

## یک روش یادگیری تجمعی جدید در اسپم فیلترینگ متنی

### سمن مثقالی<sup>۱</sup>، جواد عسگری<sup>۲</sup>

<sup>۱</sup> کارشناسی ارشد مهندسی کامپیوتر هوش مصنوعی، مرکز آموزش‌های الکترونیکی دانشگاه صنعتی اصفهان، اصفهان، ایران.

<sup>۲</sup> دکتری مهندسی برق - کنترل، هیأت علمی دانشگاه صنعتی اصفهان، دانشکده مهندسی برق و کامپیوتر، اصفهان، ایران.

نام نویسنده مسئول:

جواد عسگری

تاریخ دریافت: ۱۳۹۹/۷/۲

تاریخ پذیرش: ۱۳۹۹/۹/۴

### چکیده

اسپم‌های متنی، پیام‌هایی ناخواسته هستند که امروزه به صورت ایمیل یا پیام کوتاه دریافت می‌شوند. با توجه به افزایش حجم اسپم‌های تولیدی و با توجه به میزان ایمیل‌های متفاوتی نظیر شغلی، شخصی و سایر موارد که به طور روزانه دریافت می‌کنیم، بسیار مهم است که بتوانیم ایمیل‌های اسپم را شناسایی نماییم. از این رو هر پلتفرم ارسال و دریافت پیامی باید مجهز به یک سیستم تشخیص اسپم قوی باشد تا بتواند اسپم‌ها را در بدو ورود تشخیص داده و فیلتر کند. امروزه روش‌های متعددی برای تشخیص اسپم ارائه شده و اغلب در تشخیص اسپم-ها موفق عمل می‌کنند. اما چالشی که در این حوزه هست، وجود False Positive (FP) در تشخیص‌ها است. یعنی پیام‌های مشروع که به اشتباه به عنوان اسپم شناخته شده و فیلتر می‌شوند. در این مقاله یک روش جدید یادگیری تجمعی به منظور اسپم فیلترینگ ارائه می‌شود. در این روش برخلاف دیگر روش‌های یادگیری تجمعی که زیرمجموعه‌ها را بدون توجه به مکان نمونه‌ها انتخاب می‌کنند، هر زیرمجموعه از مکان مشخصی انتخاب می‌شود و برای تعیین برجسب نهایی متن، بین یادگیرنده‌هایی که توسط زیرمجموعه‌ها آموزش داده شده‌اند، رأی-گیری اکثریت برگزار می‌شود. نتایج نشان می‌دهند روش پیشنهادی به طور قابل ملاحظه‌ای دقت اسپم فیلترینگ را افزایش داده و FP را کاهش می‌دهد.

**واژگان کلیدی:** اسپم فیلترینگ، یادگیری ماشین، یادگیری تجمعی، ایمیل، پیام کوتاه.

## مقدمه

ایمیل‌ها، شیوه‌ی معمول تعداد زیادی از حملات اینترنتی هستند. بیشتر این حملات از نوع هرزنامه هستند. فیلترهای متنوعی مانند فیلتر بیز، فیلترهای مبتنی بر مجموع و فیلترهای مبتنی بر یادگیری ماشینی برای تشخیص هرزنامه مورد استفاده قرار می‌گیرند. بنابراین مدیریت ایمیل به دلیل مورد سوء استفاده قرار گرفتن، مشکلی مهم و روز افزون برای افراد و سازمان‌ها شده است. هرزنامه اصولاً به عنوان ارسال ایمیل دسته‌ای ناخواسته، شناخته می‌شود. به بیان دیگر ایمیل هرزنامه‌ای است ناخواسته برای چند دریافت کننده. یک تعریف معمول، محدود به ایمیل تبلیغاتی ناخواسته است. تعریفی که ارسال‌های غیر تبلیغاتی مانند اعلانات سیاسی و مذهبی، حتی ناخواسته، را به عنوان هرزنامه در نظر نمی‌گیرد. ایمیل تاکنون معمول‌ترین فرم هرزنامه اینترنتی بوده است. بر طبق داده تخمین زده شده در تحقیق فریس، ۱۵ تا ۲۰ درصد ایمیل سازمان‌های شرکتی آمریکا هرزنامه است. نیمی از کاربران، روزانه تعداد ۱۰ ایمیل یا بیشتر هرزنامه دریافت می‌کنند. در حالی که تعدادی از آن‌ها تا صدها ایمیل ناخواسته دریافت می‌کنند. به عنوان مثال، می‌توان گفت امروزه برخی از سایت‌های تجاری نیز مورد هدف تولیدکنندگان اسپم‌ها قرار می‌گیرند و از آن جایی که انتشار و انعکاس نظرات واقعی از سوی کاربران در خصوص استفاده از محصولات یک شرکت می‌تواند در بردارنده اطلاعاتی ارزشمند برای دیگر مشتریان باشد، لذا شناسایی اسپم‌ها و تشخیص نظرات جعلی در زمان کوتاه امری مهم تلقی می‌گردد.

