ارتباطی با یکدیگر ندارند [20]. آخرین لایه تماماً متصل به لایه خروجی معروف است و معمولاً نقش نمایش دهنده امتیاز هر طبقه را ایفا می کند.

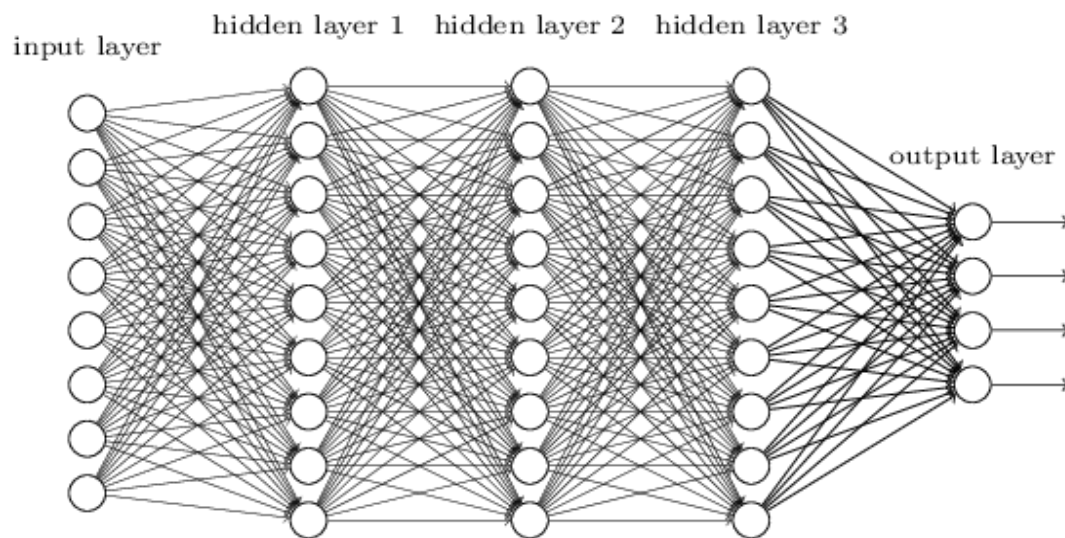
شبکه های عصبی معمولی برای تصاویر بخوبی مقیاس پذیر نیستند. بعنوان مثال فرض کنیم دارای مجموعه ای از تصاویر با اندازه ای  $32*32*3$  هستیم (۳۲ پیکسل عرض، ۳۲ پیکسل ارتفاع و ۳ کانال رنگ ). بنابراین یک نورون با اتصال کامل<sup>۱</sup> در لایه مخفی اول یک شبکه عصبی معمولی به  $3072$  وزن نیاز خواهد داشت. این مقدار شاید در نظر اول مقدار قابل توجهی بنظر نیاید اما بطور واضح مشخص است که این معماری تماماً متصل قابل استفاده برای تصاویر بزرگتر نخواهد بود مخصوصاً در صفحات وب که هر تصویر اندازه متفاوت دارد. برای مثال یک تصویر با اندازه متعارف تر مثل  $200*200*3$  باعث میشود که یک نورون  $200*200*3 = 120,000$  وزن داشته باشد. علاوه بر این ما قطعاً خواهان تعداد بیشتری از این نورون ها خواهیم بود. مشخص است این تعداد اتصال باعث افزایش نمایی پردازش شده و با افزایش تعداد پارامترها باعث ایجاد overfitting خواهد شد [20].

---

<sup>۱</sup> hidden layer

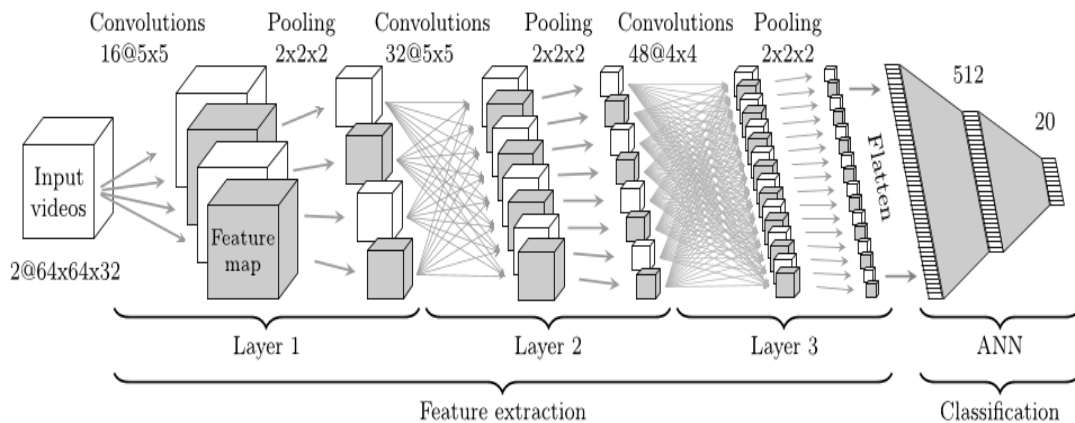
## توده های سه بعدی از نورونها<sup>۱</sup>

شبکه های عصبی کانولوشن از این واقعیت که ورودی شامل تصاویر است استفاده کرده و معماری شبکه را به روش معقولی محدود کردند. برخلاف یک شبکه عصبی، لایه های یک شبکه عصبی کانولوشن شامل نورونهایی است که در سه بعد عرض، ارتفاع و عمق قرار گرفته اند. کلمه عمق در اینجا اشاره به بُعد سوم یک توده فعال سازی<sup>۲</sup> دارد و به معنای عمق یک شبکه عصبی کامل که به معنای تعداد لایه های موجود در آن است نمی باشد. هر نورون در هر لایه بجای اتصال با تمام نورون ها در لایه قبل تنها به ناحیه کوچکی از لایه قبل از خود متصل است [20]. با رسیدن به انتهای معماری شبکه ما اندازه تصویر را کاهش می دهیم بگونه ای که در انتها تصویر کامل ورودی ما به یک بردار حاوی امتیاز طبقه ها کاهش پیدا میکند (۸۰۰ پرس و جوی انتخابی) و ما با یک بردار که حاوی امتیاز هر طبقه است مواجه خواهیم بود. شکل ۳-۹ یک شبکه چند لایه و در شکل ۳-۱۰ یک شبکه کانولوشن<sup>۳</sup> بعدی را نمایش می دهند.



شکل ۳-۹ شبکه عصبی چند لایه [20]

<sup>۱</sup> 3d volumes of neurons    <sup>۲</sup> activation volume



شکل ۳-۱۰ شبکه کانولوشن ۳ بعدی [20]

حالا به معرفی لایه‌های شبکه کانولوشن می‌پردازیم.

- **لایه کانولوشن:** در این لایه ما بردار ورودی را در بردار ویژگی‌ها<sup>۱</sup> کانوال می‌کنیم. در حقیقت در

این لایه هر قسمت از لایه قبلی با یک پنجره مشخص کانولوشن شده و نتیجه به لایه بعدی داده

می‌شود [20].

- **لایه RELU یا لایه نرمال کننده<sup>۱</sup>:** ممکن است در لایه کانولوشن یکسری از داده‌ها مقادیر نامتعارفی

(مثلا منفی) بگیرند. پس بهتر است بردار نرمالایز شود. در این جا برای نرمال از RELU استفاده شده

است. مطابق با فرمول ۳-۱ که جهت جلوگیری از کسب مقادیر منفی توسط لایه‌های مختلف [20].

$$X_i = \text{Maximum}(X_i, 0) \quad (\text{فرمول ۳-۱})$$

- **لایه Pool<sup>۲</sup>:** در این مرحله ما به کاهش اندازه بردار ورودی می‌پردازیم ( مثلا میانگین گیری ۹ پیکسل

مجاور هم ). این کار باعث دستیابی به ویژگی‌های سطح بالا<sup>۳</sup> می‌شود. (در پردازش تصویر منظور از

ویژگی های سطح پایین همان شناسایی لبه‌ها<sup>۴</sup> و شناسایی منحنی‌ها<sup>۵</sup> و ... است و ویژگی های سطح

---

curve<sup>۵</sup>      edges<sup>۴</sup>      high level features<sup>۳</sup>      pooling layer<sup>۲</sup>      rectified linear unit layer<sup>۱</sup>

بالا مثلاً شناسایی دست، بینی و همینطور در متن منظور گروه‌ها و کلمات کلیدی و مهم هستند). یکی دیگر از اهداف این لایه نگاشت بردار ورودی به بردار خروجی و نتیجه می‌باشد (در این پروژه همان ۸۰۰ پرس‌وجوی خروجی).

- **لایه اتصال کامل:** این لایه همان امتیاز دهنده بردار ورودی به بردار خروجی و داخل موتور جستجو می‌باشد و تصویر بر حسب درصد شباهت با پرس‌وجوها در یکی یا چند تا از این پرس‌وجوها قرار می‌گیرد.

#### ۴-۲-۳ استخراج ویژگی از متن

در این پروژه علاوه بر خود تصاویر ما از متون صفحات وب هم برای کمک به طبقه‌بندی و شناسایی اشیاء استفاده می‌کنیم. بخاطر ماهیت بدون ناظر بستر وب، متون نزدیک به تصویر در صفحات وب می‌توانند کمک خوبی برای شناسایی پرس‌وجوها باشند. حال به معرفی چالش‌های استخراج یک لیست پرس‌وجو از داخل متن می‌پردازیم

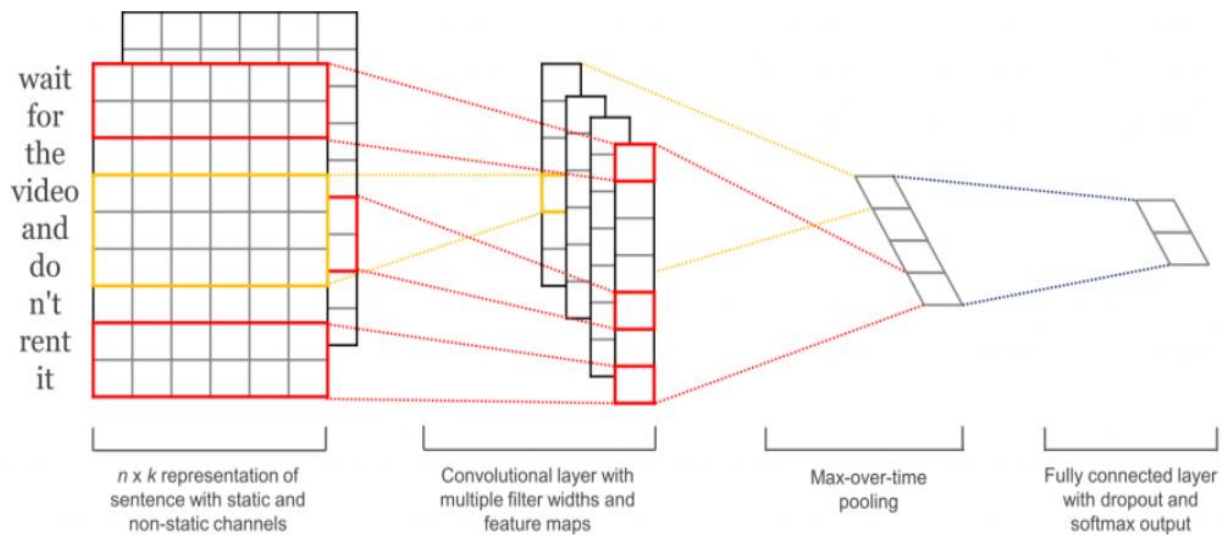
#### چالش‌ها

- چگونه تشخیص دهیم متن درون صفحه وب با تصویر داخل صفحه وب مرتبط است.
- چگونگی تشخیص اهمیت و اولویت متن با تصویر: در نزدیک تصویر (در سطوح مختلف درخت html) متونی یافت می‌شود به چه نحوی باید به هر کدام از این متون اولییتی داده شود.
- چگونگی ایجاد بردار ورودی از متن استخراج شده از صفحات وب.

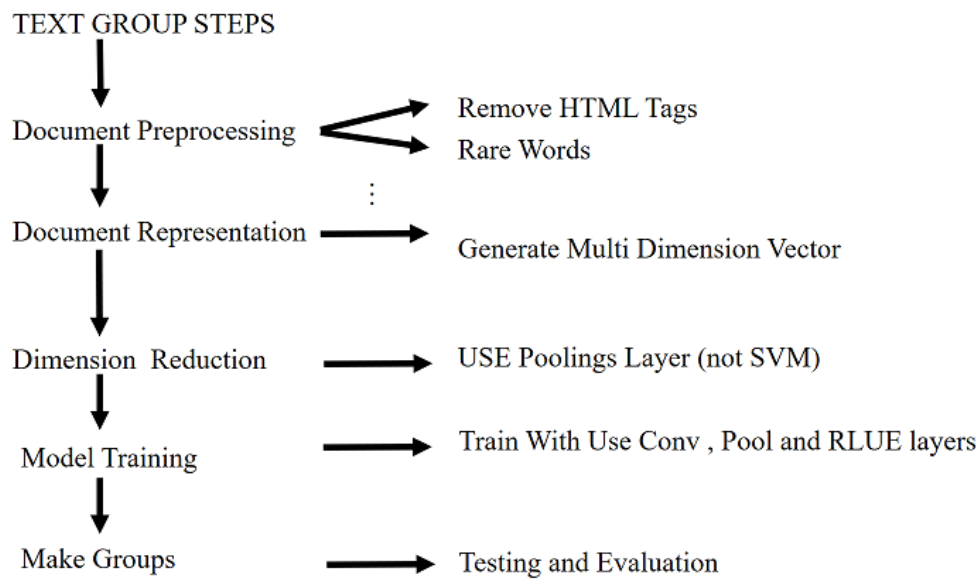
## نوآوری

در اینجا ما به استخراج متون نزدیک تصاویر پرداخته و سپس این متون را در یک پاراگراف جمع کرده و پیش پردازش‌های را روی این متن انجام می‌دهیم مانده حذف کلمات تکراری و حذف کاراکترهای ویرایش مانده خط جدید (n) و... در مرحله بعدی به ایجاد بردار ویژگی از متن پرداخته و این بردارها را به یک شبکه عصبی کانولوشن داده و نتیجه آن یک بردار حاوی ۸۰۰ طبقه مختلف از اشیاء می‌باشد. (در این پروژه ما فقط از ۸۰۰ گروه مختلف استفاده کرده‌ایم) [21]. این فاز جهت ایجاد طبقات اصلی پرس‌وجو در موتور جستجو اجرا می‌شود. در این فاز از ۱۰ دامنه و بیش از ۱۰۰۰۰ صفحه وب استفاده شده است. خروجی این مرحله لیست گروه‌های اولیه مورد جستجو می‌باشد.

سپس با استفاده از شبکه عصبی کانولوشن و یادگیری عمیق ما میزان شباهت هر متن به دست آمده نزدیک تصویر را با تمام پرس‌وجوهای داخل موتور جستجو بررسی می‌کنیم. از طرفی به مرور زمان و گرفتن امتیاز هر متن میزان اهمیت مکان متن در یک صفحه وب خاص و درمورد یک پرس‌وجوی خاص مشخص می‌شود. مطابق شکل ۳-۱۱ متن ورودی به یک شبکه عصبی چند لایه کانولوشن به صورت یک بردار داده می‌شود و خروجی آن میزان شباهت متن با هر پرس‌وجو می‌باشد. برای بردار ویژگی ورودی هم بجای تصویر ما از ماتریس متنی استفاده می‌کنیم هر سطر این ماتریس از حروف یک کلمه تشکیل شده است. بردار خروجی این شبکه نیز به تعداد پرس‌وجوهای پایگاه داده و مقدار شباهت متن ورودی با آنها می‌باشد. در شکل ۳-۱۲ نیز مراحل استخراج پرس‌وجوهای اصلی موتور جستجو نمایش داده شده است.



شکل ۳-۱۱ بردار ورودی شبکه کانولوشن متن [21]

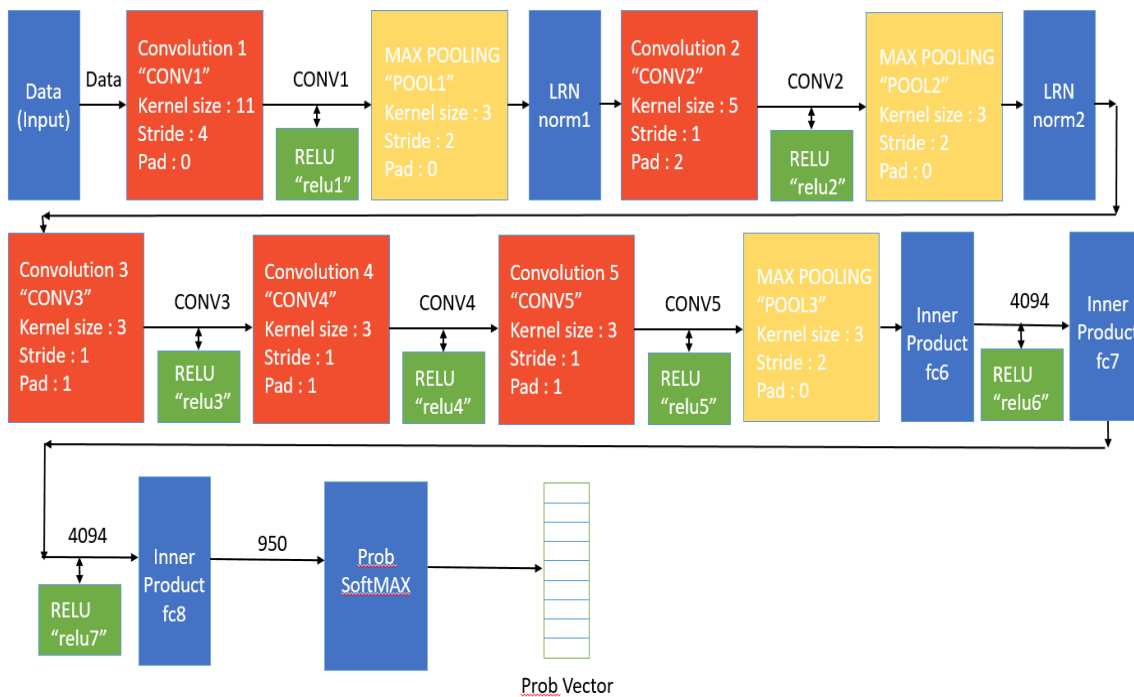


شکل ۳-۱۲ بردار ورودی شبکه کانولوشن متن [21]

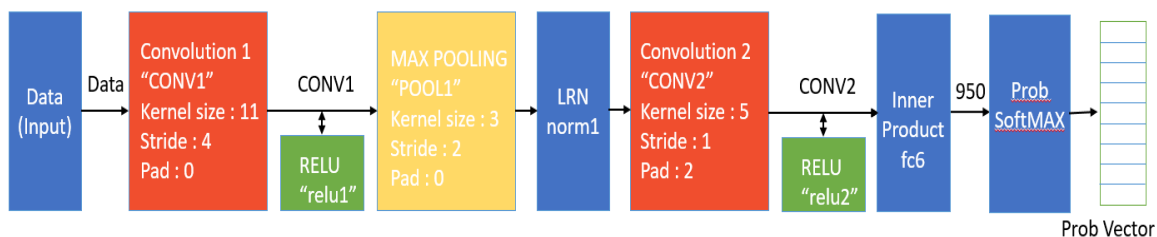
### معرفی فریم ورک Caffe

در این قسمت به معرفی یکی از برنامه‌های بسیار عالی جهت یادگیری عمیق با نام Caffe (کافی) می‌پردازیم. این برنامه برای سیستم عامل‌های (ویندوز، لینوکس و یونیکس) و زبان‌های برنامه نویسی مختلف (c++، python و ...) ارائه گردیده است. این برنامه توسط آقای Yangqing Jia توسعه داده شد و در حال حاضر بصورت

کاملاً رایگان توسط دانشگاه برکلی توسعه داده می‌شود. در این پروژه برای یادگیری تصاویر و متون صفحات وب از دو شبکه عصبی عمیق استفاده کرده‌ایم. برای شبکه عصبی تصویر مطابق با شکل ۳-۱۳ از یک شبکه عصبی بزرگ استفاده شده است که حاوی چندین لایه کانولوشن، Pooling و RELU به همراه روش‌های مانده softmax می‌باشد. همچنین برای انتخاب و یادگیری متن (ایجاد ۸۰۰ گروه پرس‌وجو) هم از یک شبکه عصبی کوچکتر مطابق با شکل‌های ۳-۱۴ استفاده گردیده است. برای این برنامه بعد از مشخص کردن مراحل و لایه‌های شبکه عصبی آنها را در فایل‌های با پسوند prototxt ذخیره می‌کنیم. در این فایل به معرفی لایه‌های مختلف تعریف شده در شبکه عصبی کانولوشن پرداخته و پارامترهای یادگیری را مشخص می‌کنیم. در شکل ۳-۱۵ یک نمونه از این تعریف را می‌بینید.



