

## ارزیابی تأثیر منشأ ویژگی‌ها بر میزان دقت تشخیص وب‌هرز توسط الگوریتم‌های طبقه‌بندی

فریبا مستشارنژاد<sup>۱</sup>، سیدرضا کامل<sup>۲</sup>

<sup>۱</sup> دانشجوی کارشناسی ارشد نرم افزار، دانشکده فنی، مهندسی، دانشگاه آزاد اسلامی، مشهد، ایران

<sup>۲</sup> استادیار گروه نرم افزار، دانشکده فنی، مهندسی، دانشگاه آزاد اسلامی، مشهد، ایران

### چکیده

امروزه با توجه به رشد اطلاعات در وب، موتورهای جستجو به عنوان یک ابزار برای ورود به دنیای وب مورد توجه قرار گرفته‌اند. آنها فهرستی از نتایج مرتبط با پرسش کاربر را در اختیار او قرار می‌دهند. از آنجا که اکثر کاربران تنها نتایج صفحه نخست و از آن میان فقط روی سه یا پنج پیوند اولیه را مورد بازدید قرار می‌دهند، حضور یک صفحه در نتایج بالای موتورهای جستجو به معنای بازدیدکننده بیشتر و نیز درآمد بیشتر است. در این میان وب‌هرز یک روش غیرقانونی و غیراخلاقی به منظور افزایش رتبه صفحات اینترنتی توسط فریب الگوریتم‌های موتورهای جستجو می‌باشد. از آنجا که کیفیت نتایج برای موتورهای جستجو اهمیت بسیاری دارد، روش‌های مختلفی برای تشخیص صفحات وب‌هرز ارائه شده است.

تاکنون بررسی‌های فراوانی بر روی مجموعه داده UK-WEBSpam-2007 صورت گرفته و الگوریتم‌های طبقه‌بندی جدید و ترکیبی به نتایج خوبی رسیده‌اند اما هدف ما بررسی عملکرد الگوریتم‌های کلاسیک بر روی این مجموعه داده است و اینکه نشان دهیم از اینگونه الگوریتم‌ها نمی‌توان به صورت خام برای تشخیص وب‌هرز استفاده کرد و روش‌های ترکیبی جدید گزینه مناسب‌تری در این خصوص است. در این مقاله قصد داریم تأثیر روش‌های منتخب طبقه‌بندی را بر میزان تشخیص این صفحات با در نظر گرفتن چگونگی انتخاب ویژگی‌ها، بررسی نماییم. بدین منظور از مجموعه داده UK-WEBSpam-2007 استفاده کرده و ۱۲ روش مختلف طبقه‌بندی را برای تشخیص صفحات وب‌هرز از دیگر صفحات بر روی حالات مختلف انتخاب ویژگی اعمال کردیم. بهترین نتیجه از اعمال الگوریتم‌های طبقه‌بندی بر ترکیب ویژگی‌های مبتنی بر محتوا و ویژگی‌های مبتنی بر پیوند بدست آمد.

### کلمات کلیدی

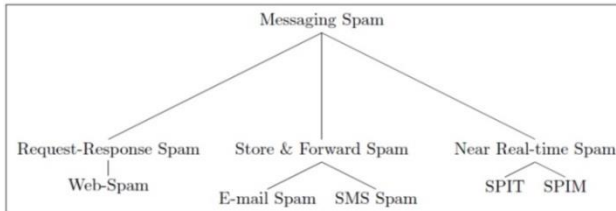
وب‌هرز، طبقه‌بندی، ویژگی‌های مبتنی بر محتوا، ویژگی‌های مبتنی بر پیوند، ویژگی‌های مبتنی بر پیوند تبدیل یافته

این مجموعه داده است و اینکه نشان دهیم از اینگونه الگوریتم ها نمی توان به صورت خام برای تشخیص وب هرز استفاده کرد و روش های ترکیبی جدید گزینه مناسب تری در این خصوص است.

در بخش ۲ ابتدا انواع روش هایی که صفحات وب هرها برای افزایش رتبه خود در موتورهای جستجو استفاده می کنند مورد بررسی قرار گرفته است. در بخش ۳ انواع ویژگی های قابل استخراج از صفحات وب توضیح داده شده است. در بخش ۴ توضیح مختصری از ۱۲ روش منتخب طبقه بندی استفاده شده بیان می شود. در بخش ۵ معیارهای ارزیابی مورد استفاده معرفی شده اند. در بخش ۶ روش های طبقه بندی مورد آزمایش قرار گرفته و نتایج بدست آمده از اعمال آنها بر حالات مختلف انتخاب ویژگی نشان داده شده است و در نهایت در بخش ۷ نتیجه گیری بررسی های انجام شده، عنوان می شود.

## ۲- بخش بندی مقاله

رایج ترین شکل اسپم، هرزنامه<sup>۳</sup> است. جایگاه وب هرها در دسته بندی انواع پیام های ناخواسته در شکل (۱) نشان داده شده است.