یکی از روش‌های رایج و مؤثر تشخیص هرزنامه، استفاده از الگوریتم‌های یادگیری ماشینی است. به عنوان یکی از شاخه‌های وسیع و پرکاربرد هوش مصنوعی و علوم کامپیوتر، یادگیری ماشینی به تنظیم و اکتشاف شیوه‌ها و الگوریتم‌هایی می‌پردازد که براساس آن‌ها رایانه‌ها و سامانه‌ها توانایی یادگیری پیدا می‌کنند. یادگیری تحت نظارت، یک روش عمومی در یادگیری ماشینی است که در آن به یک سیستم، مجموعه‌ای از جفت‌های ورودی - خروجی ارائه شده و سیستم تلاش می‌کند تا تابعی از ورودی به خروجی را فرا گیرد. به روش‌های عمومی یادگیری ماشینی اغلب روش‌های طبقه‌بندی نیز گفته می‌شود. برخی از این روش‌ها عبارتند از درخت تصمیم، ماشین بردار پشتیبان، نزدیک‌ترین همسایه و شبکه‌های عصبی مصنوعی [۱-۴].

تحقیقات نشان می‌دهند با ترکیب روش‌های طبقه‌بندی پایه و اجرای رأی‌گیری وزن‌دار بین آن‌ها، روش ترکیبی می‌تواند دقت بالاتری نسبت به طبقه‌بندی بیزین،  $k$ - نزدیک‌ترین همسایه و ماشین بردار پشتیبان داشته باشد. این تحقیق بر روی مجموعه داده‌های نظرات متنی در وبسایت آمازون آزمایش شده است و نتایج حاکی از بهبود طبقه‌بندی به کمک یادگیری دسته‌جمعی بوده است [۷]. در این مقاله نیز یک روش مبتنی بر یادگیری تجمعی به منظور فیلترینگ اسپم متنی ارائه می‌شود. ادامه‌ی مقاله به صورت زیر تنظیم شده است.

در بخش ۲ به معرفی تعدادی از مطالعات مرتبط پرداخته می‌شود. در بخش ۳ روش پیشنهادی شرح داده می‌شود. در بخش ۴ نتایج آزمایش‌ها ارائه شده و در نهایت در بخش ۵ به نتیجه‌گیری پرداخته می‌شود.

## مطالعات مرتبط

اولین ابزار ریاضیاتی که برای سیستم‌های فیلتر هرزنامه مورد استفاده قرار گرفت، الگوریتم بیز بود که توسط پژوهشگران زیادی از آن بهره برده شد [۵-۸]. طبقه‌بند بیز متکی بر قضیه‌ی معروف بیز است که اولین مقالات در مورد آن را می‌توان در سال‌های نزدیک ۱۹۶۰ یافت [۹]. در [۱۰] در برخی فیلترها مانند «قاتل هرزنامه» از احتمالات تداخلی پیشنهادی توسط آر فیش در سال ۱۹۵۰، استفاده شده است. برای کشف هرزنامه پیشنهاد داده شد که علاوه بر احتمال «هرزنامه بودن» ایمیل، احتمال «مجاز بودن» ایمیل نیز محاسبه شود.