شکل ۳-۱۳ بلوک دیاگرام یادگیری تصویر در کافی



شکل ۳-۱۴ بلوک دیاگرام یادگیری متن در کافی

```

1  name: "CaffeNet"
2  layer {
3      name: "data"
4      type: "Input"
5      top: "data"
6      input_param { shape: {
7          dim: 10 dim: 1 dim: 255 dim: 255
8      } }
9  }
10 layer {
11     name: "conv1"
12     type: "Convolution"
13     bottom: "data"
14     top: "conv1"
15     convolution_param {
16         num_output: 96
17         kernel_size: 11
18         stride: 4
19     }
20 }

```

شکل ۳-۱۵ نمونه‌ای از فایل prototxt در کافی

همانطور که در تصویر ۳-۱۵ مشخص است، هر لایه در شبکه کانولوشن با لایه های بالای و پایینی خود

ارتباط دارد به عنوان مثال در شکل لایه کانولوشن اول ما با نام conv1 از طرفی به لایه ورودی متصل

است.

در مرحله بعدی و قبل از شروع یادگیری ما باید مجموعه داده آموزش و آزمون خود را بسازیم که برای این

منظور ابتدا تصاویر را از پایگاه داده استخراج کرده (تصاویر دریافت شده در صفحات وب) و برای سادگی

یک پیش پردازش مطابق زیر انجام می‌دهیم

- یکی کردن اندازه تصاویر :

به دلیل عدم یکی بودن اندازه تصاویر و نوع آنها (ممکن است بعضی تصاویر ۳ کاناله و بعضی تک کاناله



و خاکستری باشند) در این پیش پردازش و به منظور ایجاد ثابت کردن اندازه پنجره ورودی تصاویر را به اندازه ثابت ۲۵۵\*۲۵۵ به منظور ثابت کردن پنجره ورودی تبدیل می کنیم. البته می توان تصویر را به قطعه های مختلف با اندازه ثابت تبدیل کرد و آنرا به عنوان نتیجه به شبکه یادگیری فرستاد که به علت محدودیت سخت افزاری این کار انجام نگردید.

- **خاکستری کردن تصاویر:**

به منظور ساده سازی پردازش و کاهش اندازه ورودی بجای استفاده از یک تصویر با ۳ کانال آنرا به یک کانال خاکستری تبدیل کردیم ( البته ممکن است تصاویر ۳ کاناله و خاکستری در صفحات وب وجود داشته باشند که به منظور یکی کردن غالب تصاویر همه را به فرمت خاکستری تبدیل کردیم).

- **تغییر همه تصاویر به یک فرمت :**

در بستر وب تصاویر با فرمت های مختلف وجود دارند (فرمت های مانده PNG، JPG، WEBM و...) به این منظور تصاویر را به یک فرمت مشخص که در اینجا فرمت انتخابی ما JPEG است به دلیل فشردگی سازی خوب و بالای تصویر و عدم حذف زیاد جزئیات از داخل آن انتخاب کرده ایم.

- **حذف تصاویر تکراری :**

در این بخش تصاویر تکراری حذف می گردند. جهت حذف به راحتی می توان تصاویری را که بیش از ۱۰ بار در پایگاه داده تکرار شده اند از داخل آن حذف کرد. دلیل این امر این است که به احتمال بالا این تصاویر یکی از موارد زیر می باشند.

- یک تصویر عمومی در صفحه وب می باشد مانده لوگوی یک دامنه خاص که با نرخ تکرار بالای در صفحات وب قرار می گردد.

- تصویر تبلیغاتی می باشد که در صفحات زیادی و به تکرار دیده می شود.

بعد از این پیش پردازش در گام بعدی ما متن ها و تصاویر به همراه فایل آموزش به برنامه کافی جهت آموزش می دهیم. در این پروژه تعداد ۱۰۰۰۰ تکرار را برای این شبکه عصبی به همراه ۸۰۰ گروه مشخص تعیین کرده ایم. و بعد از هر ۲۰۰ تکرار یک مرحله تست را هم داریم. مشخصات کامل آموزش در جدول ۳-۳ قابل مشاهده است. در انتها و بعد از فاز آموزش برنامه کافی یک فایل حاوی تمام وزنها (یادگیری شده) و لایه های مختلف شبکه عصبی کانولوشن را با پسوند caffeprototype به عنوان خروجی می دهد. که برای تست تصاویر ورودی و متن ها در مراحل بعدی استفاده می شود. به این صورت که از این مرحله به بعد تصاویر و متون بدست آمده از صفحات وب به این شبکه داده شده و خروجی میزان شباهت متن و یا تصویر ورودی به پرس وجوهای صفحه وب را نمایش می دهد (بردار خروجی).

جدول ۳-۳ توضیحات فاز آموزشی

ردیف	نوع	توضیحات
۱	مشخصات سخت افزار آموزش	<ul style="list-style-type: none"> <li>• 2 عدد کارت گرافیکی مدل Geforce 840M با حافظه ۸ گیگا بایت (تکنیک SLI) و ۸۰۰ واحد گودا</li> <li>• یک عدد پردازنده core i5 4600U با فرکانس کاری ۱۸۰۰ GHZ و ۴ هسته</li> <li>• ۱۲ گیگا بایت حافظه مدل DDR3</li> </ul>
۲	تعداد تکرار فاز آموزشی	۱۰۰۰۰ بار تکرار
۳	تعداد گروه مشخص شده اولیه	۸۰۰ گروه مشخص گردید که به عنوان پرس وجوهای موتور جستجو مشخص شده اند
۴	تعداد تصاویر	۵۰۰۰۰ تصویر از بیش از ۱۹۲۰ دامنه و ۲۳۳۵۰ صفحه وب متفاوت
۵	فاز تست بعد از چند آموزش	۲۰۰
۶	زمان آموزش	زمان تکرار ۱۰۰۰۰ تکرار در این پروژه حدوداً چهل روز به طول انجامید

## ۵-۲-۳ صف اولویت آدرس‌ها

همانطور که قبلاً هم گفته شد موتور جستجو جهت دریافت هر صفحه وب باید درخواستی به خدمات دهنده آن صفحه بفرستد متن درخواست هم باید آدرس صفحه مورد بررسی باشد. در زیر لیست مشکلات و راه حل‌های این قسمت را بررسی می‌کنیم.

- جمع آوری آدرس صفحات وب : به چه صورت آدرس یک صفحه وب را بدست آورد و متوجه شد آیا این آدرس معتبر است. (حاوی اطلاعات است یا ارجاع به یک صفحه تبلیغاتی یا حتی ویروسی باشد )
- راه حل : برای این منظور ابتدا لیستی از دامنه‌های مشخص را به موتور جستجو نمایش می‌دهیم (در حافظه آن ذخیره می‌کنیم). موتور جستجو بعد از استخراج آدرس از درون این صفحات و بررسی اطلاعات آنها توسط قسمت تجزیه کننده خودش به استخراج آدرس‌های جدید می‌پردازد. برای تشخیص عدم مشکل و تبلیغاتی بودن یک صفحه وب هم می‌شود به بررسی سابقه دامنه‌ای که آدرس از آن آمده است پرداخت به این صورت که اگر صفحات داخل این دامنه فاقد محتوای مفید درباره پرس‌وجوهای مورد بررسی موتور جستجو باشند این صفحات از لیست صفحات وب حذف شوند و به همین صورت ارزش دامنه این صفحات و صفحاتی از دامنه‌های دیگر که به این دامنه ارجاع داده‌اند (به این دامنه و محتوای داخلش لینک پیوندی داده‌اند) هم کاهش می‌یابد. در زیر چند مزیت کسب آدرس‌های معتبر را بررسی می‌کنیم [5].

- جلوگیری از صفحات تبلیغاتی : بعضی صفحات وب فقط جنبه تبلیغاتی دارند و هیچگونه اطلاعات مفیدی برای موتور جستجو ندارند و از طرفی به علت ورود اطلاعات نا صحیح به زیان موتور جستجو هستند به دلیل افزایش پردازش داده‌های که فاقد محتوای مفید هستند.
- کاهش فضای جستجو : با حذف صفحات وب تکراری بار پردازشی موتور جستجو کاهش می‌یابد.

- چگونگی تشخیص اولویت یک آدرس : به دلیل رشد نمایی آدرس صفحات وب، چه راهکاری برای تشخیص اهمیت آدرس صفحه وب در مورد یک پرس‌وجوی خاص در موتور جستجو است.

**راه حل :** سیستم رتبه بندی مهمترین قسمت یک موتور جستجو است و تفاوت اصلی در موتورهای جستجو در همین بخش نهفته است. سیستم رتبه بندی برای پاسخگویی به کاربران پارامترهای زیادی را در نظر می‌گیرد تا بهترین پاسخ را ارائه کند. روش هایی رتبه بندی صفحات وب به دو نوع زیر تقسیم می‌شوند [1].

- مبتنی بر ساختار<sup>۱</sup> : مانده الگوریتم های page ranking algorithm، الگوریتم کلونی
- مبتنی بر محتوا<sup>۲</sup>.

در روش های مبتنی بر ساختار به ساختار و مقادیر تگ های یک صفحه وب امتیاز می‌دهیم. به عنوان مثال اگر یک صفحه وب دارای تگ‌های تصویر زیادی است می‌شود این احتمال است که این صفحه وب دارای اطلاعات مفیدی درباره تصویر باشد (یا حتی شامل اطلاعات غیر مفید و تبلیغاتی می‌باشد). در روش‌هایی مبتنی بر محتوا بیشتر تمرکز روی محتوا و صحت داده ها می‌باشد. در این پروژه ما از ترکیب این دو روش استفاده کرده‌ایم که رتبه صفحه وب ترکیبی از غالب صفحه وب و متن و تصاویر مشابه با پرس‌وجو های داخل موتور جستجو بدست می‌آید در بخش‌های آیند به طور کامل در این مورد بحث می‌کنیم.

۱-۵-۲-۳ روش آقای ساجین گوپتا<sup>۳</sup>

---

<sup>۱</sup> structure-based

<sup>۲</sup> content-based

<sup>۳</sup> sachin gupta

آقای گوپتا با استفاده از فرمول ۲-۳ به امتیاز دهی به آدرس‌های وب توسط موتور جستجو پرداختند [22].

$$PR(u) = (1 - d) + d * \sum_{v \in B(u)} (PR(v) * W^{in}(v, u) * W^{out}(v, u)) \quad (\text{فرمول ۲-۳})$$

حال به معرفی پارامترهای این فرمول می‌پردازیم.

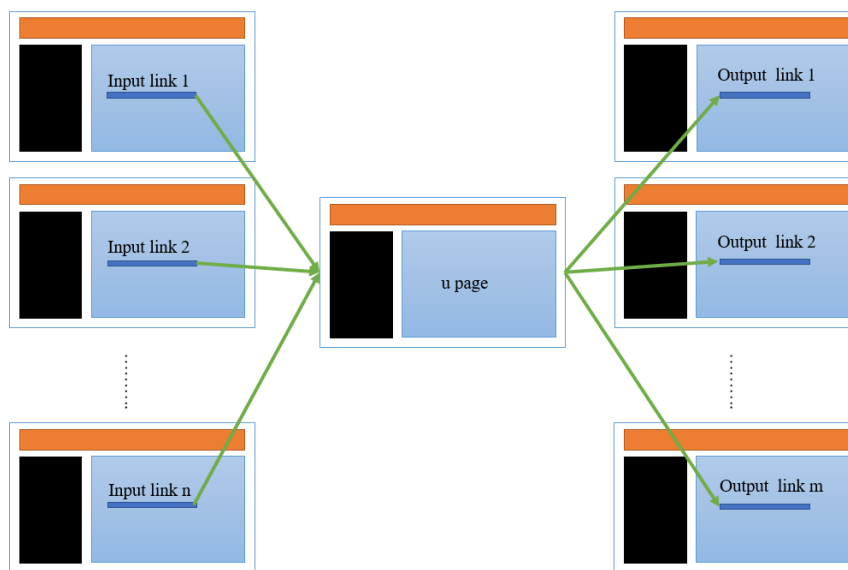
- $W^{in}$ : لینک‌های ورودی به صفحه وب  $u$  (تعداد تمام صفحاتی که به صفحه  $u$  ارجاع دارند و از داخل آن صفحات می‌شود وارد این صفحه شد) و تمام لینک‌های ورودی ارجاع داده شده از لینک خروجی  $v$  (تمام لینک‌های که صفحه  $v$  به صفحه  $u$  یا پدران و فرزندان صفحه  $u$  متصل شده‌اند).
- $W^{out}$ : لینک‌های خروجی از صفحه  $u$  و تمام لینک‌های خروجی از همه صفحات ارجاع داده شده از لینک ورودی  $v$ .
- $d$ : یا فاکتور میرایی<sup>۱</sup>: که یک ثابت در این فرمول می‌باشد و آقای گوپتا پیشنهاد داده‌اند از عدد ۰٫۸۵.
- در آن استفاده شود. البته در مقاله اشاره به آموزش و تغییر این فاکتور به مرور زمان شده است.
- $PR^r$ : تابع محاسبه امتیاز جهت محاسبه امتیاز لینک  $u$  استفاده می‌شود.
- صفحات وب: در اینجا  $u$  همان لینک صفحه وب جاری است و  $v$  مجموع لینک‌های هستند که به نحوی به صفحه  $u$  متصل هستند مثلاً یا از صفحه  $u$  به آنها ارجاع داریم که به آن لینک ورودی گویند یا صفحه  $v$  به صفحه  $u$  یک ارجاع داده است که به آن لینک خروجی گویند.
- تابع  $B(u)$ : این تابع از مجموعه لینک‌های موجود لینک‌های را انتخاب می‌کند که یا به صفحه  $u$  ارجاع داده‌اند یا صفحه  $u$  به آنها ارجاع داده است. در حقیقت ما باید لیست لینک‌های ورودی و خروجی هر صفحه وب را بیابیم. جهت انجام این کار ما باید تمام صفحات وب و لینک‌های متصل

<sup>۱</sup> damping factor

<sup>r</sup> process ranking

به آنها را بیابیم و سپس به محاسبه مقادیر این اتصالات بپردازیم. به این منظور از یک جدولی به نام جدول رابطه<sup>۱</sup> در پایگاه داده (در فصل چهارم به طور کامل درباره پایگاه داده<sup>۲</sup> استفاده شده در موتور جستجو صحبت می کنیم) استفاده کرده ایم که وظیفه آن ذخیره لیست ارجاعات به صفحات وب می باشد.

در شکل ۱۶-۳ به مفاهیم ارجاع ورودی و خروجی صفحات وب اشاره می کنیم.



شکل ۱۶-۳ لینک های ورودی و خروجی صفحه وب [1]

## ۲-۵-۳ روش موتور جستجوی گوگل

البته این موتور جستجوی گوگل به مرور زمان امتیاز دهی خود را به صفحات وب را تغییر می دهد. ولی یکی از روش هایی امتیاز دهی به صفحات وب این موتور جستجو مطابق فرمول ۳-۳ می باشد [1].

$$PR(u) = (1-d) + d \sum_{i=1}^n \frac{PR(T^i)}{C(T^i)} \quad (\text{فرمول ۳-۳})$$

<sup>۱</sup> relation table

<sup>۲</sup> data base

حال به معرفی پارامترهای این فرمول می‌پردازیم.

- $d$  یا فاکتور میرایی: که یک ثابت در این فرمول می‌باشد و از عدد  $0,85$  در آن استفاده شده است.
- $T^i$ : مجموعه صفحاتی که به صفحه  $u$  ارجاع داده اند.
- $C(T^i)$ : تعداد لینک های که از صفحه  $T^i$  خارج شده اند (منظور تعداد صفحات وبی است که صفحه  $T^i$  به آنها ارجاع داده است).
- $PR$ : جهت محاسبه امتیاز لینک از این تابع استفاده می‌شود.
- تابع  $B(u)$ : این تابع از مجموعه لینک های موجود لینک های را انتخاب می‌کند که یا به صفحه  $u$  لینک ورودی دارند یا لینک خروجی.