شکل (۱): رده بندی سطح بالای هرزنامه ها

در سال ۱۹۹۵، Lycos اولین موتور جستجوی تجاری وارد دنیای وب شد. از همان زمان مبارزه با وب هرها تحت عنوان spamdexing و با استفاده از روش های یادگیری ماشین تبدیل به یکی از مباحث مورد علاقه مجامع دانشگاهی شد [6]. از سال ۲۰۰۵، کارگاه های آموزشی AIRWeb به محلی برای تبادل اطلاعات در این زمینه تبدیل شد [7]. محققان وب هرها را به سه دسته اصلی تقسیم می کنند:

- اسپم های مبتنی بر محتوا
- اسپم های مبتنی بر پیوند
- اسپم های مبتنی بر پنهان سازی صفحه

## ۲-۱- اسپم های مبتنی بر محتوا

اسپم محتوا اولین و گسترده ترین شکل وب هرها است. دلیل گستردگی این شکل از وب هرز در این حقیقت نهفته است که موتورهای جستجو برای رتبه دهی به صفحات وب از مدل های بازبازی اطلاعات مبتنی بر محتوای صفحه از قبیل مدل فضای بردار [8]، BM25 [9] و مدل های آماری زبان [10] استفاده می کنند. به همین دلیل اسپمها نقاط ضعف این مدل ها را تحلیل کرده و از آنها سوءاستفاده می کنند.

اسپم محتوا را به ۵ زیرگروه تقسیم می کنیم [11].

- اسپم بدنه<sup>۴</sup>: در این روش عبارات اسپم در بدنه صفحه قرار می گیرند. به عنوان نمونه اگر یک اسپم بخواهد در صورت جستجوی یک سری کلمات محدود و از پیش تعیین شده به رتبه بالایی دست پیدا کند، کافی است آن کلمات و مشتقات آن را در صفحه تکرار کند.

## ۱- مقدمه

امروزه با توجه به رشد اطلاعات در وب، موتورهای جستجو به عنوان یک ابزار برای ورود به وب مورد توجه قرار گرفته اند. آنها فهرستی از نتایج مرتبط با پرسش کاربر را در اختیار او قرار می دهند. روش قانونی برای افزایش رتبه سایتها در فهرست نتایج موتورهای جستجو افزایش کیفیت صفحات سایتها است، اما این روش زمان بر و پرهزینه است. روش دیگر استفاده از روش های غیرقانونی و غیراخلاقی برای افزایش رتبه در موتورهای جستجو است.

تحقیقاتی که در مرجع [1] و [2] صورت گرفت نشان می دهد که تقریباً ۸۵٪ کاربران تنها نتایج صفحه نخست را مورد بازدید قرار داده و از آن میان فقط روی سه یا پنج پیوند اولیه کلیک می کنند. در نتیجه حضور یک صفحه در نتایج بالای موتورهای جستجو به معنای بازدیدکننده بیشتر و نیز

درآمد بیشتر است. به همین دلیل موتورهای جستجو برای آنکه بهترین نتایج ممکنه را در اختیار کاربران قرار دهند از الگوریتم های بسیار پیچیده ای استفاده می کنند.

در این میان وب هرز یک روش غیرقانونی و غیراخلاقی به منظور افزایش رتبه صفحات اینترنتی توسط فریب الگوریتم های موتورهای جستجو می باشد. هدف وب هرها تغییر رتبه صفحات وب در نتایج جستجو است. به این صورت که رتبه صفحات را بیش از آنچه که باید نشان دهد تا ترجیحاً در میان ۱۰ نتیجه اول ظاهر شود. وب هرها باعث کاهش کیفیت نتایج جستجو و در نتیجه اتلاف وقت کاربر می شوند. افزایش تعداد این صفحات، افزایش تعداد صفحات جستجو شده توسط خزگرها<sup>۱</sup> و مرتب شده توسط اندیس-گزارها<sup>۲</sup> را در پی دارد. این امر موجب اتلاف منابع موتور جستجو و افزایش زمان پاسخگویی به کاربر می شود.

گاهی وب سایت های اسپم به عنوان وسیله ای برای گسترش بدافزار، دسترسی به محتوای ممنوعه و حملات فیشینگ به کار گرفته می شود. به عنوان مثال در مرجع [3] حدود ۱۰۰ میلیون صفحه با استفاده از الگوریتم PageRank رتبه دهی شدند و مشخص شد ۱۱ صفحه از ۲۰ نتیجه ارائه شده توسط موتورهای جستجو، سایت های هرزه نگاری بودند که به خاطر تغییر محتوا و لینک رتبه بالایی گرفته بودند. این امر باعث می شود که شرکت سازنده موتور جستجو مقادیر زیادی از منابع محاسباتی و ذخیره سازی را هدر دهد.

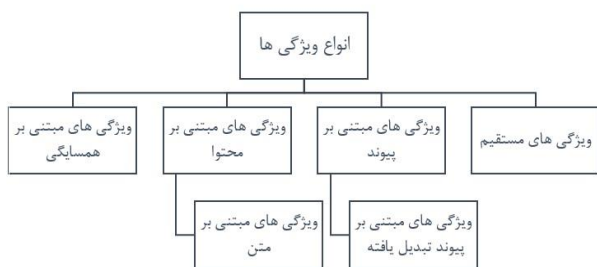
در تحقیقی که در سال ۲۰۰۵ صورت گرفت [4]، تخمین زده شد که وب هرها سالانه حدود ۵۰ میلیون دلار باعث اتلاف منابع مالی می شوند. این مقدار در تحقیقی مشابه [5] که در سال ۲۰۰۹ صورت گرفت به ۱۳۰ میلیون دلار افزایش پیدا کرد. از جمله دلایل این رشد سریع تعداد وب هرها می توان به ساده شدن ابزار ایجاد وب (از قبیل ویکی و وب های رایگان، بسترهای بلاگ نویسی و ...) و کاهش هزینه نگهداری وب سایت اشاره کرد. تشخیص وب هرها برای موتورهای جستجو از اهمیت بسزایی برخوردار بوده و در بازار رقابتی با سایر موتورهای جستجو جزو اولویتهای کاری آنها محسوب می-شود.