موارد دیگری که برای تشخیص هرزنامه ارائه شدند برنامه‌های رتبه‌بندی صفحه با زنجیره مارکف و مدل مارکف پنهان بودند که در [۱۱، ۱۲] معرفی شدند. همچنین روش تخمین پیچیدگی کلموگروف در [۱۳] قابل مشاهده است. یک روش کاملاً جدید نیز در مدل تجزیه و تحلیل دیجیتالی متن ایمیل ارائه شده است [۱۴]. در این روش، ایمیل پس از اعمال روش‌های پردازش دیجیتالی و تعریف احتمالات مثبت کاذب، به عنوان یک سیگنال  $x(n)$  در نظر گرفته می‌شود. روش‌های تجزیه و

تحلیل به صورت خوشه‌بندی برای تشخیص ایمیل هرزنانه نیز در [۱۵-۱۸] مورد استفاده قرار گرفتند. از سال ۲۰۰۹، با شروع از [۱۹] می‌توان روش‌های داده‌کاوی و همچنین روش‌های ترکیبی مانند فیلتر مشارکتی و فیلتر محتوا محور را مشاهده کرد. در پروژه‌های تحت عنوان «دیگ عسل»، در [۲۳،۲۴] روش خوشه‌بندی طیفی بر روی مجموعه‌ای از پیام‌های هرزنانه با هدف ردیابی شبکه‌های اجتماعی سازندگان هرزنانه اجرا شد. در این پروژه، یک شبکه از سازندگان هرزنانه به عنوان گرافی از گره‌ها و یال‌ها نمایش داده شد که هر گره معادل با یک سازنده‌ی هرزنانه بوده و هر یال بین دو گره نیز ارتباط دو سازنده‌ی هرزنانه را نشان می‌داد.

### روش پیشنهادی

روش‌های مبتنی بر یادگیری تجمعی مانند Bagging، اگرچه باعث افزایش دقت الگوریتم‌های پایه می‌شوند ولی در انتخاب زیرمجموعه‌های آموزشی توجهی به مکان نمونه‌ها ندارند و در بسیاری از حالات، زیرمجموعه‌های انتخابی صرفاً نماینده‌های ضعیفی از کل مجموعه داده هستند. گاهی با رأی‌گیری‌های وزن‌دار تأثیر زیرمجموعه‌های ضعیف کمتر می‌شود. اما در این مقاله ایده این است که زیرمجموعه‌های آموزشی را از وضعیت‌های مختلف انتخاب کنیم. به این ترتیب می‌توانیم اطمینان حاصل کنیم که مدل‌ها پس از یادگیری تفاوت زیادی با هم دارند و همین باعث خواهد شد که نظرات متفاوتی برای برچسب نمونه‌ها حاصل شود.

در نهایت نتایج نشان خواهند داد که یادگیری تجمعی پیشنهادی کارایی بالاتری از روش‌های مرسوم مانند Bagging و Boosting دارد. لازم به ذکر است که مراحل متن‌کاوی و پردازش زبان طبیعی در این مقاله مد نظر قرار نگرفته است. از این رو ورودی‌های مسئله به صورت بردارهای ویژگی فرض می‌شوند. اینکه چگونه از متن، بردارهای ویژگی استخراج شده است، خارج از حوزه‌ی بررسی این مقاله است.

### مراحل یادگیری تجمعی پیشنهادی

ابتدا توزیع نمونه‌های هر کلاس در هر بُعد از داده‌ها به دست می‌آید. در صورتی که نمونه‌های آموزشی دارای  $N$  بُعد باشند، با توجه به اینکه مسائل از نوع دو-کلاسی است، در این مرحله  $2 \times N$  توزیع به دست می‌آید. سپس برای هر توزیع مقدار میانگین محاسبه می‌شود. برای هر زیرمجموعه در هر بُعد به صورت مجزا یک طبقه‌بند معمولی مانند ماشین بردار پشتیبان یا  $k$ -نزدیک‌ترین همسایه اجرا و مدل‌های طبقه‌بندی پس از آموزش استخراج می‌شود.