یکی از مزایای روش موتور جستجوی گوگل تغییر اولویت یک صفحه وب با افزایش اهمیت صفحات متصل به آن می‌باشد. و از معایب این روش هم می‌شود به این مورد اشاره کرد که بین صفحاتی که واقعاً اطلاعات خوب و مهمی دارند و صفحاتی که فاقد اطلاعات با ارزشند و فقط به منظور افزایش بازدید از یکسری کلمه کلیدی استفاده کرده اند تفاوتی نمی‌گذارد.

### ۳-۵-۲ روش **page ranking developed algorithm**

این روش گونه توسعه یافته روش PageRank است و نحوه امتیاز دهی به صفحات وب در این روش به فرم زیر است [5].

- اهمیت دادن به لینک های ورودی و خروجی صفحه وب جاری
- امتیاز رتبه بندی مبتنی بر محبوبیت صفحات توزیع شده. ( این محبوبیت یا از طریق بازدید کاربران و کلیک روی لینک آن در نتایج صفحات وب می‌باشد یا به فرم افزایش امتیاز صفحه در موتور جستجو

درباره یک پرس‌وجوی خاص) روشهای مختلفی برای پیاده سازی این الگوریتم وجود دارد مثل شبکه عصبی و مارکوف.

**نوآوری :** در این پروژه برای هر دامنه یک امتیاز در نظر گرفته می‌شود. در ابتدا که صفحات وب فاقد امتیاز هستند طبق الگوریتم page ranking می‌باشد ولی به مرور زمان و بازدید از صفحات امتیاز صفحات با توجه به نوع پرس‌وجوی پرسیده شده از آن تغییر می‌کند. هر صفحه امتیازی از یک پرس‌وجوی خاص می‌گیرد (همان پرس‌وجوهای داخل موتور جستجو). پس برای هر آدرس یک جدول از امتیازات پرس‌وجو را ذخیره می‌کنیم. در مورد این روش به فرم کامل در فصل چهارم صحبت می‌کنیم.

### ۳-۳ مدیریت پایگاه داده

منظور از این بخش اطلاعاتیست که بعد از دریافت توسط موتور جستجو جهت استفاده‌های بعدی در موتور جستجو ذخیره می‌شوند. حذف اطلاعات و ویژگی‌های تکراری یا بی ارزش در پایگاه داده جهت افزودن اطلاعات جدید. به علت محدود بودن پایگاه داده قسمت مدیریت پایگاه داده باید اطلاعات تکراری یا فاقد ارزش را حذف کند. در موتور جستجو باید لیستی از سوالات و پرس‌وجوهای پرسیده شده از موتور جستجو ذخیره شوند.

**چالش :** چه روشی برای اندیس گذاری نتایج باید استفاده شود. در حالت خلاصه منظور از اندیس کردن نتایج نوعی نگاشت نتایج به پرس‌وجوی موجود در موتور جستجو می‌باشد. و این اندیس ها چگونه در پایگاه داده ذخیره شوند.

### نوآوری :

باید نتایج بدست آمده در صفحه وب به نحوی اندیس گذاری شوند که با صرف کمترین هزینه و وقت



نتایج آماده نمایش شوند. از طرفی یک نتیجه ممکن است با چندین موضوع مرتبط باشد به عنوان مثال فرض کنید یک تصویر با موضوعیت برج در پایگاه داده داریم و دو پرسوجوی "نمایه زیبا از برج میلاد" و "برج آزادی" از موتور جستجو پرسیده شود. پس در هر دو این پرسوجوها باید در لیست نتایج نمایش داده شود. این کار بعد از مراحل یادگیری موتور جستجو و یا انتخاب کاربران بازدید کننده از سایت تغییر می‌کند. در این جا با عبور پرسوجوها از شبکه عصبی موتور جستجو درصد شباهت پرسوجوی جدید با لیست پرسوجوهای داخل موتور جستجو مشخص می‌شود. باید برای هر پرسوجو یک اندیس داشته باشیم. و از

طرفی بعد از مشخص شدن اهمیت پرسوجو در یک صفحه وب به دادن امتیاز به تصاویر و متون نزدیک تصاویر در صفحه وب می‌پردازیم.

#### ۴-۳ مدیریت پرسوجو

ایجاد و طبقه بندی پرسوجوها در یک موتور جستجو از اهمیت بالایی برخوردار می‌باشد. ممکن است کاربر سوالی را از موتور جستجو بپرسد که در لیست طبقات پرسوجوی موتور جستجو وجود نداشته باشد. ولی متن پرسوجوی جدید با یک یا چندتا از پرسوجوهای ذخیره شده در موتور جستجو شباهت داشته باشد. در این مرحله ما ابتدا پرسوجوی جدید را بعد از مرحله پیش پردازش و حذف کلمات تکراری و علامت‌های ویرایشی به یک بردار تبدیل کرده و به شبکه عصبی موتور جستجو می‌دهیم و خروجی این قسمت میزان شباهت پرسوجوی جدید با پرسوجوی جدید با پرسوجوهای داخل موتور جستجو است.

#### کاربردها

- بررسی مشکلات نحوی پرسوجو:

ممکن است شما در پرس‌وجوی خود کلمه ای را به اشتباه وارد کرده باشید مثلاً در جمله show best dug کلمه dog را به اشتباه dug وارد کرده‌اید. که وظیفه موتور جستجو تشخیص این اشکالات نحوی، لغوی و ارسال مجدد صحیح پرس‌وجو می‌باشد. که به راحتی در شبکه عصبی میزان شباهت dog با dug مشخص می‌گردد و می‌شود نتیجه را به عنوان کلمه صحیح به کاربر نمایش داد.

- نمایش لیست پرس‌وجوهای مشابه و پر بازدید :

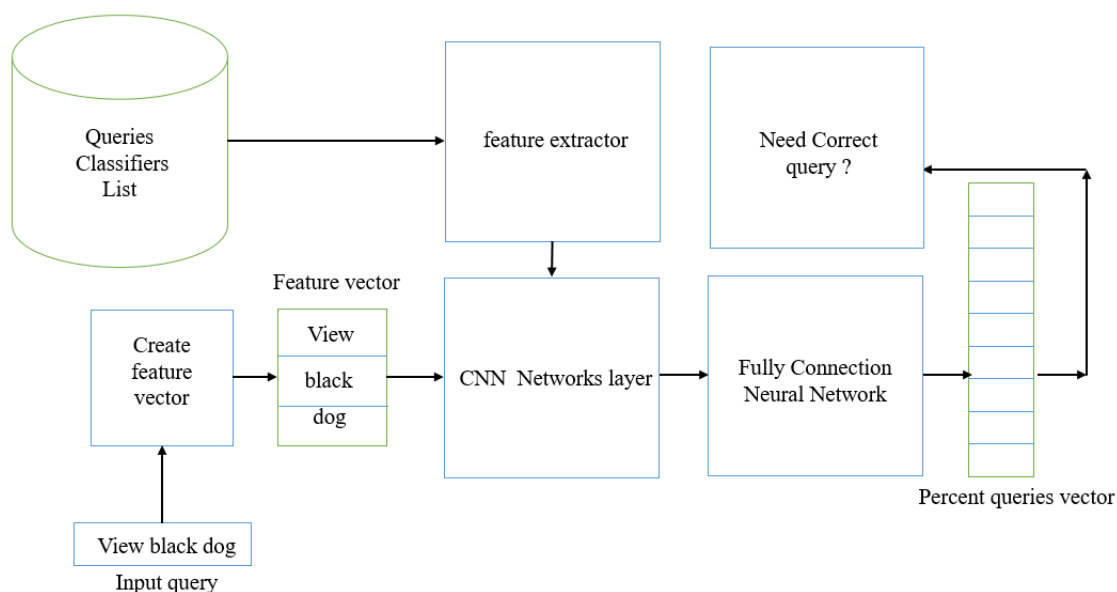
هدف از این قسمت نمایش لیست چند پرس‌وجوی مشابه با پرس‌وجوی خود که نتایج خوبی داشتند است.

- چگونگی بررسی شباهت سوال پرسیده شده از موتور جستجو با پرس‌وجوهای موجود در هر طبقه و دادن یک

امتیاز به عنوان امتیاز شباهت به آن و یا حتی ایجاد یک طبقه بند جدید.

## نوآوری :

در این روش و با استفاده از تکنیک یادگیری عمیق به بررسی شباهت متن پرس‌وجوی داده شده به موتور جستجو با پرس‌وجوهای ذخیره شده در موتور جستجو پرداخته و در نهایت خروجی این شبکه عصبی لیست احتمال شباهت متن پرس‌وجوی ورودی با طبقات مختلف پرس‌وجو می‌باشد که در نهایت پرس‌وجو یا به یک یا چند طبقه خاص بر حسب احتمال شباهت متصل می‌گردد و یا یک طبقه جدید ایجاد می‌گیرد. همانطور که در شکل ۳-۱۷ مشاهده می‌شود.



شکل ۳-۱۷ شبکه CNN برای محاسبه طبقه بندی پرسوجوی جدید [21].

همانطور که در شکل ۳-۱۷ مشخص شده، ابتدا ما بردار ویژگی پرسوجوی ورودی را به دست آورده و آنرا به شبکه عصبی چند لایه کانولوشن می‌دهیم و سپس نتیجه بدست آمده به شبکه عصبی جهت طبقه بندی پرسوجو

داده می‌شود که در نهایت یک بردار حاوی احتمال شباهت پرسوجوی ما از سایر پرسوجو ها بدست می‌آید. اندازه این بردار به اندازه طبقات پرسوجوی موجود در پایگاه داده می‌باشد. در صورت شباهت پرسوجوی جدید با یکی یا چند نتیجه در بردار ویژگی این پرسوجو به آنها اضافه می‌شود و اگر پرسوجوی امتیاز خوبی برای شباهت با این پرسوجو را نداشته باشد یک طبقه جدید ایجاد می‌گردد و سر انجام بردارهای ویژگی هم برای مراحل بعدی به روز رسانی می‌گردند. در مرحله آخر هم می‌شود با نتایج بدست آمده (بردار احتمال خروجی شباهت متن با پرسوجو های داخل موتور جستجو) و درصد شباهت پرسوجو با لیست نتایج به شباهت پرسوجوی ورودی با پرسوجوهای موجود رسید و همینطور به بررسی مشکلات نحوی و لغوی پرسوجو پرداخت البته می‌شود قبل از ایجاد بردار ویژگی به بررسی مشکلات لغوی هم پرداخت (انجام یک پیش پردازش).



## فصل چهارم:

معرفی الگوریتم موتور جستجوی

وب

## ۴-۱ مقدمه

در فصل قبل ما درباره چالش ها و راه حل های پیش روی یک موتور جستجو بحث کردیم حالا و در ادامه به بررسی مراحل پیاده سازی موتور جستجوی با محوریت تصویر به صورت مرحله به مرحله پرداخته و هر قسمت را معرفی می کنیم.

## ۴-۲ دریافت لیست صفحات وب

در این مرحله مدیریت صف آدرس درخواستی به پایگاه داده فرستاده و لیست تمام آدرس های وب را از آن می گیرد. برای قرار دادن آدرس ها در لیست فوق باید ابتدا اولویت هر آدرس مشخص گردد. بعضی از این آدرس ها جدید بوده و تا بحال پردازش نشده اند و فاقد اولویت هستند پس چگونه اولویت این صفحات را مشخص کنیم. در صورتی که صفحه وب از قبل در پایگاه داده مقدار اولویت نداشته باشد (به تازگی این آدرس بدست آمده است) باید برای بدست آوردن اولویت آن مراحل زیر را انجام داد. فرض کنید صفحه جاری که فاقد امتیاز است با نام unPriorityP شناخته می شود.

- اگر آدرس unPriorityP جز یک دامنه خاص باشد که در پایگاه داده موجود است می توان لیست کلیه صفحاتی که با این صفحه ارتباط دارند را یافت و میانگین این امتیازات را به عنوان امتیاز صفحه unPriorityP قرارداد [1].

در صورتی که آدرس unPriorityP جز یک دامنه خاص در پایگاه داده نیست ابتدا کلیه صفحاتی که دارای اولویت هستند را گروه بندی می کنیم. گروه بندی به این صورت می باشد که کلیه صفحاتی که متعلق به یک دامنه اند را در یک گروه قرار می دهیم به عنوان مثال دو صفحه [www.hardwarenews.com/cpu](http://www.hardwarenews.com/cpu) و [www.hardwarenews.com/ram](http://www.hardwarenews.com/ram) هر دو به یک دامنه اصلی به نام [www.hardwarenews.com](http://www.hardwarenews.com) اشاره دارند پس هر دو در یک گروه قرار می گیرند مانند شکل ۴-۱.

- حال در هر گروه به دنبال صفحاتی می‌گردیم که با این صفحه ارتباط دارند یعنی به این صفحه ارجاع داده‌اند. سپس برای هر صفحه در گروه به محاسبه امتیازش می‌پردازیم که به فرم جمع کردن تمام اتصالات به همراه وزنشان می‌باشد.
- در مرحله آخر با میانگین‌گیری امتیاز هر صفحه در یک گروه خاص امتیاز آن گروه خاص به صفحه وب مشخص می‌گردد. مطابق فرمول ۱-۴.

$$PR(\text{group}_i) = (\sum_{P \in \text{group}_i} \text{ScoreLink}(P) * P(\text{priority})) / \text{count}(\text{group}_i) \quad (\text{فرمول } 1-4)$$



شکل ۱-۴ گروه بندی دامنه ها

حال به معرفی پارامترهای فرمول می‌پردازیم :

- PR : تابع محاسبه امتیاز یک گروه خاص (  $group_i$  ) به صفحه بدون امتیاز می‌باشد.
- P : صفحه که متعلق به یک گروه خاص می‌باشند. (دامنه یکسانی با سر گروه دارد)
- ScoreLink : وظیفه این تابع محاسبه امتیاز ارتباطی بین صفحه P با صفحه unPriorityP می‌باشد. یعنی محاسبه تعداد لینک ارتباطی صفحه P و صفحه unPriorityP . در صورت عدم ارتباط این دو صفحه نتیجه این تابع صفر می‌باشد.
- priority : همان مقدار ارزش صفحه P است.
- count : این تابع وظیفه محاسبه تعداد صفحاتی است که در یک دامنه‌اند و با صفحه unPriorityP ارتباط دارند.

• حال که امتیاز هر گروه به صفحه unPriorityP مشخص گردید. باید با یک روش به محاسبه امتیاز اولیه صفحه وب unPriorityP پردازیم.

- روش میانگین گیری : در این روش اهمیت و اولویت صفحه با میانگین گیری از صفحات مرتبط به این صفحه بدست می‌آید (البته وزن اتصال هم در این مورد اهمیت دارد) [1].
- روش تصادفی<sup>۱</sup> تکی : در این روش یکی از دامنه‌ها را به تصادف انتخاب کرده و ارزش آنرا به این صفحه می‌دهیم. البته می‌شود از روش‌هایی مانند چرخ رولت<sup>۲</sup> هم استفاده کرد ( اهمیت هر دامنه بر حسب امتیاز ارتباطش با صفحه unPriorityP می‌باشد ).

روش تصادفی چند گانه : این روش تقریباً شبیه روش تصادفی تکی است با این تفاوت که ابتدا ما n گروه از دامنه‌ها را به فرم تصادفی انتخاب می‌کنیم و با میانگین‌گیری از آنها نتیجه را به عنوان امتیاز به صفحه می‌دهیم. در این پروژه از روش تصادفی چندگانه استفاده می‌کنیم.

---

<sup>۱</sup> stochastic method

<sup>۲</sup> roulette wheel



حالا تمام صفحات امتیاز دارند و به ترتیب اولویت به صف آدرس اضافه می گردند.

حال که ارزش دامنه همه صفحات مشخص شد بهتر است به روش پیشنهادی که برای دادن امتیاز به یک دامنه خاص در این پروژه استفاده شده است نیز اشاره داشته باشیم.

#### ۱-۲-۴ محاسبه امتیاز دامنه

حالا به محاسبه امتیاز یک دامنه خاص و درباره یک موضوع پرس و جوی خاص می پردازیم. در این قسمت با استفاده از یادگیری ماشین و رگرسیون به آموزش و محاسبه ارزش یک دامنه می پردازیم.

در این روش ابتدا به معرفی چند پارامتر ورودی می پردازیم.

- $\text{domain}[i,j]$  : منظور محاسبه امتیاز دامنه  $i$  ام و پرس و جوی  $j$  ام است.
- $\text{imagesquery}[i,j]$  : این پارامتر تعداد تصاویر دیده شده در دامنه  $i$  ام که با پرس و جوی  $j$  ام شباهت دارد را نمایش می دهد.
- $\text{seeimages}[i]$  : این پارامتر تعداد تصاویر دیده شده در دامنه  $i$  ام را نمایش می دهد.
- $\text{textquery}[i,j]$  : این پارامتر تعداد متون دیده شده در دامنه  $i$  ام که با پرس و جوی  $j$  ام شباهت دارد را نمایش می دهد.
- $\text{seetexts}[i]$  : این پارامتر تعداد متون دیده شده در دامنه  $i$  ام و نزدیک تصاویر را نمایش می دهد.
- $\text{successpages}[i]$  : این پارامتر تعداد صفحات وبی را در دامنه  $i$  ام که به طور صحیح از خدمات دهنده آن دریافت شده است را نمایش می دهد.
- $\text{allpages}[i]$  : این پارامتر تعداد صفحات وبی را در دامنه  $i$  ام قرار دارد را نمایش می دهد.