تاکنون بررسی های فراوانی بر روی مجموعه داده UK-WEBSpam-2007 صورت گرفته و الگوریتم های طبقه بندی جدید و ترکیبی به نتایج خوبی رسیده اند اما هدف ما بررسی عملکرد الگوریتم های کلاسیک بر روی

- پنهان سازی: در این روش درخواست مشاهده وب سایت بررسی می-شود. اگر از سوی خزشگرهای موتورهای جستجو بود، محتوای متفاوتی نسبت به کاربر عادی به آن نمایش داده می شود.
- تغییر مسیر: وقتی کاربر از کلیک کردن بر روی پیوند نتایج موتور جستجو درخواست مشاهده صفحه را می کند، به صفحه دیگری منتقل می شود.

### ۳- انواع ویژگی ها

به طور کلی ویژگی های قابل استخراج از صفحات وب به ۴ دسته اصلی تقسیم می شوند که در شکل (۲) نشان داده شده است.



شکل (۲): دسته بندی ویژگی های مختلف صفحات وب

دادگان مورد استفاده در این مقاله [12] UK-WEBSPAM-2007 می باشد که معتبرترین مجموعه داده در زمینه شناسایی صفحات وب هرز است. این دادگان شامل مجموعه صفحاتی است که از نتایج یک خزشگر در دامنه .uk بدست آمده است و شامل ۱۰۵/۹ میلیون صفحه از ۱۱۴۵۲۹ میزبان می باشد. از این صفحات سه دسته ویژگی استخراج شده است. ویژگی های مبتنی بر محتوا، ویژگی های مبتنی بر پیوند و ویژگی های مبتنی بر پیوند تبدیل یافته.

### ۳-۱- ویژگی های مبتنی بر محتوا

این دسته از ویژگی ها بر روی محتوای صفحات وب تمرکز دارند. برخی از این ویژگی ها در (جدول ۱) بیان شده است. در زبان انگلیسی برای ویژگی های مبتنی بر محتوا ۹۶ ویژگی استخراج گردیده است.

جدول (۱): برخی از ویژگی های مبتنی بر محتوا

ردیف	نام ویژگی
۱	تعداد کلمات در صفحه
۲	تعداد کلمات در عنوان
۳	میانگین طول کلمه
۴	بخشی از متن قابل مشاهده
۵	نرخ تراکم (فشرده سازی) صفحه
۶	۱۰۰ فراخوانی برتر دیتاست
۷	بخشی از متن لینک داده شده
۸	بخشی از متن لینک داده شده
۹	۲۰۰ دقت (درستی) برتر دیتاست

### ۳-۲- ویژگی های مبتنی بر پیوند

این دسته از ویژگی ها بر روی پیوندهای موجود در صفحات وب تمرکز دارند. برخی از این ویژگی ها در (جدول ۲) بیان شده است. این ویژگی ها هم در صفحه اصلی میزبان<sup>۷</sup> و هم در صفحه ای از میزبان که بالاترین رتبه<sup>۸</sup> را دارد

- اسپم عنوان<sup>۹</sup>: برخی موتورهای جستجو وزن بیشتری برای عنوان اسناد قائل هستند. به همین دلیل اسپمها عبارات اسپم را در عنوان صفحه قرار می دهند.
- اسپم فرابرجسب<sup>۱۰</sup>: توضیحات فرابرجسب های HTML این امکان را برای طراحان صفحات فراهم می کند که توضیحات کوتاهی راجع به صفحه قرار دهند. اگر کلمات نامربوط در این مکان قرار بگیرند و الگوریتم های موتورهای جستجو بر اساس این توضیحات صفحه را مورد بررسی قرار دهند، آنگاه این صفحات به خاطر کلمات نامربوطی که دارند، رتبه بالاتری دریافت می کنند.
- اسپم آدرس اینترنتی: برخی الگوریتم ها از عبارات موجود در URL سایت به منظور محاسبه رتبه استفاده می کنند. به همین دلیل اسپمها برای

- صفحات، لینک های طولانی شامل کلمات عبارت جستجو شده ایجاد می کنند.
- اسپم متن Anchor: بخش Anchor یک لینک حاوی خلاصه ای از متن صفحه ای است که به آن اشاره می کند. بنابراین اسپمها با قرار دادن عبارات دلخواه خود (که معمولاً ارتباطی با محتوای صفحه ندارد) می توانند الگوریتم هایی که برای عبارات این بخش ارزش بیشتری قائل می شوند را فریب دهند و رتبه خود را افزایش دهند.
- با ظهور الگوریتم های رتبه دهی پیوند محور، تا حدود زیادی بر مسئله اسپم-های محتوا محور غلبه شد. هر چند اسپمها نیز به زودی راهکار مقابله با این روش را در ایجاد مزارع پیوند و دیگر روش ها یافتند. با این حال این جنگ و گریز همچنان ادامه دارد.