### نحوه ارزیابی روش

برای ارزیابی عملکرد روش ارائه شده در این تحقیق، از الگوریتم 10-fold cross-validation استفاده شده است. در این روش در پایگاه‌داده، ۱۰ درصد از نمونه‌ها را برای آزمون و ۹۰ درصد مابقی را به‌عنوان داده‌های آموزش به الگوریتم یادگیر می‌دهیم و دقت ۱۰ درصد آزمون را اندازه می‌گیریم. سپس ۱۰ درصد دیگر از نمونه‌ها را انتخاب کرده و ۹۰ درصد مابقی را به‌عنوان داده‌های آموزش در نظر می‌گیریم. این عملیات را ۱۰ بار تکرار کرده تا تمام نمونه‌ها هم به‌عنوان آموزش و هم به‌عنوان آزمون در نظر گرفته شوند. از دقت‌های به دست آمده از هر دور میانگین گرفته و به‌عنوان دقت هر الگوریتم معرفی می‌نماییم. به منظور ارزیابی نتایج آزمایش بر روی مجموعه‌داده‌ی تحقیق از سه معیار صحت متوسط، دقت و فراخوانی استفاده می‌شود. در ادامه به تعریف هر کدام از این معیارهای ارزیابی می‌پردازیم.

**صحت متوسط:** نسبت تعداد نمونه‌های درست طبقه‌بندی شده به تمام نمونه‌ها که توسط رابطه‌ی ۱ محاسبه می‌شود:

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN} = \frac{N_T}{N} \quad (1)$$

که در آن

TP: تعداد نمونه‌های مثبتی که درست تشخیص داده شده‌اند.

TN: تعداد نمونه‌های منفی که درست تشخیص داده شده‌اند.

FP: تعداد نمونه‌های مثبتی که به اشتباه منفی تشخیص داده شده‌اند.

FN: تعداد نمونه‌های منفی که به اشتباه مثبت تشخیص داده شده‌اند.

دقت: عبارت است از نسبت نمونه‌های مثبت درست طبقه‌بندی شده به کل نمونه‌های مثبت موجود که در رابطه ۲ دیده

می‌شود.

$$Precision = \frac{TP}{TP + FP} \quad (۲)$$

فراخوانی: عبارت است از نسبت نمونه‌های مثبت درست طبقه‌بندی شده به کل نمونه‌هایی که مثبت تشخیص داده شده -

اند که در رابطه ۳ مشاهده می‌گردد. لازم به توضیح است که برخی از نمونه‌هایی که مثبت تشخیص داده شده‌اند اشتباه هستند و در مجموعه‌ی FN قرار می‌گیرند.

$$Recall = \frac{TP}{TP + FN} \quad (۳)$$

### نتایج تجربی

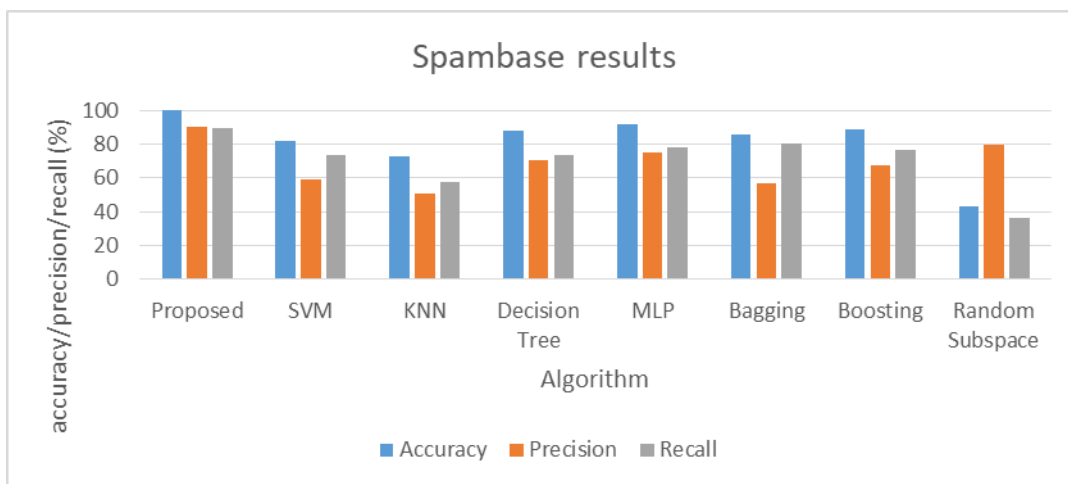
برای ارزیابی روش پیشنهادی از مجموعه داده‌ی Spambase استفاده شده است. این مجموعه داده از پایگاه UCI قابل دسترسی است. این مجموعه داده شامل ۴۶۰۱ نمونه از ایمیل‌های متنی اسپم و غیراسپم است. هر ایمیل متنی با یک بردار با ۵۷ ویژگی متناظر شده است که اغلب ویژگی‌ها فراوانی کلمات ویژه در متن ایمیل می‌باشند. مابقی ویژگی‌ها شامل برچسب نمونه (اسپم یا غیراسپم)، طول متن ایمیل، تعداد کلمات استفاده شده و تعداد کلمات استفاده شده به صورت حروف بزرگ است.