حال با استفاده از این پارامترها به ایجاد ۳ پارامتر ورودی مطابق پارامترهای فرمول ۴-۲ می پردازیم.

$$X_0 = 1 \quad X_1 = \frac{\text{imagesquery}[i,j]}{\text{seeimages}[i]} \quad X_2 = \frac{\text{textquery}[i,j]}{\text{seetexts}[i]} \quad X_3 = \frac{\text{successpages}[i,j]}{\text{allpages}[i]} \quad (\text{فرمول ۲-۴})$$

در مرحله بعدی و با استفاده از وزن‌های ذخیره شده در پایگاه داده درباره دامنه نام و پرسوجوی نام به محاسبه امتیاز دامنه می‌پردازیم. مطابق با فرمول ۳-۴.

$$\text{domain}[i,j] = X_0 * W_0 + X_1 * W_1 + X_2 * W_2 + X_3 * W_3 \quad (\text{فرمول ۳-۴})$$

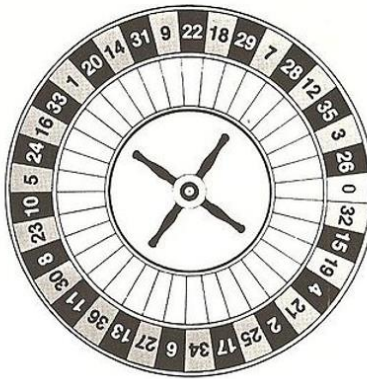
حال به محاسبه احتمال لیست دامنه ها می‌پردازیم در اینجا با استفاده از روش softmax به محاسبه این احتمال مطابق فرمول ۴-۴ می‌پردازیم. در اینجا ما فرض کرده‌ایم تعداد دامنه ها D تا است.

$$\text{domainrank}[i,j] = \frac{e^{\text{domain}[i,j]}}{\sum_{d=1}^D e^{\text{domain}[d,j]}} \quad (\text{فرمول ۴-۴})$$

در فرمول ۴-۴ منظور از domainrank[i,j] محاسبه احتمال انتخاب دامنه نام درباره پرسوجوی نام می‌باشد. حالا با استفاده از چرخ رولت به انتخاب یکی از این دامنه ها می‌پردازیم مطابق با شکل ۲-۴ [4].

$$\text{Probs} = [ \text{domainrank}[1,j] , \text{domainrank}[2,j] , \dots , \text{domainrank}[D,j] ]$$

roulette wheel



شکل ۲-۴ انتخاب دامنه درمورد یک پرسوجوی خاص [4]

بعد از انتخاب دامنه و صفحه وب به دریافت اطلاعات می‌پردازیم. و در نهایت بعد از دریافت اطلاعات به محاسبه نتیجه صفحه وب و به روز رسانی وزن ها و مقدار امتیاز دامنه در یک پرس‌وجوی خاص می‌پردازیم. در این پروژه وزنها و مقدار امتیاز قبلی در پایگاه داده ذخیره می‌گردد.

### ۳-۴ مازول دریافت کننده صفحه وب

بعد از پر شدن صف آدرس و مقدار دهی اولویت به کلیه صفحات باید مازول دریافت کننده را آگاه کرد تا اولین آدرس را از صف اولویت دریافت کرده و با ارسال درخواستی به خدمات دهنده آدرس به دریافت صفحه وب بپردازد. بعد از دریافت صفحه توسط مازول دریافت کننده و قبل از ارسال به مرحله بعد یکسری پیش پردازش مطابق مراحل زیر صورت می‌گیرد.

- **بررسی صحت و درستی صفحه وب :** همانطور که در فصل قبل توضیح داده شد آیا صفحه وب دریافتی صحیح است و به طور کامل از طرف خدمات دهنده<sup>۱</sup> به موتور جستجو ارسال گردیده است. در صورت بروز خطا یکسری کد مخصوص ایجاد می‌شود که موتور جستجو متوجه صحت و درستی نتایج یا خرابی آنها می‌گردد و این نتایج در امتیاز صفحه وب و دامنه آن تاثیر گذار می‌باشد.
- **تغییر امتیاز صفحه وب :** در صورت اشکال در صفحه وب می‌شود امتیاز صفحه وب (و حتی دامنه آن) را کاهش داد. (دلیل عدم حذف آنی صفحه بجای کاهش امتیاز این است که ممکن است صفحه وب موقتا غیر فعال شده باشد و در آینده دوباره روی خدمات دهنده قرار بگیرد). البته در اینجا کلیه آدرسها در پایگاه داده هم که با این صفحه ارتباط لینکی دارند با کاهش امتیاز مواجه می‌شوند.
- **بررسی قالب صحیح صفحه وب :** آیا قالب صفحه با استاندارد صفحات وب HTML سازگار است یا مشکلی برای آن به وجود آمده است. در اینجا مازول دریافت کننده قالب صفحه وب را با استاندارد صفحه وب مقایسه

می‌کند و در صورت عدم مشکل این محتوا را برای قسمت تجزیه کننده می‌فرستد و در غیر این صورت دوباره به کاهش امتیاز صفحه وب و همینطور صفحات وب مرتبط با آن می‌پردازد (صفحاتی که به صفحه وب دریافتی ارجاع داده‌اند).

## ۴-۴ قسمت تجزیه کننده


در مرحله قبل ما صفحه وب را دریافت کرده و صحت آنرا مورد بررسی قرار دادیم. وظیفه اصلی این قسمت تجزیه کردن صفحه وب و استخراج داده‌های مورد نیاز برای پردازش اطلاعات می‌باشد و همینطور بررسی آنها از قبیل متون، تصاویر و لینک‌های ارتباطی.

۴-۴-۱ استخراج تصاویر و متن از صفحه وب : در این مرحله ما به دریافت تصاویر در صفحه وب و لیست متون نزدیک به تصویر می‌پردازیم. منظور از متون نزدیک به تصویر لیست متونی است که در یکی از سطوح

### صفحه وب

روزی روزگاری: فولکلور رویا... | نقد و بررسی بازی Folklore

در خصوص بازی‌های خاص و مخاطبان آن‌ها و این‌که صرفاً خاص بودن یک بازی دلیلی قطعی بر خوب یا بد بودن آن نیست بازها و بارها صحبت کرده‌ایم و آن را از زوایای مختلف بررسی کرده‌ایم. به واقع در این صنعت عظیم بازی‌های رایانه‌ای و تعداد بیشمار بازی خیلی تکراری که در آن وجود دارند. باز هم برخی از بازی‌ها هستند که وقتی به هر ۸ نسل بازی‌های رایانه‌ای نگاه می‌کنی...




پرونده: تاریخچه...

---

حماسه‌ای از جنس سرب و پلسماسا | نقد و بررسی Halo: Reach

بدون دره‌ای شک یکی از بهترین عناوینی که باعث شد کنسول مایکروسافت در مقابل رقابای کنسولی خود حرفی برای گفتن داشته باشد. مجموعه بازی‌های هیلو است. البته عرضه هیلو‌ها هم داستان خود را دارند و از عنوانی که قرار بود در ابتدا برای رایانه‌های مک عرضه شود به طور ناگهانی تغییر جهت داد و انحصاری مایکروسافت شد. البته بهترین دلیلی که باعث شد هیلو به این جایی...



پرونده: تاریخچه...

### درخت HTML

```

<div class="PostBlock">
  <div class="PostBlock">
    <div class="PostHeadLine">
      <h3>
        <span class="TabBox PostTab">
          <p>در خصوص بازی‌های خاص و مخاطبان آن‌ها و این‌که صرفاً خاص بودن یک بازی دلیلی قطعی بر خوب یا بد بودن آن نیست بازها و بارها صحبت کرده‌ایم و آن را از زوایای مختلف بررسی کرده‌ایم. به واقع در این صنعت عظیم بازی‌های رایانه‌ای و تعداد بیشمار بازی خیلی تکراری که در آن وجود دارند. باز هم برخی از بازی‌ها هستند که وقتی به هر ۸ نسل بازی‌های رایانه‌ای نگاه می‌کنی...</p>
        </span>
      </h3>
    </div>
    <div class="PostInfo clear">
      <div class="right">
        <div class="left">
          <span>
            <a href="http://gamefa.com/403532/4d8b14d9988d8b24db48c-4d8b14d9988d8b24db48f4d8a76d8b14db48c-4d9414d9988d9484da9d984d9988d8b1-4d8b14d9988d8b48c4d8a7-4d94564d982d88af-4d9458-4d8a8d8b14d8b14d8b34db48c/" title="روزی روزگاری: فولکلور رویا... | نقد و بررسی بازی Folklore">
              <img alt="Folklore game cover" data-bbox="785 595 870 635"/>
            </a>... ادامه مطلب...</span>
          </div>
        </div>
      </div>
    </div>
  </div>
  <div class="PostBlock">
    <div class="PostBlock">
      <div class="PostHeadLine">
        <h3>
          <span class="TabBox PostTab">
            <p>بدون دره‌ای شک یکی از بهترین عناوینی که باعث شد کنسول مایکروسافت در مقابل رقابای کنسولی خود حرفی برای گفتن داشته باشد. مجموعه بازی‌های هیلو است. البته عرضه هیلو‌ها هم داستان خود را دارند و از عنوانی که قرار بود در ابتدا برای رایانه‌های مک عرضه شود به طور ناگهانی تغییر جهت داد و انحصاری مایکروسافت شد. البته بهترین دلیلی که باعث شد هیلو به این جایی...</p>
          </span>
        </h3>
      </div>
      <div class="PostInfo clear">
        <div class="right">
          <div class="left">
            <span>
              <a href="http://gamefa.com/404163/4d8fad94855d8a74d8b34d9487-4d8a74db48c-4d8b2-4d8a8d9484da9d984d9988d8b1-4d8b14d9988d8b48c4d8a7-4d94564d982d88af-4d9458-4d8a8d8b14d8b14d8b34db48c-ha/" title="حماسه‌ای از جنس سرب و پلسماسا | نقد و بررسی Halo: Reach">
                <img alt="Halo: Reach game cover" data-bbox="785 660 870 700"/>
              </a>... ادامه مطلب...</span>
            </div>
          </div>
        </div>
      </div>
    </div>
  </div>
</div>

```

شکل ۴-۳ متون نزدیک تصویر

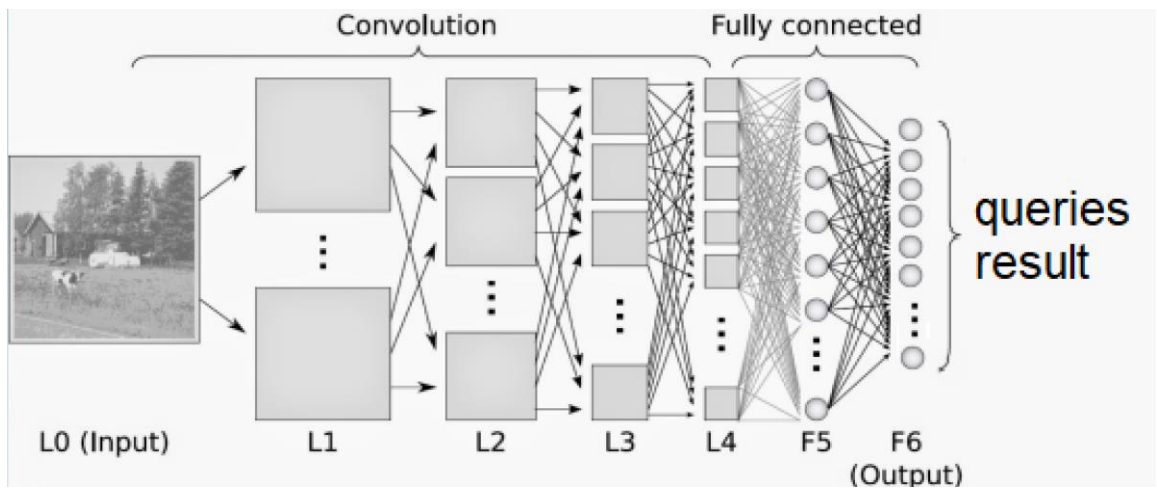
صفحه وب (درخت HTML) و نزدیک به تصویر قرار دارند. این متون بر حسب یک روش خاص از صفحه وب استخراج شده اند (در مراحل قبلی پردازش صفحه وب با روش هایی یادگیری عمیق صفحه وب متوجه اولویت متون نزدیک به صفحه شده است به عنوان مثال در یک صفحه وب ممکن است متوجه شده باشیم که متن موجود در داخل سر صفحه وب<sup>۱</sup> و یا خود تصویر صفحه از اهمیت بالای برخوردارند ولی در یک صفحه وب دیگر ممکن است متن با اهمیت نزدیک صفحه در پاراگراف بالای آن باشد و ... ) یا به فرم تصادفی ( اگر اولین بار است یک سری متون را از داخل آن استخراج می کنیم ) شاید راهکار منطقی تر استفاده از الگوی صفحه ای است که بیشترین شباهت را با صفحه فوق دارد (معمولا صفحات دیگری از دامنه که با صفحه فوق ارتباط دارند ). از طرفی هر صفحه یک الگوی استخراج متون نزدیک به تصویر مخصوص به خود دارد که با استفاده از یادگیری عمیق به مرور زمان یاد می گیرد. در شکل ۳-۴ می توان تصویر و متون نزدیک به آن را مشاهده کرد. بعد از دریافت تصویر و متون داخل صفحه وب در این مرحله ما به ایجاد بردار تصویر و بردار متن پرداخته و هر کدام را به شبکه عصبی کانولوشن خودشان برای پردازش می دهیم.

## ۲-۴-۴ بررسی تصویر و متن داخل صفحه وب

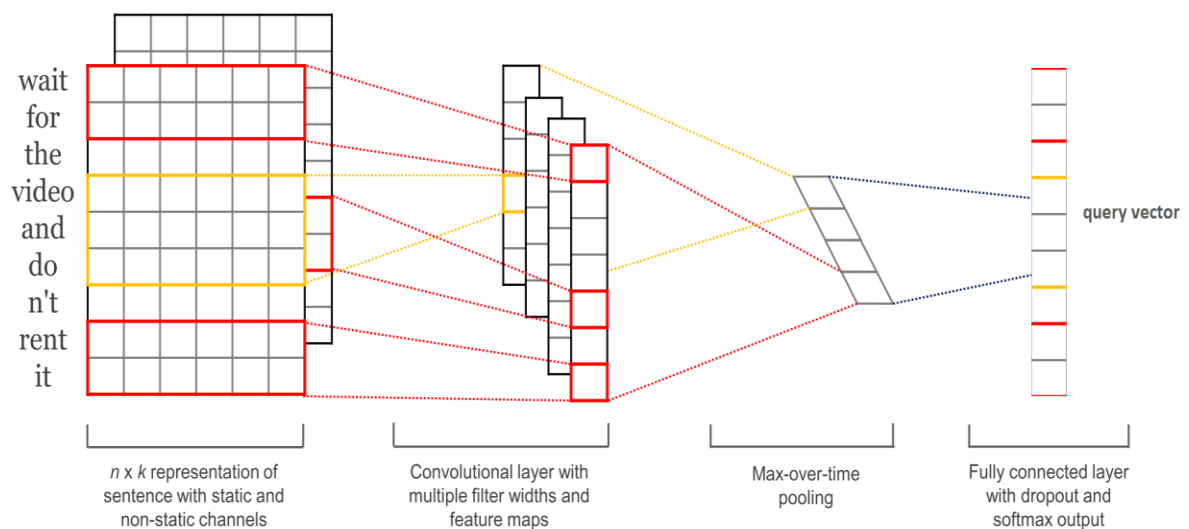
در گام بعدی ما به ایجاد بردار ویژگی تصویر و متون نزدیک تصویر پرداخته و آنها را به شبکه های کانولوشن خودشان مطابق شکل های ۴-۴ و ۵-۴ می دهیم. خروجی برای هر کدام از این شبکه ها یک بردار احتمال<sup>۲</sup> به طول کلیه پرس و جوهای موتور جستجو می باشد. و احتمال شباهت هر پرس و جو با تصویر و متن مشخص می گردد. در این مرحله با توجه به متن استخراج شده و بردار احتمال به روز رسانی مکان تصاویر مهم در صفحه وب جاری می پردازیم. (به عنوان مثال ممکن است به مرور زمان مکان پر اهمیت متن نزدیک یک تصویر از پاراگراف نزدیک آن به برجسب فرزند آن منتقل گردد).

<sup>۱</sup>header

<sup>۲</sup>stochastic vector



شکل ۴-۴ شبکه کانولوشن تصویر [20]



شکل ۵-۴ شبکه کانولوشن متن [21]

حال بر حسب بردار ویژگی بدست آمده از روی تصویر می‌شود به اهمیت گروه تصویر پی برد یک راه حل برای این منظور اتصال تصویر به گروه ها با توجه به وزن بدست آمده است هر چقدر وزن یک پرس‌وجو بالاتر باشد احتمال شباهت تصویر با آن گروه هم به همان نسبت افزایش می‌یابد. حال با بروز رسانی بردارهای ویژگی پرس‌وجوی خاص با استفاده از تصویر و وزن ارتباطی بین آنها می‌پردازیم.

در گام بعدی با استفاده از بردار ویژگی بدست آمده از متن می‌شود به بررسی پرس‌وجو ها پرداخت و برای ایجاد یک پرس‌وجوی جدید یا الحاق پرس‌وجو ها اقدام کرد. پرس‌وجو ها هم همانطور که قبلا گفته شد با هم اتصال دارند با یکسری وزن خاص حالا که بردار متن بدست آمده می‌شود وزن پرس‌وجوها را به روز رسانی کرد و ممکن است پرس‌وجوی جدید ایجاد شده که یا به یک پرس‌وجوی سطح بالاتر متصل می‌شود یا خود تشکیل یک طبقه پرس‌وجوی جدید را می‌دهد.

حال به به روز رسانی امتیاز و اولویت صفحه وب جاری و صفحات مرتبط با آن می‌پردازیم ( تمام صفحات مرتبط با صفحه جاری چه از لحاظ ارتباط لینکی و هم از لحاظ یکسان بودن دامنه).

### ۳-۴-۴ استخراج لینک های اتصالی

دراین مرحله ما به استخراج کلیه لینک‌های اتصالی صفحه وب جاری به صفحات وب دیگر می‌پردازیم. در حقیقت هدف از این کار محاسبه امتیاز صفحه وب جاری (یا به روز رسانی) و صفحات مرتبط با آن می‌باشد. بعد از بدست آوردن لیست لینک ها در چند مرحله زیر به اقداماتی می‌پردازیم.

- **نرمالایز کردن :**

همانطور که در فصل دوم اشاره شد جهت کاهش بار پردازشی موتور جستجو با استفاده از فرایند نرمال سازی باعث کاهش و مرتب سازی آدرس‌های وب می‌شویم. و این کار باعث کاهش ترافیک شبکه نیز خواهد گردید.

- **به روز رسانی اولویت**

در این مرحله هر کدام از این ارجاعات استخراج شده را با صفحات موجود در پایگاه داده مقایسه می‌کنیم. و اگر در پایگاه داده این آدرس وجود داشته باشد بر حسب امتیازات بردار ویژگی صفحه جاری و صفحه بدست آمده اولویت صفحه تغییر می‌کند. دقت شود فرض کنید در لیست آدرس‌ها، آدرسی به

نام newURL داریم بدیهی است تمام صفحاتی در پایگاه داده که به نحوی با newURL ارتباط دارند ( چه از لحاظ داشتن دامنه یکسان، دارا بودن رجاع به صفحات داخل پایگاه داده و یا خود آدرس و همینطور شباهت بالا در پرس و جوهای موتور جستجو ) دچار تغییر امتیاز خواهند شد.