### ۲-۲- اسپم های مبتنی بر پیوند

- دو دسته عمده از لینک اسپم ها وجود دارد: اسپم لینک های خروجی و اسپم لینک های ورودی.
- اسپم لینک های خروجی: این روش آسان ترین و ارزان ترین روش افزایش رتبه صفحات وب است. از آنجا که اسپم مالک صفحه وب خود است و محتوای آن را می تواند تغییر دهد، به راحتی می تواند لینک هر صفحه ای که بخواهد را در صفحه خود قرار دهد یا حتی محتوای سایت هایی از قبیل DMOZ و Yahoo! Directory که حاوی فهرستی از آدرس صفحات اینترنتی هستند را در صفحه خود کپی کند. این روش الگوریتم های رتبه دهی همچون HITS را نشانه گرفته و در پی افزایش رتبه قطب خود هستند.
- اسپم لینک های ورودی: در این روش اسپمها الگوریتم رتبه دهی PageRank را نشانه گرفته اند و بنابراین سعی در افزایش تعداد لینک های ورودی دارند. در اینجا بسته به اینکه اسپم مالک صفحه باشد یا تنها به آن دسترسی داشته باشد، روش ها متفاوت است.

### ۳-۳- اسپم های مبتنی بر پنهان سازی صفحه

در این روش وب هرزها محتوای متفاوتی نسبت به آنچه کاربر می بیند را به موتورهای جستجو نمایش می دهند. انواع روش های آن عبارتند از:

گروه قرار می‌گیرد. این الگوریتم برای تفکیک دو کلاس از هم، از یک صفحه استفاده می‌کند به طوری که این صفحه از هر طرف بیشترین فاصله را تا هر دو کلاس داشته باشد. نزدیک‌ترین نمونه‌های آموزشی به این صفحه "بردارهای پشتیبان" نام دارند.

• Adaboost: این الگوریتم در واقع یک متا الگوریتم است که بمظور ارتقاء عملکرد، و رفع مشکل رده‌های نامتوازن همراه دیگر الگوریتم‌های یادگیری استفاده می‌شود. در این الگوریتم، طبقه‌بند هر مرحله جدید به نفع نمونه‌های غلط طبقه‌بندی شده در مراحل قبل تنظیم می‌گردد. این الگوریتم از کل مجموعه داده به منظور آموزش هر دسته‌کننده استفاده می‌کند، اما بعد از هر بار آموزش، بیشتر بر روی داده‌های سخت تمرکز می‌کند تا به درستی کلاسه بندی شوند.

• KNN: روش K نزدیک‌ترین همسایه یک گروه شامل K رکورد از مجموعه رکوردهای آموزشی که نزدیک‌ترین رکوردها به رکورد آزمایشی باشند را انتخاب کرده و بر اساس برتری رده یا برچسب مربوط به آن‌ها در مورد دسته رکورد آزمایشی مزبور تصمیم‌گیری می‌نماید. به عبارت ساده‌تر این روش رده‌ای را انتخاب می‌کند که در همسایگی انتخاب شده بیشترین تعداد رکورد متناسب به آن دسته باشند. بنابراین رده‌ای که از همه رده‌ها بیشتر در بین K نزدیک‌ترین همسایه مشاهده شود، به عنوان رده رکورد جدید در نظر گرفته می‌شود.

• MLP: یکی از مرسوم ترین انواع شبکه‌های عصبی، شبکه عصبی پرسپترون چند لایه (MLP) است. هنگام کار با شبکه‌های عصبی MLP با دو مسئله روبه رو هستیم: انتخاب معماری مناسب و انتخاب الگوریتم آموزشی مناسب. معماری مناسب به معنی انتخاب بهینه تعداد لایه‌ها، تعداد نرون‌ها در هر لایه و نوع تابع تحریک هر نرون می‌باشد و معماری بهینه شبکه‌های عصبی مبتنی بر مجموعه داده‌ها و ویژگی‌های آنهاست. الگوریتم‌های آموزشی متنوعی جهت آموزش شبکه‌های عصبی به کار می‌رود. متداول‌ترین الگوریتم آموزشی این شبکه‌ها، الگوریتم پس انتشار خطا می‌باشد. در الگوریتم پس انتشار خطا در هر مرحله مقدار خروجی محاسبه شده جدید، با مقدار واقعی مقایسه شده و با توجه به خطای به دست آمده به اصلاح وزن‌های شبکه پرداخته می‌شود. به نحوی که در انتهای هر تکرار اندازه خطای حاصله کمتر از میزان به دست آمده در تکرار قبلی باشد.

• Random Forest: این دسته‌بند یک روش بر مبنای مدل دسته جمعی است. این مدل، روش یادگیری برای طبقه‌بندی و رگرسیون است که با ساخت بسیاری از درخت‌های تصمیم‌گیری در زمان آموزش و انتخاب بهترین درخت از میان درختان تولید شده کار می‌کند. مدل دسته جمعی بر مبنای دقت و میزان اهمیت متغیرها کار می‌کند.

• Classification & Regression Trees (CART): درخت تصمیم‌گیری CART یک روش تقسیم بندی بازگشتی باینری است. در این روش درختان به حداکثر رشد خود می‌رسند و سپس اصلاح می‌شوند که کار پیچیده‌ای است. هدف این مکانیزم تولید تنها یک درخت نیست بلکه تولید یک سری درختان اصلاح شده تو در تو است که همه

محاسبه می‌شوند. در زبان انگلیسی برای ویژگی‌های مبتنی بر پیوند ۴۱ ویژگی استخراج گردیده است.