جدول ۱ و شکل ۱ نتایج به دست آمده برای روش پیشنهادی را در مقایسه با سایر الگوریتم‌های عنوان شده نشان می‌دهند. لازم به ذکر است که الگوریتم یادگیرنده‌ی روش پیشنهادی،  $k$ - نزدیک‌ترین همسایه بوده است. همانطور که دیده می‌شود، روش پیشنهادی با اختلاف چشمگیری، عملکرد بهتری نسبت سایر الگوریتم‌ها داشته است. همچنین دیده می‌شود که روش پیشنهادی دارای FP برابر با صفر است. یعنی از بین ۴۶۰۱ ایمیل متنی، هیچ ایمیل مشروعی یافت نشده که به اشتباه اسپم تشخیص داده شود. در ادامه مشاهده می‌شود که شکل ۲، مقایسه مقادیر FP را برای الگوریتم‌های مختلف نمایش می‌دهد. با توجه به شکل ۲، پس از روش پیشنهادی الگوریتم یادگیری تجمعی Random Subspace، کمترین مقدار FP را دارد که مقدار آن ۱۴۰ است. اما با دقت در جدول ۱ مشاهده می‌شود که این الگوریتم بیشترین مقدار را در FN دارد. یعنی بسیاری از ایمیل‌های اسپم را به اشتباه مشروع در نظر گرفته اند.

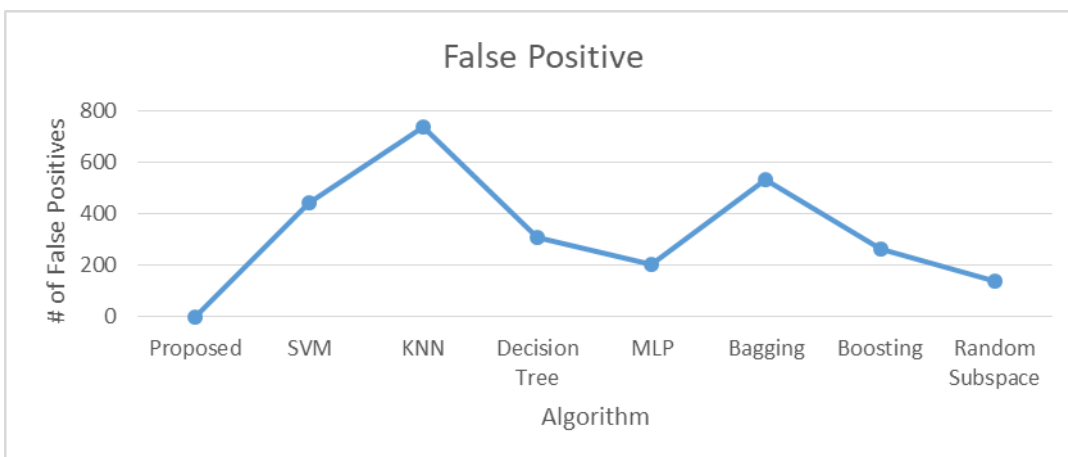
همچنین شکل ۳، زمان اجرای روش پیشنهادی را بر روی یک سیستم معمولی خانگی در مقایسه با دیگر الگوریتم‌ها نشان می‌دهد و مشاهده می‌شود زمان اجرای بسیار طولانی دارد. اما زمان اجرای آن از ماشین بردار پشتیبان کمتر بوده، در حالی که دقت آن در مقایسه با ماشین بردار پشتیبان بالاتر است.