## ۵-۴ قسمت کنترل پایگاه داده<sup>۱</sup>

وظیفه اصلی این قسمت ذخیره و بازیابی اطلاعات در پایگاه داده می باشد. از طرفی به دلیل درخواست های زیاد و مکرر موتور جستجو از یک طرف و تغییر و به روز رسانی های بسیار زیاد موتور جستجو روی پایگاه داده از طرف دیگر باید این قسمت از لحاظ تراکنشی<sup>۲</sup> بسیار قوی و با اطمینان عمل کند. و کوچکترین اشکالی در این امر منجر به گمراه شدن موتور جستجو می گردد به عنوان مثال در هنگام به روز رسانی بردار پرس و جوی یک آدرس وب اگر یک فرایند دیگر اقدام به به روز رسانی همان آدرس بکند اطلاعات با ارزش از بین خواهد رفت. در شکل ۴-۶ می توانید بلوک دیاگرام قسمت کنترل پایگاه داده را ببینید. در ادامه به بررسی قسمت های مختلف این کنترلر می پردازیم.

### • بررسی درخواست<sup>۳</sup>

این قسمت بررسی می کند آیا درخواست پرسیده شده از پایگاه داده معتبر است یا نه. دلیل این امر هم جلوگیری از تزریق سوالات تزریقی<sup>۴</sup> و درخواست های که فاقد مجوز لازم هستند (به عنوان مثال هر ماژول در موتور جستجو مجوز خاص خود را دارد مثلاً ماژول دریافت کننده صفحات وب مجوز تغییر بردار اولویت هیچ صفحه وبی را ندارد)

---

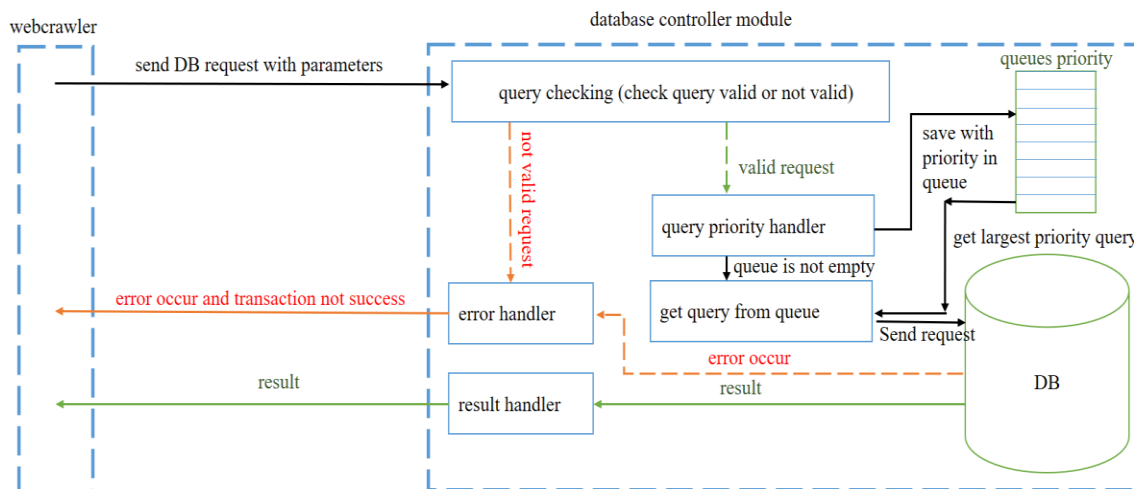
<sup>۱</sup> database controller module

<sup>۲</sup> transaction

<sup>۳</sup> query checking

<sup>۴</sup> sql injection





شکل ۴-۶ بلوک دیاگرام ماژول کنترل پایگاه داده

#### • بررسی اولویت<sup>۱</sup>

این قسمت بعد از تایید درخواست توسط قسمت بررسی قرار دارد. و به بررسی اولویت درخواست می‌پردازد و بر حسب اولویت در صف درخواست ها قرار می‌دهد. این قسمت از اهمیت بسیاری در این ماژول قرار دارد و باید با یک سیستم هوشمند به تشخیص اهمیت درخواست بپردازد. به عنوان مثال اگر دو درخواست در صف داریم که یکی دریافت لیست اولویت کلیه صفحات وب و دیگری تغییر در مقدار اولویت یک صفحه وب باشد بدیهی است که ابتدا بهتر است اولویت آدرس تغییر کرده و سپس کل اولویت آدرس‌ها را به موتور جستجو بفرستیم.

#### • اجرای درخواست<sup>۲</sup>

این قسمت بالاترین درخواست را از صف اولویت استخراج کرده و دستورات را جهت اجرا به پایگاه داده می‌فرستد. و پاسخ را دریافت کرده و نتیجه را به موتور جستجو می‌فرستد البته اگر تراکنش نا موفق بود خطا را برای مدیریت به قسمت کنترل خطا می‌فرستد.

<sup>۱</sup> priority handler

<sup>۲</sup> run handler

- کنترل خطا<sup>۱</sup>

وظیفه این قسمت ارسال خطا به موتور جستجو را بر عهده دارد چنانچه خطای در نتیجه پایگاه داده رخ بدهد و یا درخواست موتور جستجو معتبر نباشد این قسمت وظیفه کنترل این خطا را بر عهده دارد.

#### ۴-۶ پایگاه داده موتور جستجو :

در موتور جستجو برای ذخیره داده‌ها از یک پایگاه داده رابطه‌ای استفاده کرده‌ایم. دلایل استفاده از پایگاه داده رابطه‌ای را می‌توانید در لیست زیر ببینید.

- به علت ارتباط بالای جداول پایگاه داده با یکدیگر (به عنوان مثال جدول تصاویر ارتباط نسبتاً بالای با جدول صفحات وب دارد و هر تصویر به یک صفحه وب متصل می‌باشد) یک پایگاه داده رابطه‌ای بسیار بهتر و سریعتر می‌تواند جداول را با یکدیگر مرتبط کند.

- مدل رابطه‌ای پایگاه داده به طور پیش فرض بسیار سریعتر از مدل‌های دیگر پایگاه داده مانده مدل ساختاری می‌باشد. و به دلیل حجم بالای داده در یک موتور جستجو سرعت پاسخگویی پایگاه داده و ارتباط آن با موتور جستجو از اهمیت بالای برخوردار است.

در این پروژه از یک پایگاه داده رابطه‌ای به نام webCrawlerDB استفاده می‌کنیم. در این پایگاه داده از چندین جدول که با یکدیگر رابطه دارند استفاده گردیده است. حال به معرفی جداول پایگاه داده می‌پردازیم.

۱- **جدول دامنه**<sup>۲</sup>: همانطور که از نام این جدول مشخص می‌باشد وظیفه این جدول ذخیره دامنه‌هایست که موتور جستجو پیدا می‌کند. در جدول ۴-۱ می‌توانید لیست ستون‌های این جدول را ببینید. و در شکل ۴-۷ دیاگرام این جدول را مشاهده می‌کنید.

---

<sup>۱</sup>error handler

<sup>۲</sup>domain

domain		
id	char(36)	PRIMARY KEY UNIQUE
domain	varchar(1000)	
vote	float	

شکل ۴-۷ جدول دامنه

جدول ۴-۱ مشخصات جدول دامنه

ردیف	نام ستون	نوع	توضیحات
۱	id	char(36)	کلید اصلی جدول و مقداری منحصر به فرد برای هر سطر جدول .
۲	domain	varchar(1000)	ذخیره دامنه جدول
۳	vote	float	وظیفه این ستون جدول ذخیره امتیاز دامنه می باشد. برای فرایند انتخاب دامنه در موتور جستجو از این بخش استفاده می شود.

## ۲- جدول امتیاز دامنه<sup>۱</sup>:

این جدول وظیفه نگهداری وزن های هر دامنه جهت به روز رسانی امتیاز آنها را دارد. برای به روز رسانی امتیاز هر دامنه با استفاده از رگرسیون خطی برای هر دامنه از این جدول استفاده می گردد که وظیفه ذخیره پارامترهای و به روز رسانی دامنه را به عهده دارد. در جدول ۴-۲ می توانید لیست ستون های این جدول را ببینید. و در شکل ۴-۸ دیاگرام این جدول را مشاهده می کنید.

<sup>۱</sup> domain rank

DOMAINRANK		
id	char(36)	PRIMARY KEY UNIQUE
domainID	char(36)	FOREIGN KEY DOMAIN(id)
preHypothesis	float	
newHypothesis	float	
W0	float	
W1	float	
W2	float	
W3	float	
learningRatio	float	

شکل ۴-۸ جدول امتیاز دامنه  
جدول ۴-۲ مشخصات جدول امتیاز دامنه

ردیف	نام ستون	نوع	توضیحات
۱	id	char(36)	کلید اصلی جدول و مقداری منحصر به فرد برای هر سطر جدول
۲	domainID	char(36)	کلید اصلی جدول دامنه (جهت ارتباط این دو جدول با یکدیگر)
۳	preHypothesis	float	ذخیره امتیاز قبلی دامنه در این سلول قرار دارد.
۴	newHypothesis	float	ذخیره امتیاز کنونی دامنه در این ستون قرار دارد
۵	W <sub>۰</sub>	float	ذخیره وزن پارامتر اول در دامنه
۶	W <sub>۱</sub>	float	ذخیره وزن پارامتر دوم در دامنه
۷	W <sub>۲</sub>	float	ذخیره وزن پارامتر سوم در دامنه
۸	W <sub>۳</sub>	float	ذخیره وزن پارامتر چهارم در دامنه
۹	learningRatio	float	نرخ یادگیری برای هر دامنه در روش رگرسیون

### ۳- جدول صفحات وب<sup>۱</sup>:

وظیفه این جدول ذخیره آدرس تمام صفحات وب داخل یک دامنه خاص می‌باشد. هر دامنه شامل چندین صفحه با آدرس مختلف می‌باشد که همه آنها در این جدول قرار دارند. در جدول ۴-۳ می‌توانید لیست

<sup>۱</sup>Page table

ستون‌های این جدول را ببینید. و در شکل ۹-۴ دیاگرام این جدول را مشاهده می‌کنید.

PAGE		
id	char(36)	PRIMARY KEY UNIQUE
domainID	char(36)	FOREIGN KEY DOMAIN(id)
uri	varchar(2000)	
vote	float	

شکل ۹-۴ جدول صفحه وب

جدول ۳-۴ مشخصات جدول دامنه

ردیف	نام ستون	نوع	توضیحات
۱	id	char(36)	کلید اصلی جدول و مقداری منحصر به فرد برای هر سطر جدول
۲	domainID	char(36)	کلید اصلی جدول دامنه ( جهت ارتباط این دو جدول با یکدیگر)
۳	uri	Varchar(2000)	وظیفه این ستون جدول ذخیره آدرس منحصر بفرد و نرمالایز شده هر صفحه وب می‌باشد.
۴	vote	float	وظیفه این ستون جدول ذخیره امتیاز صفحه وب می‌باشد. برای فرایند انتخاب صفحه وب در موتور جستجو از این بخش استفاده می‌شود.

#### ۴- جدول تاریخچه<sup>۱</sup>:

این جدول وظیفه نگهداری تاریخچه دریافت یک صفحه وب را دارد. بعد از دریافت هر صفحه وب محتوای صفحه وب در این جدول به همراه تاریخ دریافت ذخیره می‌شود. در جدول ۴-۴ می‌توانید لیست ستون‌های این جدول را ببینید. و در شکل ۱۰-۴ دیاگرام این جدول را مشاهده می‌کنید.

<sup>۱</sup> history table

HISTORY		
id	char(36)	PRIMARY KEY UNIQUE
pageID	char(36)	FOREIGN KEY PAGE(id)
response	int	
content	text	
seeImages	smallint	
getImage	smallint	
seeLinks	smallint	
normalLinks	smallint	
date	timestamp	

شکل ۴-۱۰ جدول تاریخچه

## ۵- جدول رابطه :

این جدول وظیفه ذخیره رابطه بین صفحات وب را بر عهده دارد. در حقیقت در این جدول لیست از ارتباطات صفحات مختلف وب را داریم. به عنوان مثال اگر صفحه وب X با صفحه وب Y ارتباط داشته باشد در این جدول یک سطر حاوی این ارتباط قرار دارد. در جدول ۴-۵ می‌توانید لیست ستون‌های این جدول را ببینید و در شکل ۴-۱۱ دیاگرام این جدول را مشاهده می‌کنید.

RELATION		
id	char(36)	PRIMARY KEY UNIQUE
inputlink	char(36)	FOREIGN KEY PAGE(id)
outputlink	char(36)	FOREIGN KEY PAGE(id)

شکل ۴-۱۱ جدول رابطه

جدول ۴-۴ مشخصات جدول تاریخچه

ردیف	نام ستون	نوع	توضیحات
۱	id	char(36)	کلید اصلی این جدول و مقداری منحصر به فرد برای هر سطر جدول است
۲	pageID	char(36)	کلید اصلی جدول صفحه. ( جهت ارتباط این دو جدول با یکدیگر ) هر صفحه وب شامل چندین تاریخچه از خود می باشد.
۳	response	int	امتیاز نتیجه دریافت هر صفحه وب در یک زمان مشخص می باشد. همینطور که قبلا گفته شد به عنوان مثال اگر صفحه وب با موفقیت دریافت شود این امتیاز ۲۰۰ می باشد.
۴	content	text	وظیفه این ستون جدول ذخیره خود صفحه وب می باشد که در غالب HTML می باشد.
۵	seeImages	smallInt	نمایش تعداد تصاویر دیده شده در یک صفحه وب خاص در یک زمان مشخص می باشد.
۶	getImages	smallInt	نمایش تعداد تصاویری که در یک صفحه وب خاص و در یک زمان مشخص با موفقیت دریافت شده اند.
۷	seeLinks	smallInt	تعداد آدرس صفحات وب دیده شده در یک صفحه وب مشخص و در یک زمان مشخص.
۸	normalLink	smallInt	تعداد آدرس صفحات وب نرمال شده در یک صفحه وب مشخص و در یک زمان مشخص.
۹	date	date	تاریخ دریافت صفحه وب

جدول ۵-۴ مشخصات جدول رابطه

ردیف	نام ستون	نوع	توضیحات
۱	id	char(36)	کلید اصلی این جدول و مقداری منحصر به فرد برای هر سطر جدول است
۲	inputLink	char(36)	کلید اصلی جدول صفحه ( جهت ارتباط این دو جدول با یکدیگر )
۳	outputLink	char(36)	کلید اصلی جدول صفحه ( جهت ارتباط این دو جدول با یکدیگر )

## ۶- جدول تصویر :

این جدول وظیفه ذخیره تمامی تصاویر یک صفحه وب مشخص را دارد. در جدول ۴-۶ می‌توانید لیست ستون‌های این جدول را ببینید. و در شکل ۴-۱۲ دیاگرام این جدول را مشاهده می‌کنید.

جدول ۴-۶ مشخصات جدول رابطه

ردیف	نام ستون	نوع	توضیحات
۱	id	char(36)	کلید اصلی این جدول و مقداری منحصر به فرد برای هر سطر جدول است
2	pageID	char(36)	کلید اصلی جدول صفحه وب (جهت ارتباط تصویر با صفحه وب)
۳	blob	text	این ستون جهت ذخیره محتوای تصویر دریافتی توسط موتور جستجو استفاده می‌شود.
۴	w	int	این ستون در بردارنده تعداد پیکسل سطرهای یک تصویر مشخص است
۵	h	int	این ستون در بردارنده تعداد پیکسل ستون‌های یک تصویر مشخص است
۶	src	text	این ستون دربرگیرنده آدرس تصویر دریافتی می‌باشد.

IMG		
id	char(36)	PRIMARY KEY UNIQUE
pageID	char(36)	FOREIGN KEY PAGE(id)
blob	text	
w	int	
h	int	
src	text	

شکل ۴-۱۲ جدول تصویر



## ۷-جدول متن<sup>۱</sup>:

این جدول وظیفه ذخیره تمامی متن‌های نزدیک یک تصویر مشخص به همراه امتیاز آن متن می‌باشد. در جدول ۷-۴ می‌توانید لیست ستون‌های این جدول را ببینید. و در شکل ۴-۱۳ دیاگرام این جدول را مشاهده می‌کنید.

TEXT		
id	char(36)	PRIMARY KEY UNIQUE
imgID	char(36)	FOREIGN KEY IMG(id)
text	text	
distance	Varchar(100)	
tag	text	

شکل ۴-۱۳ جدول متن

در انتها و در شکل ۴-۱۴ می‌توانید لیست کامل جدول‌های پایگاه داده به همراه ارتباط آنها با یکدیگر را ببینید.

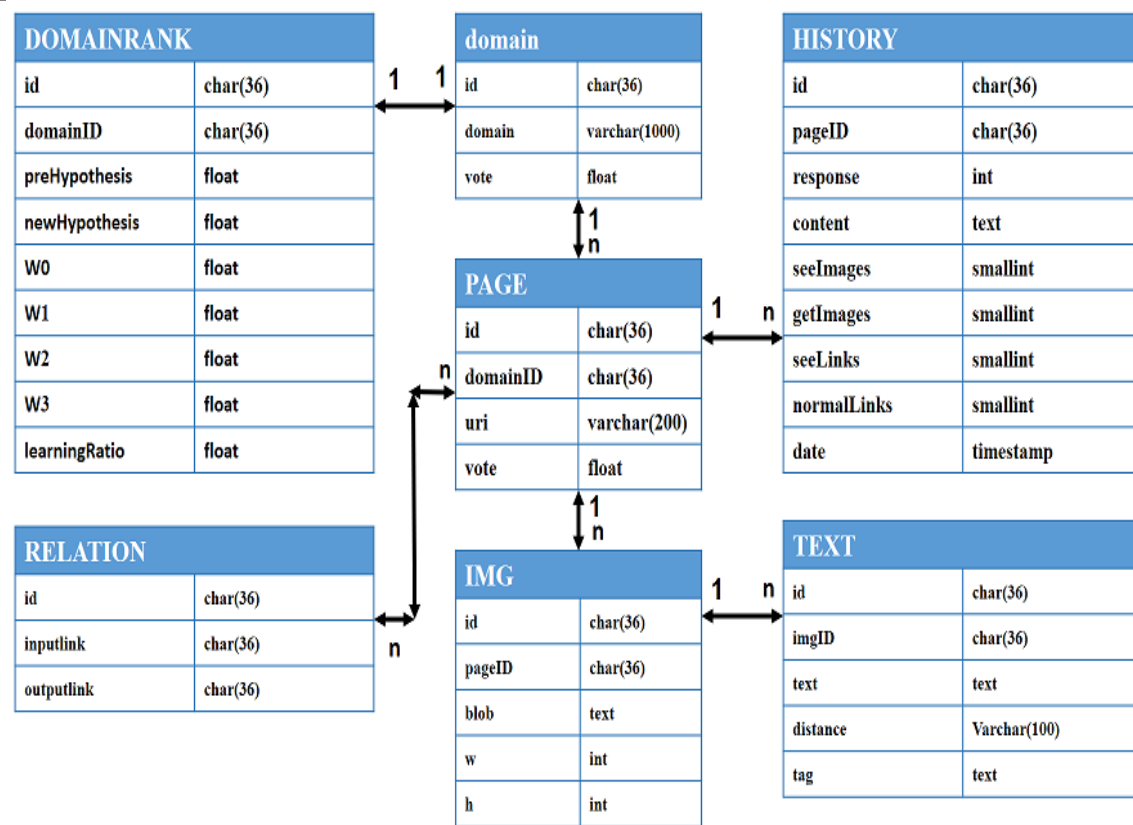
در ارتباطات عبارات ۱ و n نمایشگر تعداد رابطه بین جداول است به عنوان مثال و مطابق شکل هر سطر از جدول page با تعداد زیادی از جدول history ارتباط دارد.