#### جدول (۲): برخی از ویژگی‌های مبتنی بر پیوند

ردیف	نام ویژگی
۱	تعداد پیوندهای خروجی صفحه اصلی
۲	تعداد پیوندهای ورودی صفحه اصلی
۳	نسبت تعداد پیوندهای خروجی داخلی به کل پیوندهای خروجی
۴	نسبت پیوند های خروجی صفحات دیگر به صفحه اصلی
۵	صفحاتی که به فاصله ۲ پیوند با صفحه اصلی هستند
۶	صفحاتی که به فاصله ۳ پیوند با صفحه اصلی هستند
۷	صفحاتی که به فاصله ۴ پیوند با صفحه اصلی هستند
۸	تعداد میزبان های متفاوت که در فاصله ۲ پیوند با صفحه اصلی هستند

#### ۳-۳- ویژگی‌های مبتنی بر پیوند تبدیل یافته

این دسته از ویژگی‌ها شامل تبدیلات عددی ساده و ترکیب ویژگی‌های مبتنی بر پیوند می‌شود. برخی از این ویژگی‌ها در (جدول ۳) بیان شده است. در زبان انگلیسی برای ویژگی‌های مبتنی بر پیوند تبدیل یافته ۱۳۸ ویژگی استخراج گردیده است.

#### جدول (۳): برخی از ویژگی‌های مبتنی بر پیوند تبدیل یافته

ردیف	نام ویژگی
۱	لگاریتم پیوندهای ورودی به صفحه اصلی
۲	لگاریتم پیوندهای خروجی از صفحه اصلی
۳	لگاریتم میانگین پیوندهای ورودی از صفحات خروجی به صفحه اصلی
۴	نسبت بین Indegree/PageRank
۵	نسبت بین TrustRank/PageRank
۶	لگاریتم ویژگی‌های مختلف

#### ۴- الگوریتم‌های طبقه‌بندی

در این مقاله ۱۲ الگوریتم طبقه‌بندی مورد استفاده قرار گرفته است که اکثر آنها جز برترین الگوریتم‌های طبقه‌بندی شناخته شده می‌باشند [13] و به شرح ذیل معرفی می‌گردند:

- C4.5: در روش درخت تصمیم، نمونه‌ها به یک الگوریتم C4.5 که یکی از درخت‌های تصمیم‌گیری است اعمال می‌شوند، سپس درخت هرس شده حاصل از الگوریتم C4.5 را گرفته و کلیه ویژگی‌هایی که در آن وجود دارد را بعنوان جواب مساله باز می‌گردانیم. الگوریتم C4.5 تکمیل شده الگوریتم ID3 می‌باشد.
- Naive Bayes: تئوری بیز یکی از روش‌های آماری برای رده‌بندی به شمار می‌آید. در این روش کلاس‌های مختلف، هر کدام به شکل یک فرضیه دارای احتمال در نظر گرفته می‌شوند. هر رکورد آموزشی جدید، احتمال درست بودن فرضیه‌های پیشین را افزایش و یا کاهش می‌دهد و در نهایت، فرضیاتی که دارای بالاترین احتمال شوند، به عنوان یک کلاس در نظر گرفته شده و برچسبی بر آن‌ها زده می‌شود [14].
- SVM: ماشین بردار پشتیبان یک گروه از الگوریتم‌های طبقه‌بندی نظارت شده هستند، که پیش‌بینی می‌کند یک نمونه در کدام کلاس یا

$$TPR(Recall) = \frac{TP}{(TP + FN)} \quad (2)$$

$$PPR(Precision) = \frac{TP}{(TP + FP)} \quad (3)$$

$$Accuracy = \frac{(TP + TN)}{(P + N)} \quad (4)$$

$$F - Measure = \frac{2TP}{(2TP + FP + FN)} \quad (5)$$

در روابط مذکور، TP تعداد موارد طبقه‌بندی مثبت صحیح (ماشین یادگیری به درستی طبقه‌بندی می‌کند)، TN تعداد موارد طبقه‌بندی منفی صحیح (ماشین یادگیری به درستی طبقه‌بندی می‌کند)، FP تعداد موارد طبقه‌بندی مثبت نادرست (ماشین یادگیری به درستی طبقه‌بندی نمی‌کند)، FN تعداد موارد طبقه‌بندی منفی نادرست (ماشین یادگیری به درستی طبقه‌بندی نمی‌کند)، سطح زیر منحنی ROC است که عددی بین صفر و یک است و هرچه به یک نزدیکتر باشد بهتر است.

## ۶- پیاده‌سازی و نتایج

الگوریتم‌های اشاره شده در بخش ۴، توسط نرم افزار weka [14] بر روی داده‌گان مورد بحث اعمال شده است. مجموعه داده مورد استفاده WEBSpAM-UK2007 می‌باشد که در آن مجموعه تست از مجموعه

آموزشی جدا است. در این مقاله سعی شده است علاوه بر استفاده از ویژگی‌های مبتنی بر محتوا و ویژگی‌های مبتنی بر پیوند از دسته سوم ویژگی‌ها که مبتنی بر پیوند تبدیل یافته و همچنین ترکیب ویژگی‌ها استفاده شود. ۹۶ ویژگی مبتنی بر محتوا و ۴۱ ویژگی مبتنی بر پیوند و ۱۳۸ ویژگی مبتنی بر پیوند تبدیل یافته در مجموعه داده ذکر شده مد نظر قرار گرفته است.

جدول (۴): نرخ تشخیص وب هرز با ویژگی‌های مبتنی بر محتوا

روش	Accuracy	Precision	Recall	F-Measure	ROC Area
C4.5	92.25%	0.289	0.23	0.256	0.57
SVM	94.20%	0	0	0	0.5
AdaBoost	94.20%	0	0	0	0.782
KNN	92.97%	0.36	0.274	0.312	0.622
Naïve Bayes	14.32%	0.058	<b>0.903</b>	0.109	0.653
CART	94.20%	0	0	0	0.5
Logistic Regression	93.43%	0.317	0.115	0.169	0.737
Logit Boost	94.20%	0.5	0.097	0.163	0.78
Model Tree	94.61%	0.643	0.159	0.255	<b>0.794</b>
MLP	<b>94.66%</b>	<b>0.655</b>	0.168	0.268	0.758
Random Forest	94.61%	0.583	0.248	<b>0.348</b>	0.788
LAD Tree	94.35%	0.579	0.097	0.167	0.762

جدول (۵): نرخ تشخیص وب هرز با ویژگی‌های مبتنی بر پیوند.

روش	Accuracy	Precision	Recall	F-Measure	ROC Area
Naïve Bayes	92.46%	0.133	0.049	0.072	0.643
SVM	94.06%	0	0	0	0.5
Logistic Regression	<b>93.82%</b>	0.222	0.016	0.031	0.619
MLP	92.41%	0.075	0.025	0.037	0.037
KNN	89.93%	0.117	<b>0.107</b>	0.112	0.528
AdaBoost	94.06%	0	0	0	0.702
Logit Boost	94.06%	0	0	0	<b>0.734</b>
C4.5	93.72%	0.231	0.025	0.044	0.609
LAD Tree	93.48%	0.269	0.057	0.095	0.71
Model Tree	94.06%	0	0	0	0.609
Random Forest	93.38%	<b>0.281</b>	0.074	<b>0.117</b>	0.615

آن درختان بهینه داوطلب هستند و انتخاب درخت بهینه تنها پس از ارزشیابی داده‌ها صورت می‌گیرد.

- logistic regression: رگرسیون لجستیک یکی از ابزارهای آماری است که به منظور مدل‌سازی و تحلیل داده‌ها از آن استفاده می‌شود. رگرسیون لجستیک دارای شکل کلی زیر است:

$$\log\left(\frac{\pi}{1-\pi}\right) = \alpha + \sum_{i=1}^k \beta_i \cdot x_i \quad (1)$$

در این مدل،  $\pi$  احتمال تعلق فرد به سطح اول متغیر وابسته است،  $x_i$  متغیر مستقل  $i$  ام و  $\beta_i$  ضریب برآورد شده مدل برای متغیر مستقل  $i$  ام است. از مزایای استفاده از مدل رگرسیون لجستیک علاوه بر مدل‌سازی مشاهده‌ها، امکان پیش‌بینی احتمال تعلق هر فرد به هر یک از سطوح متغیر وابسته و همچنین امکان محاسبه‌ی مستقیم نسبت شانس با استفاده از ضرایب مدل است.

- Logit boost: الگوریتم افزایشی (Logit boost) به طریقی عمل می‌کند که همانند یک روش گروهی خود را برای انجام هدف اصلی تجهیز می‌کند و به لحاظ مفهوم شبیه به Bagging است. افزایشی، آموزش دهنده ضعیف بعدی را بر اساس خطاهای آموزش دهنده قبلی آموزش می‌دهد. در شرایط عادی، افزایشی نسبت به Bagging عملکرد بهتری دارد.
- Logistic Model Tree: یا به اختصار LMT یک مدل طبقه‌بندی تحت نظارت است که از ترکیب رگرسیون لجستیک و درخت تصمیم استفاده می‌کند. ایده این روش به این صورت است که به جای اینکه در برگ‌های درخت از مقادیر ثابت استفاده شود، با استفاده از الگوریتم LogitBoost برای هر یک از آنها یک مدل رگرسیون خطی ایجاد شود. این الگوریتم به جای استفاده از مقادیر تصادفی، در هر فراخوانی از نتایج گره والد خود برای شروع استفاده می‌کند. سپس این گره‌ها با استفاده از الگوریتم C4.5 تقسیم شده و در نهایت هرس می‌شوند. همچنین برای پیدا کردن تعداد تکرارهای مناسب برای اجرای LogitBoost از cross validation استفاده می‌شود.
- Least Absolute Deviation Tree: درختان رگرسیون LAD از به حداقل رساندن انحراف مطلق (Absolute Deviation) بین مقدار پیش‌بینی شده و مقدار واقعی، به عنوان معیار انتخاب استفاده می‌کنند. استفاده از این معیار باعث می‌شود که درخت نسبت به حضور داده‌های پرت مقاوم تر باشند.

## ۵- معیارهای ارزیابی

پس از اعمال الگوریتم‌های طبقه‌بندی بر ویژگی‌های مورد نظر شاخص‌های حساسیت یا نرخ یادآوری، نرخ پیش‌بینی صحیح، دقت، معیار  $F$  و سطح زیر منحنی قابل اندازه‌گیری خواهند بود.

LAD Tree	93.84%	0.387	0.106	0.167	0.773
model tree	94.51%	0.607	0.15	0.241	0.75
random forest	<b>94.97%</b>	0.703	0.23	<b>0.347</b>	0.753
CART	94.10%	0.429	0.053	0.094	0.599

CART	94.06%	0	0	0	0.5
------	--------	---	---	---	-----

جدول (۶): نرخ تشخیص وب هرز با ویژگی های مبتنی بر پیوند تبدیل یافته

روش	Accuracy	Precision	Recall	F-Measure	ROC Area
Naïve Bayes	82.48%	0.128	<b>0.336</b>	<b>0.186</b>	0.683
Svm	<b>94.06%</b>	0	0	0	0.5
logistic regression	93.33%	0.105	0.016	0.028	0.716
MLP	91.97%	0.078	0.033	0.046	0.595
KNN	90.85%	0.17	0.139	0.153	0.548
Adaboost	<b>94.06%</b>	0	0	0	0.673
logit boost	93.97%	0	0	0	<b>0.723</b>
C4.5	92.75%	0.154	0.049	0.075	0.631
LAD Tree	93.63%	0.263	0.041	0.071	0.696
model tree	<b>94.06%</b>	0	0	0	0.629
random forest	93.53%	<b>0.28</b>	0.057	0.095	0.634
CART	94.06%	0	0	0	0.5

جدول (۱۰): نرخ تشخیص وب هرز با ترکیب هر سه دسته ویژگی.