---

<sup>۱</sup>text table

جدول ۷-۴ مشخصات جدول متن

ردیف	نام ستون	نوع	توضیحات
۱	id	char(36)	کلید اصلی این جدول و مقداری منحصر به فرد برای هر سطر جدول است
۲	imgID	char(36)	کلید اصلی جدول تصویر (جهت ارتباط تصویر با صفحه وب)
۳	text	text	این ستون جهت ذخیره یکی از متون نزدیک یک تصویر مشخص می باشد.
۴	distance	Varchar(100)	این ستون در بردارنده فاصله متن از تصویر در درخت HTML می باشد
۵	tag	text	این ستون در بردارنده نام تگی است که این متن را دربر گرفته است.



شکل ۴-۱۴ لیست روابط در پایگاه داده

## فصل پنجم:

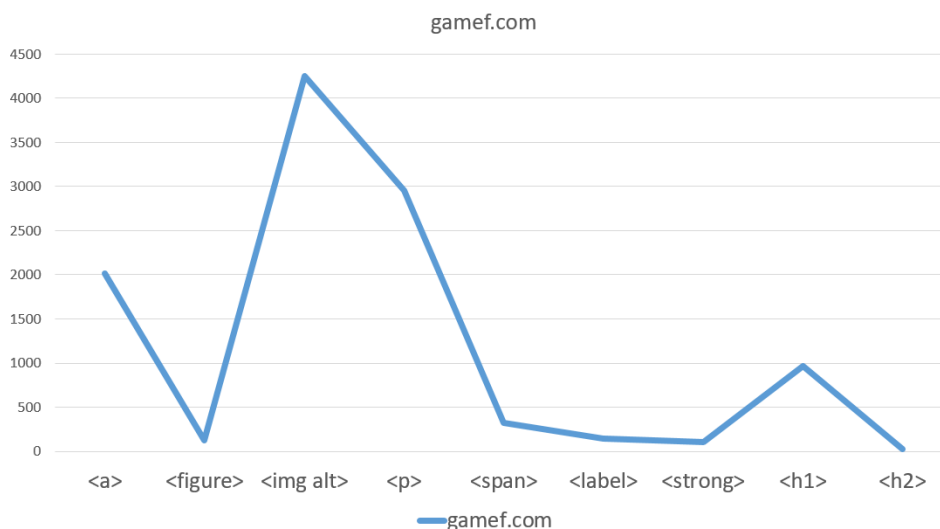
## نتیجه گیری

## ۵-۱ مقدمه

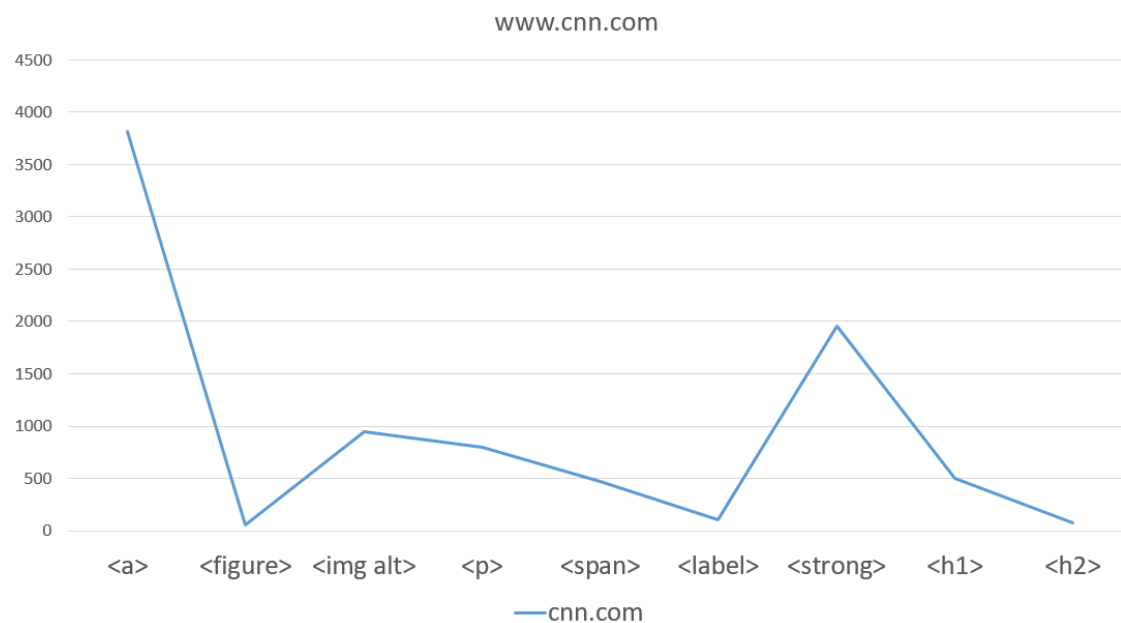
در فصول قبل ابتدا به بررسی صفحات وب و انواع آنها پرداختیم. سپس در فصل سوم به بررسی چالش ها و راهکار های برای آنها پرداختیم. و در فصل چهارم هم به بررسی اجمالی از روند کار موتور جستجو پرداختیم در این فصل که فصل نتیجه گیری می باشد حالا به بررسی چند خروجی و مقایسه روش هایی مختلف در پروژه می پردازیم.

## ۵-۲ اهمیت مکانی متون در صفحه وب

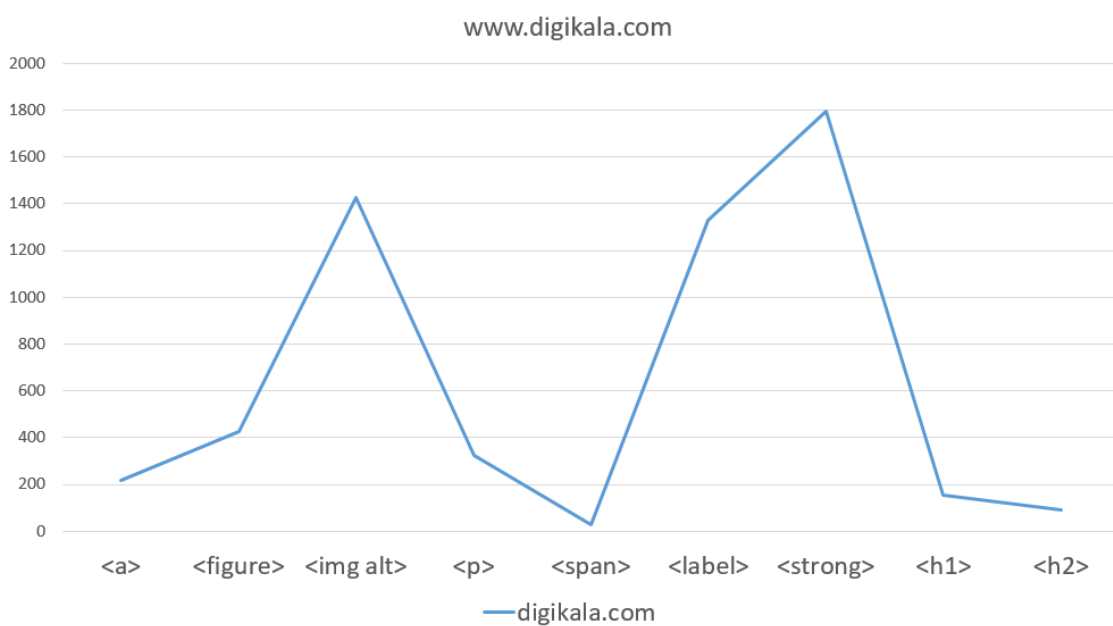
خوشبختانه صفحات وب (کلیه صفحات یک وب سایت در غالب یک دامنه) از لحاظ غالب و تصویر حالت یکسانی دارند. به عنوان مثال در شکل های ۵-۱ الی ۵-۶ به بررسی ۶ دامنه مختلف و مکان قرار گرفتن بهترین متون نزدیک تصاویر را در آنها ببینید (در اینجا منظور از تگ های مفید تگ های است که اطلاعات داخلی آنها بیشترین شباهت را با لیست پرس وجوهای داخل موتور جستجو دارند چه از نظر متنی چه از نظر تصویری). همانطور که در تصاویر مشخص است هر سایت تقریباً یک غالب مشخص دارد.



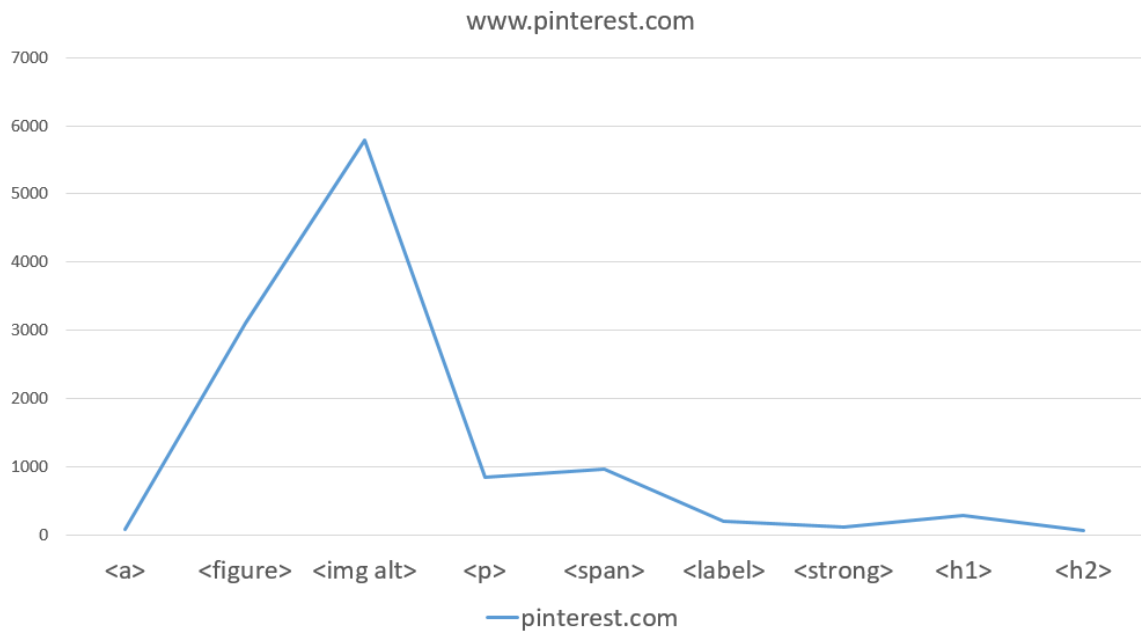
شکل ۵-۱ مقایسه مکان تصاویر در دامنه ۱-



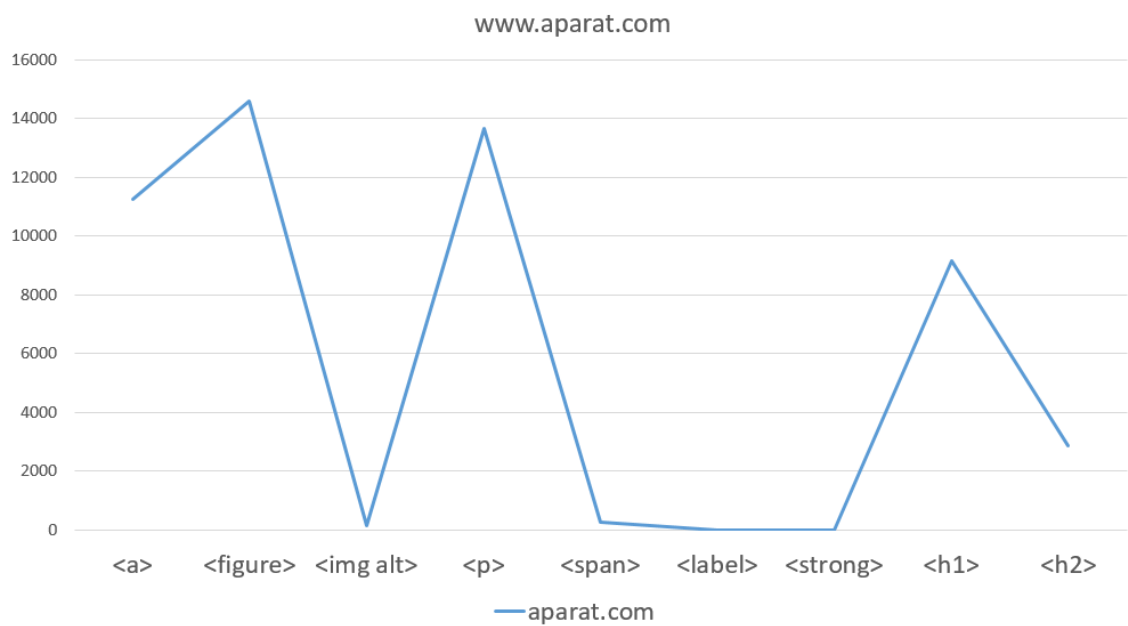
شکل ۲-۵ مقایسه مکان تصاویر در دامنه ۲-



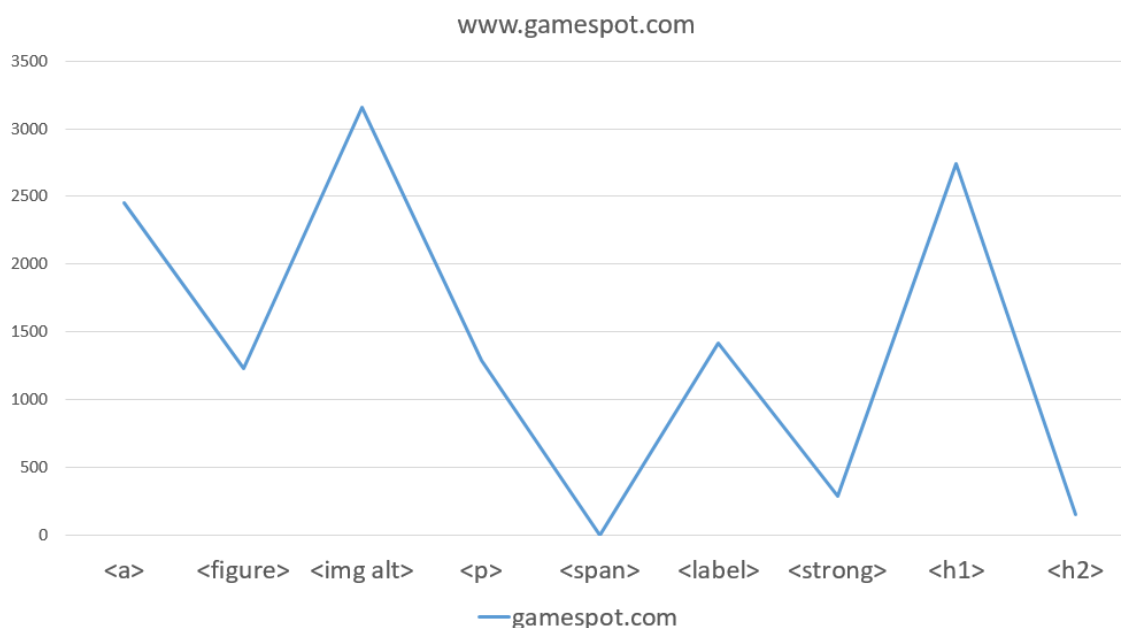
شکل ۳-۵ مقایسه مکان تصاویر در دامنه ۳-



شکل ۴-۵ مقایسه مکان تصاویر در دامنه ۴-



شکل ۵-۵ مقایسه مکان تصاویر در دامنه ۵-



شکل ۵-۶ مقایسه مکان تصاویر در دامنه ۶-

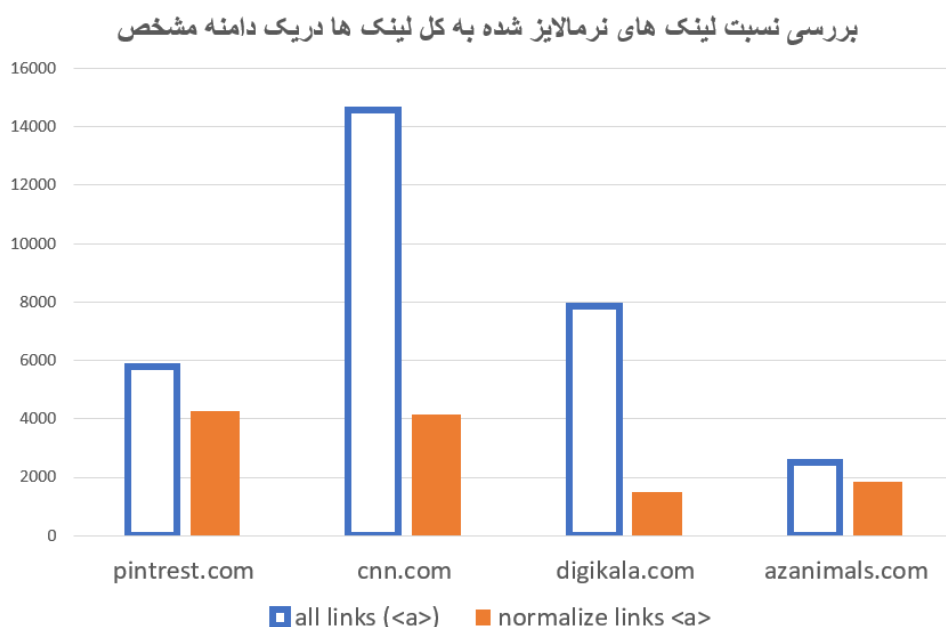
در این جا منظور بررسی متون نزدیک تصویر می باشد که اطلاعات خوبی درباره تصویر را در خود دارد و همینطور خود تصاویر. به عنوان مثال در سایت gamefa.com بیشترین متون مفید نزدیک تصویر (مشابه پرس وجوهای موتور جستجو) در غالب پاراگراف بالای آن و متن توضیحی داخل خود تصویر و توضیحات داخل اتصالات (تگ <a>) می باشد. پس می شود نتیجه گرفت که با بررسی در هر دامنه خاص به اهمیت مکان متون و تصاویر با ارزش درباره یک پرس وجوی خاص مشخص می شود. البته بهتر است این مقادیر با دیدن صفحات جدید به روز شود. ولی در کل معمولا غالب صفحات وب در یک دامنه یکسان است.