روش	Accuracy	Precision	Recall	F-Measure	ROC Area
Naïve Bayes	69.68%	0.115	<b>0.636</b>	0.195	0.678
Svm	94.21%	0	0	0	0.5
logistic regression	92.20%	0.266	0.196	0.226	0.709
MLP	93.88%	0.35	0.065	0.11	0.612
KNN	91.72%	0.189	0.131	0.155	0.548
Adaboost	94.53%	<b>0.875</b>	0.065	0.122	0.787
logit boost	<b>94.59%</b>	0.733	0.103	0.18	0.782
C4.5	91.61%	0.25	0.224	<b>0.236</b>	0.511
LAD Tree	93.94%	0.414	0.112	0.176	<b>0.788</b>
model tree	94.37%	0.556	0.14	0.224	0.759
random forest	94.26%	0.517	0.14	0.221	0.755
CART	93.88%	0.375	0.084	0.137	0.593

جدول (۷): نرخ تشخیص وب هرز با ترکیب ویژگی های مبتنی بر پیوند و مبتنی بر پیوند تبدیل یافته

روش	Accuracy	Precision	Recall	F-Measure	ROC Area
Naïve Bayes	84.24%	0.134	<b>0.304</b>	<b>0.186</b>	0.689
Svm	94.10%	0	0	0	0.5
logistic regression	92.97%	0.133	0.035	0.055	<b>0.709</b>
MLP	92.25%	0.154	0.07	0.096	0.628
KNN	91.12%	0.17	0.13	0.148	0.545
Adaboost	94.10%	0	0	0	0.71
logit boost	93.99%	0	0	0	0.724
C4.5	93.63%	0.154	0.017	0.031	0.636
LAD Tree	<b>93.69%</b>	0.25	0.035	0.061	0.702
model tree	94.10%	0	0	0	0.625
random forest	93.38%	<b>0.281</b>	0.078	0.122	0.644
CART	94.10%	0	0	0	0.5

می توان عنوان داشت که Naïve Bayes در صورتی که بر روی ویژگی های محتوا محور اعمال شود (جدول ۴)، بالاترین نرخ Recall (۹۰٪) را در میان دیگر روش ها دارد. این نرخ بالا به بهای نرخ بسیار اندک precision (۶٪) و دقت (۱۴٪) به دست آمده است. بالاترین نرخ precision و F-Measure نیز از اعمال روش Random Forest بر روی ترکیب ویژگی های محتوا و پیوند محور به دست آمده است (جدول ۹) که به ترتیب ۹۵٪ و ۳۵٪ می باشد اما در قبال نرخ Recall ۲۳٪ به دست آمده است.

در مورد کارایی ویژگی های انتخاب شده نیز با نگاهی به نتایج می توان مشاهده کرد که ویژگی های پیوند محور به خودی خود دارای ارزش تفکیک کنندگی بالایی نیستند، اما در کنار ویژگی های محتوا محور منجر به افزایش دقت آن می شوند و از آنجا که توزیع داده در این دادگان یکسان نیست و تعداد نمونه های کلاس صفحات نرمال بیش از ۱۷ برابر نمونه های کلاس صفحات اسپم است، حتی اگر کلیه صفحات اسپم به اشتباه نرمال تشخیص داده شوند، دقت طبقه بند بیش از ۹۴٪ خواهد شد.

## ۷- نتیجه گیری

در این مقاله انواع روش هایی که وب هرزها برای بالا بردن رتبه خود در موتورهای جستجو استفاده می کنند را معرفی کردیم. سپس انواع ویژگی های استخراج شده از صفحات وب به منظور تشخیص وب هرزها را بررسی کردیم و در نهایت ۱۲ روش طبقه بندی منتخب را در جستجوی بهترین نتیجه، بر روی مجموعه داده WEBSpam-UK2007 پیاده سازی کردیم. از آنجا که ویژگی های مورد استفاده نقش مهمی در بالا بردن دقت طبقه بند دارند، این روش ها به تفکیک بر روی ویژگی های مبتنی بر محتوا، مبتنی بر پیوند و مبتنی بر پیوند تبدیل یافته و همچنین ترکیب این ویژگی ها آزمایش شد. بر اساس بررسی های صورت گرفته، توزیع داده در این دادگان یکسان نیست و تعداد نمونه های کلاس صفحات نرمال بیش از ۱۷ برابر نمونه های کلاس صفحات اسپم است؛ بنابراین حتی اگر کلیه صفحات اسپم به اشتباه نرمال تشخیص داده شوند، دقت طبقه بند بیش از ۹۴٪ خواهد شد. به همین دلیل نتایج بدست آمده کمی دور از انتظار است و برخی از الگوریتم ها مثل SVM به نتایج خوبی نرسیدند و در نتیجه Accuracy به تنهایی ملاک خوبی برای