به عنوان مثالی دیگر در دامنه digikala.com همانطور که مشخص است تگ های <strong> و <label> بیشترین شباهت را با پرس وجوهای داخل موتور جستجو دارند.

### ۳-۵ تاثیر نرمال سازی در کاهش بار پردازشی صفحات وب

در اینجا به بررسی تاثیر نرمال سازی آدرس صفحات در یک دامنه خاص می پردازیم. برای این بررسی ما بیش از ۵۰۰ صفحه وب از ۳۰ دامنه متفاوت را مورد بررسی قرار داده و از همه این ۵۰۰ صفحه تمام

لینک‌های داخلی آنرا بدست آوردیم، سپس تمام این آدرس ها را نرمالایز کردیم نتیجه بدست آمده حیرت آور است. بعد از ارسال موتور جستجو به این صفحات وب و استخراج تمام لینک ها بیش از ۷۸۳۹ لینک تولیدگردید است که بعد از نرمالایزه کردن و حذف دامنه های تکراری، تعداد این لینک ها به ۲۰۹۳ عدد  $(\frac{7839}{2093} = 3.74)$  کاهش می‌یابد. که بسیار عالی می‌باشد. پس می‌توان نتیجه گرفت که برای کاهش ترافیک شبکه و همینطور کاهش بار پردازش موتور جستجو بهتر است قبل از هر فرایندی آدرس را نرمالایز کنیم. در شکل ۵-۷ به بررسی ۴ دامنه پرداختیم که برای هر دامنه ۱۰۰۰ صفحه وب را بررسی کرده و تعداد لینک بدست آمده را با تعداد لینک بعد از نرمالایز کردن مقایسه نموده‌ایم.



شکل ۵-۷ مقایسه نتیجه نرمال سازی صفحات وب

#### ۴-۵ تایین امتیاز هر صفحه وب.

تغییر امتیاز هر صفحه وب همانطور که قبلا هم گفته شد در چند مرحله اتفاق می‌افتد. ممکن است امتیاز صفحه افزایش یا کاهش بیابد و این تغییر امتیاز صفحه روی صفحاتی که با این صفحه ارتباط دارند (به



صفحه فوق لینک ورودی داده اند و از صفحه خود به این صفحه یک اتصال ایجاد کرده اند) هم تاثیر دارد و باید امتیاز آنها هم دچار تغییر گردد.

#### ۱-۴-۵ دلایل اصلی کاهش امتیاز یک صفحه وب

۱- عدم دسترسی به صفحه وب :

دلیل این اتفاق یکی از عوامل زیر می باشد

- **عدم داشتن مجوز لازم جهت ورود به صفحه وب :** ممکن است صفحه وب قرار داشته

روی یک خدمت گذار مجوز عمومی برای دسترسی به یک صفحه وب را ندهد و برای ورود

به این صفحه وب باید اقداماتی ماننده اهراز هویت صورت بگیرد. شاید یک راه حل برای

رفع این مشکل استفاده از پروتکل http بجای https باشد و یا حتی استفاده از موتورهای

جستجوی مخفی.

- **حذف کامل صفحه وب از روی خدمت گذار :** ممکن است صفحه وب انتخابی بنا به

دلایلی (ماننده انقراض و قدیمی شدن اطلاعات صفحه وب) از خدمت گذار خود حذف شده

باشد.

- **حذف موقت صفحه وب از خدمت گذار :** ممکن است صفحه وب انتخابی بنا به دلایلی

(تغییر محتوای صفحه وب ماننده غالب یا اطلاعات) به طور موقت از روی خدمت گذار خود

حذف شده باشد. و بعد از مدتی دوباره بر روی خدمت گذار خود قرار بگیرد.

**راه حل :** برای موارد بالا بهتر است بعد از بروز مشکل مدتی صفحه وب از لیست صفحات وب ما خارج

گردد

و بعد از یک مدت دوباره در لیست بررسی قرار بگیرد. همانطور که در جدول ۵-۱ می‌بینید ما یک سایت خاص (gamefa.com) را مورد بررسی قرار دادیم و همانطور که می‌بینید به عنوان مثال در صورت عدم دسترسی به یک صفحه وب معمولا غالب سایت در حال به روز رسانی می‌باشد و بعد از مدت تقریبا ۳۵ ساعت دوباره در دسترس قرار می‌گیرد. و تقریبا در صورت حذف یک مطلب بعد از ۴ بار بررسی می‌شود نتیجه گرفت مطلب به کلی از سایت و دامنه فوق حذف گردیده است. پس اگر موتور جستجو این مطلب را متوجه شود بعد از یک دوره زمانی ۳۵ ساعته دوباره آدرس فوق را در لیست بررسی خود قرار می‌دهد و در بررسی مشاهده گردید که اغلب سایت ها در این مورد یک استاندارد نسبتا ثابتی را رعایت می‌کنند. در جدول ۵-۱ ما به بررسی چند زیر دامنه از دامنه gamefa.com پرداخته و همانطور که می‌بینید زمان غیر قابل از دست دادن صفحات در آنها مشخص شده است.

جدول ۵-۱ بررسی ۸ زیر دامنه در یک دامنه خاص متوسط در دسترس بودن بعد از ۴ بار بررسی

ردیف	نام آدرس	زمان در دسترس بودن
۱	<a href="http://cinema.gamefa.com/">http://cinema.gamefa.com/</a>	۳۱ ساعت و ۲۴ دقیقه
۲	<a href="http://gamefa.com/category/platform/playstation/۳">http://gamefa.com/category/platform/playstation/۳</a>	۳۶ ساعت و ۲۹ دقیقه
۳	<a href="http://gamefa.com/category/tech/">http://gamefa.com/category/tech/</a>	۳۲ ساعت و ۱۸ دقیقه
۴	<a href="http://gamefa.com/category/tech/hardware/">http://gamefa.com/category/tech/hardware/</a>	۴۱ ساعت و ۱۲ دقیقه
۵	<a href="http://gamefa.com/category/featured/news-game/">http://gamefa.com/category/featured/news-game/</a>	۲۹ ساعت و ۱۴ دقیقه
۶	<a href="http://cinema.gamefa.com/category/news/">http://cinema.gamefa.com/category/news/</a>	۱۲ ساعت و ۳۵ دقیقه
۷	<a href="http://cinema.gamefa.com/category/tv-series-news/">http://cinema.gamefa.com/category/tv-series-news/</a>	۱۹ ساعت و ۴۹ دقیقه
۸	<a href="http://cinema.gamefa.com/writers/">http://cinema.gamefa.com/writers/</a>	۳۹ ساعت و ۲۳ دقیقه

مهمترین کاربرد این بخش در کاهش بار ترافیکی و پردازشی موتور جستجو خود را نشان می‌دهد. به عنوان مثال اگر ما بعد از مشاهده یک خطا در دریافت یک صفحه وب دوباره اقدام به دریافت این صفحه کنیم با

یک ترافیک بالا و بی‌فایده طرف خواهیم بود که مدام در حال دریافت صفحات وبی می‌باشد که به طور موقت غیر فعالند و هیچ داده‌ای نیز جمع‌آوری نمی‌شود.

## نوآوری :

در این جا بعد از بروز مشکل عدم دسترسی علاوه بر کاهش امتیاز صفحه وب و دامنه آن یک دوره استراحت برای آن در نظر گرفته می‌شود (البته برای هر دامنه یک مقدار) و بعد از اتمام زمان بروز مشکل اگر کماکان صفحه وب در دسترس نبود علاوه بر کاهش امتیاز صفحه وب و دامنه آن، زمان استراحت دامنه را هم افزایش می‌دهیم با استفاده از فرمول ۱-۵ .

$$\text{RestTime}_{\text{domain}} = \text{RestTime}_{\text{Domain}} + \alpha * \text{PageRank}_{\text{selected web page}} \quad (\text{فرمول ۱-۵})$$

حالا به معرفی پارامترهای این فرمول می‌پردازیم.

- $\text{RestTime}_{\text{domain}}$  : منظور زمان استراحت کلی یک دامنه است بر حسب میلی ثانیه بعد از وقوع یک مشکل.
- $\text{PageRank}_{\text{selected web page}}$  : امتیاز صفحه وب جاری
- $\alpha$  : ضریب میرایی که ثابتی است که به فرم تجربی بدست می‌آید و برای هر دامنه مقدار مختلفی دارد.

در اینجا به مرور زمان، زمان استراحت صفحه وب افزایش می‌یابد پس برای جلوگیری از این مشکل بعد از اتمام زمان استراحت و در صورت در دسترس بودن صفحه وب (بعد از زمان استراحت و اقدام به دریافت دوباره صفحه وب با موفقیت دریافت می‌شود) انتخابی زمان استراحت را با همان فرمول بالا کاهش می‌دهیم فقط ثابت  $\alpha$  ضریب منفی می‌گیرد. از طرفی بعد از چند دفعه غیر فعال شدن صفحه وب می‌شود نتیجه گرفت صفحه به طور کامل

از سایت حذف گردیده است و می‌شود با یک اطمینان بالای آنرا از لیست آدرسها حذف کرد. مورد آخر این که در بعضی از دامنه‌ها کل دامنه یک جا و به دلایل مختلفی از لیست خدمات‌دهنده برای یک مدت مشخص خارج می‌شود. یعنی با ملاحظه خرابی حتی یک صفحه از دامنه می‌شود به این نتیجه رسید که دامنه موقتاً از کار افتاده است (معمولاً جهت به روز رسانی دامنه).

۲- کاهش امتیاز به دلیل لینک دادن به صفحه وبی که امتیازش کاهش یافته است :

فرض کنید صفحه‌وبی به نام X از یک دامنه خبری به نام Y به یک صفحه‌وب به نام Z از یک سایت دیگر به نام W لینک اتصالی داده است (به عنوان مثال لینک خبر را از سایت فوق گرفته است) و حالا صفحه وب Z بنا به دلایلی از دامنه W حذف گردیده است. پس به نحوی موتور جستجو باید امتیاز صفحه X را هم کاهش بدهد. دلیل این امر مشکل در لینک اتصالی است که فقط یک بار پردازشی اضافه روی موتور جستجو می‌گذارد. البته ممکن است صفحه وب موقتاً از خدمت‌گذار حذف گردیده باشد پس بهتر است بعد از برگشت به صفحه وب امتیاز صفحاتی که به آن لینک داده بودند نیز افزایش یابد.

## ۵-۵ بررسی اهمیت متون نزدیک تصاویر در صفحه وب

شاید مهمترین نوآوری این پروژه را بتوان ادغام پردازش تصویر و متن‌های داخل صفحه وب با پرس‌وجوهای مورد جستجو نامید. در این مورد جدول زیر مشخص کننده صحت این مطلب می‌باشد. همانطور که می‌بینید در این جدول ۱۰ دامنه را انتخاب کرده و در هر کدام لیست تصاویری که در کنارشان حداقل یک متن مطابق با تصویر داشته باشد در مقابل تصاویری که حتی یک متن مشابه موضوع تصویر در آنها باشد را مشاهده می‌کنیم. به عنوان مثال دامنه blaxup.com تقریباً نزدیک هر تصویر ۲ متن مشابه با آن وجود دارد (مشابه پرس‌وجوهای داخل موتور جستجو). این مطلب را می‌توان در شکل ۵-۸ هم مشاهده کرد. همانطور که می‌بینید در کنار تصویر ۲ مطلب کاملاً مرتبط با تصویر وجود دارد. همانطور که قبلاً هم بررسی کردیم

غالب دامنه های وب که شامل صد ها یا حتی هزاران صفحه وب هستند معمولاً یکی است. پس می شود نتیجه گرفت که به عنوان مثال در دامنه blaxup.com و نزدیک هر تصویر حداقل دو متن مرتبط با آن وجود دارد. این روش کمک شایانی در افزایش سرعت پردازش داده دارد به عنوان مثال وقتی موتور جستجو یک صفحه وب جدید را از دامنه blaxup.com دریافت می کند با یک درجه بالای می تواند به متون نزدیک تصویر و محتوای آنها اعتماد بکند یا نه. خوشبختانه در غالب سایت می شود به این مطلب رسید و دید که این نرخ تقریباً در دامنه سایت یکسان می ماند. به عنوان مثال در جدول ۲-۵ لیست ۵ صفحه وب جدید اضافه شده به دامنه pinterest.com را می بینید که نرخ متن مرتبط با تصویر آنها بدون محاسبه نزدیک نرخ دامنه سایت می باشد ( $\frac{7480}{8412} = 0.88$ ). پس می شود حتی بدون محاسبه و با دیدن لینک صفحه وب جدید متوجه تعداد متون نزدیک یک تصویر گردید ولی برای بروز رسانی اطلاعات محاسبه اطلاعات هر صفحه و اعمال روی دامنه کلی اجباری است.



شکل ۵-۸ متون نزدیک و مرتبط به تصویر در یک دامنه خاص

جدول ۵-۲ بررسی نسبت متون مرتبط با تصاویر در ۱۰ دامنه مختلف

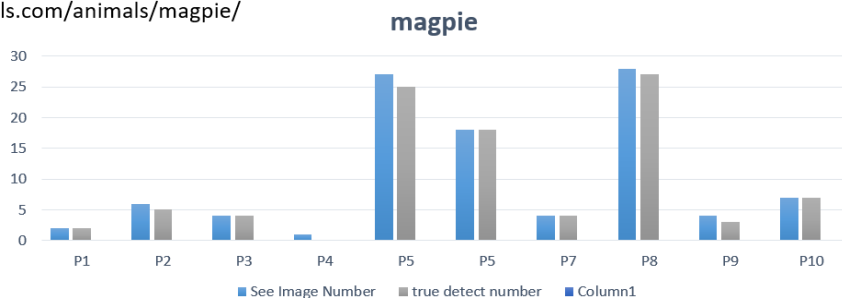
ردیف	دامنه	متن مرتبط با تصویر
۱	<a href="http://edition.cnn.com/">http://edition.cnn.com/</a>	$۱,۶۸ = \frac{3409}{2019}$
۲	<a href="http://p۳۰.download.com/">http://p۳۰.download.com/</a>	$۱,۶۵ = \frac{19127}{11489}$
۳	<a href="http://gamefa.com/">http://gamefa.com/</a>	$۲,۴۲ = \frac{8652}{3568}$
۴	<a href="http://www.blaxup.com/">http://www.blaxup.com/</a>	$۲,۰۱ = \frac{2981}{1478}$
۵	<a href="http://www.aparat.com/">http://www.aparat.com/</a>	$۱,۱۷ = \frac{21774}{18592}$
۶	<a href="http://www.intelcapital.com">www.intelcapital.com</a>	$۰,۸۷ = \frac{4198}{4785}$
۷	<a href="https://twitter.com">https://twitter.com</a>	$۱,۳ = \frac{14796}{11325}$
۸	<a href="https://www.forbes.com">https://www.forbes.com</a>	$۱,۴۴ = \frac{1420}{984}$
۹	<a href="http://pinterest.com">pinterest.com</a>	$۰,۸۸ = \frac{7480}{8412}$
۱۰	<a href="http://www.vajehyab.com">www.vajehyab.com</a>	$۰,۵۸ = \frac{1028}{1756}$

۵-۶ بررسی نتایج تست موتور جستجو

در اینجا بعد از فاز آموزشی توسط برنامه کافی خروجی ما یک شبکه وزنی بزرگ می‌باشد که جهت تست تصاویر ورودی از آنها استفاده می‌گردد. در تصاویر ۵-۹ الی ۵-۱۹ نمایش از تست تعداد ۹ پرس‌وجوی بدست آمده از شبکه عصبی خروجی (به فرم تصادفی) با ۱۰ صفحه وب که درباره این موضوع مطالبی دارند نمایش داده شده است. در تصاویر به عنوان مثال اگر در مجموع این صفحات حاوی ۴۰ تصویر مشابه با گروه انتخابی باشند. اگر

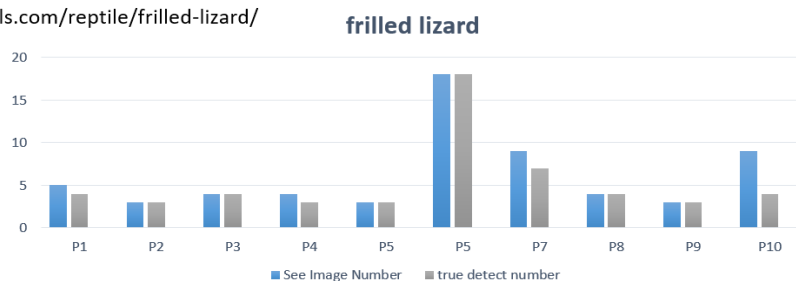
شبکه ما تعداد ۳۰ تا از این تصاویر را درست و مطابق این گروه انتخاب کرده باشد دقت این شبکه عصبی در این گروه ۷۵٪ می‌باشد.

P1 : <https://en.wiktionary.org/wiki/magpie>  
P2 : <http://www.abc.net.au/news/2017-08-31/make-friends-with-magpies-to-avoid-swooping-expert-says/8856438>  
P3 : <https://www.rspb.org.uk/birds-and-wildlife/bird-and-wildlife-guides/bird-a-z/m/magpie/>  
P4 : <https://www.collinsdictionary.com/dictionary/english/magpie>  
P5 : <http://www.arkive.org/magpie/pica-pica/>  
P6 : <http://what-when-how.com/birds/black-billed-magpie-birds/>  
P7 : <https://www.britannica.com/animal/magpie>  
P8 : <http://nzbirdsonline.org.nz/species/australian-magpie>  
P9 : <http://www.birdsinbackyards.net/species/Cracticus-tibicen>  
P10: <https://a-z-animals.com/animals/magpie/>



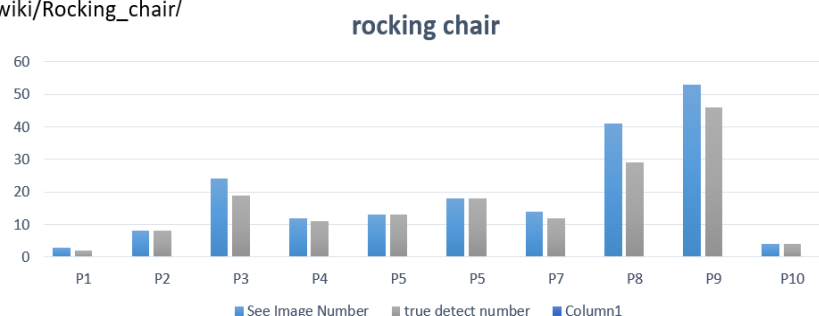
شکل ۵-۹ بررسی ۱۰ صفحه درباره کلاغ زاغی

P1 : <http://en.wikipedia.org/wiki/Chlamydosaurus/>  
P2 : <http://www.nationalgeographic.com/animals/reptiles/f/frilled-lizard/>  
P3 : <https://www.britannica.com/animal/frilled-lizard/>  
P4 : <https://a-z-animals.com/animals/frilled-lizard/>  
P5 : <http://easyscienceforkids.com/all-about-frilled-lizards/>  
P6 : <http://www.arkive.org/frilled-lizard/chlamydosaurus-kingii/>  
P7 : [http://www.outback-australia-travel-secrets.com/frilled\\_lizard.html/](http://www.outback-australia-travel-secrets.com/frilled_lizard.html/)  
P8 : <http://reptilepark.com.au/animals/reptiles/dragons/frilled-neck-lizard/>  
P9 : <http://www.reptilesmagazine.com/Care-Sheets/Frilled-Lizard-Care-Sheet/>  
P10: <https://www.aboutanimals.com/reptile/frilled-lizard/>



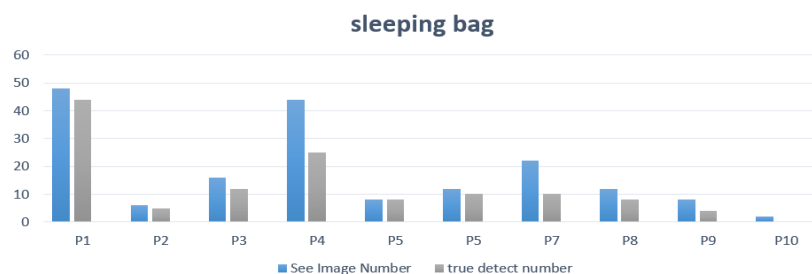
شکل ۵-۱۰ بررسی ۱۰ صفحه درباره مارمولک فرشته

P1: <https://www.lowes.com/pd/Garden-Treasures-Black-Patio-Rocking-Chair/3881543/>  
P2: <https://www.wayfair.com/furniture/sb0/rocking-chairs-c413906.html>  
P3: <https://www.houzz.com/photos/rocking-chairs/>  
P4: <http://www.ikea.com/us/en/catalog/products/S89861010/>  
P5: [http://www.garyweeks.com/rocking\\_chairs.htm/](http://www.garyweeks.com/rocking_chairs.htm/)  
P6: <https://www.polywoodoutdoor.com/styles/rocking-chairs.html/>  
P7: [https://www.amazon.com/Chairs-Rocking-Living-Room-Furniture/s?ie=UTF8&page=1&rh=n%3A3733491%2Cp\\_n\\_feature\\_three\\_browse-bin%3A3486058011/](https://www.amazon.com/Chairs-Rocking-Living-Room-Furniture/s?ie=UTF8&page=1&rh=n%3A3733491%2Cp_n_feature_three_browse-bin%3A3486058011/)  
P8: <http://shop.crackerbarrel.com/home-furniture/indoor-furniture/rocking-chairs/>  
P9: <https://www.martinandmacarthur.com/koa-furniture/living-room/rocking-chairs/>  
P10: [http://en.wikipedia.org/wiki/Rocking\\_chair/](http://en.wikipedia.org/wiki/Rocking_chair/)



جدول ۵-۱۱ بررسی ۱۰ صفحه درباره صندلی راحتی

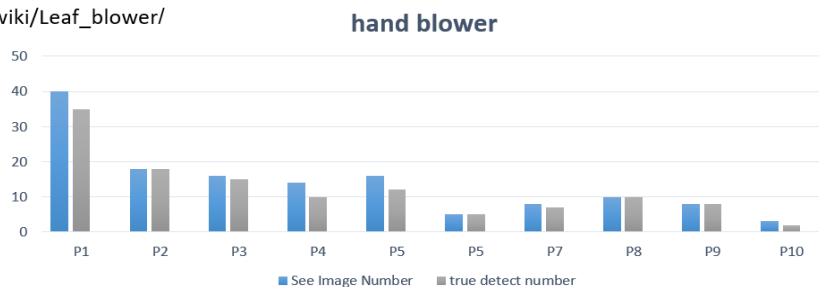
P1: <http://www.amazon.com/camping-sleeping-bags/b?ie=UTF8&node=2204498011/>  
P2: [https://en.wikipedia.org/wiki/Sleeping\\_bag/](https://en.wikipedia.org/wiki/Sleeping_bag/)  
P3: <https://www.thenorthface.com/shop/equipment-sleeping-bags/>  
P4: <https://www.rei.com/c/sleeping-bags-and-accessories?r=c&origin=web&ir=category%3Asleeping-bags-and-accessories&page=1/>  
P5: <http://www.sportsdirect.com/camping/sleeping-bags/>  
P6: <http://www.gooutdoors.co.uk/camping/sleeping-bags/>  
P7: <http://www.huskyeu.eu/sleeping-bags/>  
P8: <http://featheredfriends.com/down-sleeping-bags/three-season.html/>  
P9: <http://sleepingbags-cumulus.eu/uk/categories/sleeping-bags/down-sleeping-bags/>  
P10: <http://haeld.com/>



شکل ۵-۱۲ بررسی ۱۰ صفحه درباره کیسه خواب

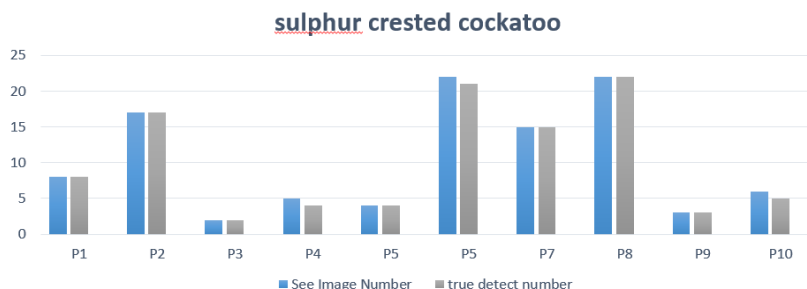


P1 : <https://www.homedepot.ca/en/home/categories/outdoors/outdoor-power-equipment/leaf-blowers.html/>  
P2 : <http://dir.indiamart.com/impcat/hand-blower.html/>  
P3 : <http://www.popularmechanics.com/home/tools/reviews/g168/which-leaf-blower-is-the-best-we-test-8-to-find-out/>  
P4 : <http://www.popularmechanics.com/home/tools/g2195/best-battery-powered-leaf-blowers/>  
P5 : <http://www.husqvarna.com/us/products/leaf-blowers/>  
P6 : <https://www.stihlusa.com/products/blowers-and-shredder-vacs/>  
P7 : <https://www.stihlusa.com/products/blowers-and-shredder-vacs/homeowner-blowers/>  
P8 : [https://www.amazon.in/ELECTRIC-HAND-BLOWER-EBC-ELECTREX/dp/B00RXQGGJH/259-5963314-6545724?\\_encoding=UTF8&%2AVersion%2A=1&%2Aentries%2A=0&portal-device-attributes=desktop/](https://www.amazon.in/ELECTRIC-HAND-BLOWER-EBC-ELECTREX/dp/B00RXQGGJH/259-5963314-6545724?_encoding=UTF8&%2AVersion%2A=1&%2Aentries%2A=0&portal-device-attributes=desktop/)  
P9 : <https://www.consumerreports.org/cro/leaf-blowers/buying-guide/>  
P10 : [http://en.wikipedia.org/wiki/Leaf\\_blower/](http://en.wikipedia.org/wiki/Leaf_blower/)



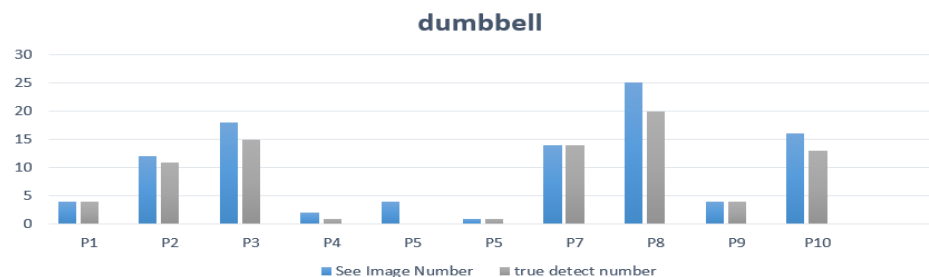
شکل ۵-۱۳ بررسی ۱۰ صفحه درباره دمنده دستی

P1 : [http://www.rollandbirdsfamily.com/available-parrots.html?gclid=Cj0KCQjwiqTnBRDVARIsAGsd9MqrUuUUAu2PqYNGx22Z4IKNFYj4n0VHvUY\\_XBVBt5MbAXgvZPA-5jsaAj5OEALw\\_wcB/](http://www.rollandbirdsfamily.com/available-parrots.html?gclid=Cj0KCQjwiqTnBRDVARIsAGsd9MqrUuUUAu2PqYNGx22Z4IKNFYj4n0VHvUY_XBVBt5MbAXgvZPA-5jsaAj5OEALw_wcB/)  
P2 : <http://www.thespruce.com/sulphur-crested-cockatoos-390714/>  
P3 : <http://www.birdsinbackyards.net/species/Cacatua-galerita/>  
P4 : [http://animaldiversity.org/accounts/Cacatua\\_galerita/](http://animaldiversity.org/accounts/Cacatua_galerita/)  
P5 : <http://cockatoo-info.com/species/sulphur-crested-cockatoo/>  
P6 : <http://www.bizrate.com/cockatoo-crested-sulphur/>  
P7 : <http://www.oiseaux-birds.com/card-sulphur-crested-cockatoo.html/>  
P8 : <http://nzbirdsonline.org.nz/species/sulphur-crested-cockatoo/>  
P9 : <http://www.billabongsanctuary.com.au/sulphu/>  
P10 : [http://en.wikipedia.org/wiki/Sulphur-crested\\_cockatoo/](http://en.wikipedia.org/wiki/Sulphur-crested_cockatoo/)



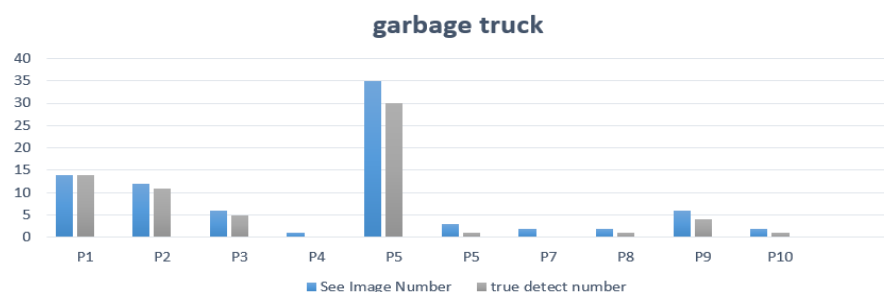
شکل ۵-۱۴ بررسی ۱۰ صفحه درباره نوعی طوطی سفید

P1 : <https://en.wikipedia.org/wiki/Dumbbell/>  
P2 : <http://www.fitness-superstore.co.uk/gym-equipment/dumbbells-kettlebells/dumbbells/>  
P3 : <https://www.amazon.com/dumbbells-dumbbell-set/b?ie=UTF8&node=3408401/>  
P4 : <https://dumbbell-exercises.com/>  
P5 : <http://anilhasehalteresmusculacao.com.br/dumbbells/>  
P6 : <https://en.wiktionary.org/wiki/dumbbell/>  
P7 : <https://www.roguefitness.com/conditioning/strength-equipment/dumbbells/>  
P8 : <http://www.sportsdirect.com/fitness-and-exercise-equipment/gym-equipment/weights-and-dumbbells/>  
P9 : <https://www.amazon.in/Dumbbells/b?ie=UTF8&node=3404717031/>  
P10: <https://www.flipkart.com/sports-fitness/exercise-fitness/dumbbells/pr?sid=dep,xnh,f92/>



شکل ۵-۱۵ بررسی ۱۰ صفحه درباره دمبل

P1 : <http://www.mingnuoclean.com/en/>  
P2 : [http://en.wikipedia.org/wiki/Garbage\\_truck/](http://en.wikipedia.org/wiki/Garbage_truck/)  
P3 : <http://www.mingnuoclean.com/en/>  
P4 : [http://en.wiktionary.org/wiki/garbage\\_truck/](http://en.wiktionary.org/wiki/garbage_truck/)  
P5 : <http://www.hyundai.com/worldwide/en/special-cv/garbage-truck/roll-pack/>  
P6 : <http://www.trackyourtruck.com/fleet-tracking-systems/industry/garbage-trucks/>  
P7 : <http://www.presstelegram.com/general-news/20170830/firefighters-put-out-garbage-truck-blaze-at-long-beach-va/>  
P8 : <http://shop.lego.com/en-LU/Garbage-Truck-10519/>  
P9 : <http://www.spoonflower.com/tags/garbage%20truck/>  
P10: <http://www.pepperl-fuchs.us/usa/en/12393.htm/>



شکل ۵-۱۶ بررسی ۱۰ صفحه درباره ماشین حمل زباله

P1 : <https://www.apple.com/lae/ipod-touch/>  
P2 : <https://www.apple.com/shop/buy-ipod/ipod-touch/>  
P3 : [https://en.wikipedia.org/wiki/IPod\\_Touch/](https://en.wikipedia.org/wiki/IPod_Touch/)  
P4 : <https://www.theguardian.com/technology/askjack/2017/aug/10/apple-discontinuing-ipods-alternative-mp3-player-podcasts-audible-audio-books/>  
P5 : [http://www.argos.co.uk/static/Browse/ID72/33012867/c\\_1/1%7Ccategory\\_root%7CTechnology%7C33006169/c\\_2/2%7C33006169%7CIPod%2C+MP3+and+headphones%7C33008472/c\\_3/3%7Ccat\\_33008472%7CIPod%7C33012867.htm/](http://www.argos.co.uk/static/Browse/ID72/33012867/c_1/1%7Ccategory_root%7CTechnology%7C33006169/c_2/2%7C33006169%7CIPod%2C+MP3+and+headphones%7C33008472/c_3/3%7Ccat_33008472%7CIPod%7C33012867.htm/)  
P6 : <https://www.walmart.com/ip/Apple-iPod-touch-16GB/42608121/>  
P7 : <http://www.ifixit.com/Device/iPod/>  
P8 : <http://www.cnet.com/products/apple-ipod-classic-7th-generation/review/>  
P9 : <http://www.businessinsider.com/apple-kills-off-ipod-nano-ipod-shuffle-2017-7/>  
P10: <https://www.tesco.com/direct/technology-gaming/ipod/cat3376298.cat?source=others/>



شکل ۵-۱۷ بررسی ۱۰ صفحه درباره ipod

## فهرست مراجع لاتین

- [1] B. Pinkerton - "WebCrawler: finding what people want", University of Washington " - 2000.
- [2] A.K.Sharma, J. P. Guspta , D. P. Agarwal- "Augment Hypertext Documents suitable for parallel crawlers" - Feb2000.
- [4] Gyanendra Kumar, Neelam Duahn - " Page Ranking Based on Number of Visits of Web Pages" - 2011.
- [5] Ali Mohammad Zareh Bidoki, Nasser Yazdani - "DistanceRank: An intelligent ranking algorithm for web pages" - June 2007
- [6] L. Wood – "Programming the Web: the W3C DOM specification" - Jan/Feb 1999
- [7] W3C. Retrieved - "HTML & CSS" - 10 Feb 2012
- [8] Dave Raggett, Arnaud Le Hors, Ian Jacobs - "HTML 4.01 Specification - Abstract". W3C Recommendation. W3C" - 10 Feb 2000.
- [9] Ravindra D.Kale – "Study of web crawler and its different types" - 2014
- [10] Liwei Wang , Yin Li – "Learning Deep Structure Preserving Image text Embeddings" - 2016
- [11] S. Hameed – "An Optimal Algorithm for HTML page Building Process" - 2004
- [12] Luis Gravano – "Query vs. Crawling based lassification of Searchable web DataBase" - 2013
- [13] Yuki Endo , Satoshi Iizuka – "Deep Prop : Extracting Deep Features from a Single Image For Edit Propagation" –Dec 2011
- [14] Dilip Kumar Sharma – "A Comparative Analysis of Web Page Ranking Algorithms" - 2001
- [15] Josip Krapac – "improving web image search results using query relative classifiers" - 2015
- [16] Anthony Rios, Ramakanth Kavuluru – "Convolution Neural Network for Biomedical Text Classification : Application in indexing Biomedical Articles" - 2015
- [17] Springer : 978-1-4471-2502-0 – "Image Acquisition" - 2001
- [18] Larry Brown – "Deep Learning For Image Classification" - 2015
- [19] Aliaksei Serveryn , Alessandro Moschitti - "Learning to Rank Short Text Pairs with Convolution Deep Neural Networks" - 2013
- [20] Isha Mahajan , Sachin Gupta - "Extended Weighted Page Rank Based on VOL by Finding User Activities Time and Page Reading Time, Storing them Directly on Search Engine DataBase Server" - 2017

**Abstract :**

With growing information on the Web, importance of a search engine which enables to provide desired information at low cost and minimum time is vital. An internet user at least once every day uses a search engine to find the required information. Of course, finding a huge-profile list of irrelevant images and results during a search is always a source of inconvenience.

One of the most important use of web search engine is extracting existing information on the web, processes them and categorizes information and eventually displays relevant results to the search subject. The result must be reported in order of similarity and relevance to the subject of the search with a particular priority to the user.

In this thesis, first different type of search engines are reviewed. Then, we present the challenges for a web search engine. We propose a method using deep learning to extract features from images and texts from web pages to categorize the results for quicker access to desired information.

**Keywords :**

search engine, deep learning, image categorization, web crawler, parallel processing, image search engine, innovations





**Shahrood University of Technology**

**Faculty of Electrical Engineering and Robotic**

**M.Sc. Thesis in Robotic Engineering**

**Identification, Classification and Learning Shapes and Images  
on the web Platform by a Search Engine**

By: Mehrdad Salehi

Supervisor:

Dr. Alireza AhmadyFard

September 2017