جدول (۸): نرخ تشخیص وب هرز با ترکیب ویژگی های مبتنی بر محتوا و مبتنی بر پیوند تبدیل یافته

روش	Accuracy	Precision	Recall	F-Measure	ROC Area
Naïve Bayes	43.89%	0.075	<b>0.761</b>	0.136	0.668
Svm	94.20%	0	0	0	0.5
logistic regression	92.92%	0.333	0.221	<b>0.266</b>	0.725
MLP	93.43%	0.333	0.133	0.19	0.691
KNN	91.68%	0.198	0.142	0.165	0.553
Adaboost	94.56%	<b>0.889</b>	0.071	0.131	0.787
logit boost	94.61%	0.75	0.106	0.186	0.779
C4.5	91.63%	0.245	0.212	0.227	0.502
LAD Tree	93.94%	0.419	0.115	0.181	<b>0.788</b>
model tree	94.46%	0.593	0.142	0.229	0.762
random forest	<b>94.71%</b>	0.692	0.159	0.259	0.76
CART	93.89%	0.375	0.08	0.131	0.599

جدول (۹): نرخ تشخیص وب هرز با ترکیب ویژگی های مبتنی بر محتوا و مبتنی بر پیوند

روش	Accuracy	Precision	Recall	F-Measure	ROC Area
Naïve Bayes	22.64%	0.061	<b>0.858</b>	0.114	0.674
Svm	94.20%	0	0	0	0.5
logistic regression	93.33%	0.302	0.115	0.167	0.733
MLP	94.51%	0.594	0.168	0.262	0.696
KNN	93.07%	0.38	0.31	0.341	0.639
Adaboost	94.56%	<b>0.889</b>	0.071	0.131	<b>0.787</b>
logit boost	94.56%	0.684	0.115	0.197	0.779
C4.5	92.30%	0.282	0.212	0.242	0.526

### زیرنویس‌ها

- <sup>1</sup> Crawler
- <sup>2</sup> Indexer
- <sup>3</sup> Email Spam
- <sup>4</sup> Body Spamming
- <sup>5</sup> Title Spamming
- <sup>6</sup> Meta Tag
- <sup>7</sup> Home Page
- <sup>8</sup> PageRank
- <sup>9</sup> Recall
- <sup>10</sup> Precision
- <sup>11</sup> Accuracy
- <sup>12</sup> Accuracy
- <sup>13</sup> AUC

ارزیابی روشهای طبقه بندی نخواهد بود. در نهایت مشاهده می شود که ویژگی ها و روش های کلاسیک طبقه بندی راه حل مناسبی برای تشخیص صفحات وب هرز نیست و می بایست از روش و ویژگی های جدیدی برای حل آن استفاده کرد.

### مراجع

- [1] C. Silverstein, H. Marais, M. Henzinger, and M. Moricz, "Analysis of a very large web search engine query log", SIGIR Forum, 33(1), Sept. 1999.
- [2] T. Joachims, L. Granka, B. Pan, H. Hembrooke, and G. Gay, "Accurately interpreting click through data as implicit feedback". In Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR'05, Salvador, Brazil, 2005.
- [3] N. Eiron, K. S. McCurley, and J. A. Tomlin, "Ranking the web frontier", In Proceedings of the 13th International Conference on World Wide Web, WWW'04, New York, NY, 2004.
- [4] D. Fetterly, M. Manasse, and M. Najork, "Detecting phrase-level duplication on the world wide web", In Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR'05, Salvador, Brazil.
- [5] R. Jennings, "Cost of spam is fattening our 2009 predictions", Ferris Research, 2009.
- [6] Davison, B.D., "Recognizing nepotistic links on the web. Artificial Intelligence for Web Search", 2000:p. 23-28.
- [7] Najork, M., Web Spam Detection. Encyclopedia of Database Systems, 1: p. 3520-3523, 2009.
- [8] G. Salton, A. Wong, and C. S. Yang, "A vector space model for automatic indexing". Common. ACM, Vol.18, Nov. 1975.
- [9] S. Robertson, H. Zaragoza, and M. Taylor, "Simple bm25 extension to multiple weighted fields". In Proceedings of the Thirteenth ACM International Conference on Information and Knowledge Management, CIKM'04, Washington, D.C., 2004.
- [10] C. Zhai, "Statistical Language Models for Information Retrieval". Now Publishers Inc., Hanover, MA, USA, 2008.
- [11] Gyongyi, Z. and H. Garcia-Molina, WebSpamTaxonomy, in First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb 2005). 2005.
- [12] Carlos Castillo, Debora Donato, Luca Becchetti, Paolo Boldi, Stefano Leonardi, Massimo Santini and Sebastiano Vigna, "A Reference collection for Web Spam", ACM SIGIR Forum, Vol. 40, No. 2, pp. 11-14, 2006.
- [13] XindongWu, Vipin Kumar, J. Ross Quinlan, Joydeep Ghosh, Qiang Yang, Hiroshi Motoda, Geoffrey J. McLachlan, Angus Ng, Bing Liu, Philip S. Yu, Zhi-Hua Zhou, Michael Steinbach, David J. Hand, Dan Steinberg, "Top 10 algorithms in data mining", Springer, 1-37, 2008.
- [14] George H. John, Pat Langley, "Estimating Continuous Distributions in Bayesian Classifiers". In: Eleventh Conference on Uncertainty in Artificial Intelligence, San Mateo, 338-345, and 1995.
- [15] <http://www.cs.waikato.ac.nz/~remco/weka.pdf>.