جدول ۱- نتایج تجربی

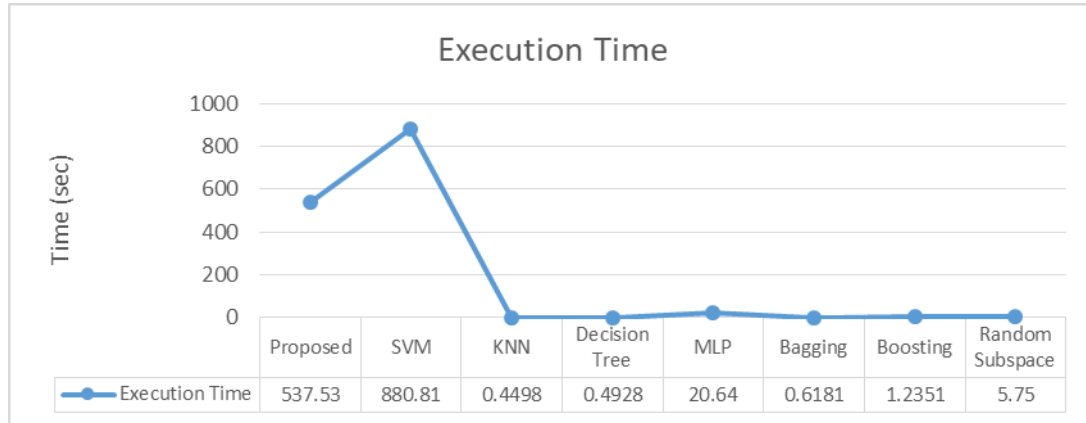
Algorithm	Accuracy (%)	Precision (%)	Recall (%)	TP	TN	FP	FN
<b>Proposed</b>	<b>99.93</b>	<b>90</b>	<b>89.71</b>	<b>1812</b>	<b>2785</b>	<b>0</b>	<b>3</b>
SVM	82.17	59.07	73.24	1371	2409	441	379
kNN	72.86	50.81	57.42	1076	2276	736	512
<b>Decision Tree</b>	87.95	70.90	73.42	1502	2544	310	244
<b>MLP</b>	91.73	74.82	78.37	1609	2611	203	177
<b>Bagging</b>	85.93	57.05	80.24	1281	2672	531	116
<b>Boosting</b>	88.82	67.15	76.48	1549	2537	263	251
<b>Random Subspace</b>	43.45	80.00	36.44	1672	326	140	2462



شکل ۱- نتایج تجربی



شکل ۲- مقایسه FP به دست آمده در الگوریتم‌های مختلف



شکل ۳- مقایسه زمان اجرای الگوریتم‌ها

همانطور که از مراحل روش پیشنهادی مشخص است، روش پیشنهادی هزینه زمانی بالایی دارد. زمان اجرای یادگیری روش پیشنهادی با افزایش ویژگی‌ها افزایش می‌یابد و در مقایسه با دیگر الگوریتم‌ها زمان اجرای بسیار بالاتری دارد.

### نتیجه‌گیری و پیشنهادات آتی

در این مقاله یک روش یادگیری جمعی جدید معرفی شد که در آن برخلاف دیگر روش‌های یادگیری جمعی، زیرمجموعه‌های آموزشی از مکان‌های مشخص انتخاب می‌شود. نتایج نشان می‌دهد که روش پیشنهادی اگرچه زمان اجرای طولانی دارد، ولی دقت آن به طور چشمگیری بالاتر از سایر روش‌ها می‌باشد. همچنین دیده می‌شود که در مجموعه داده‌ی بکار گرفته شده در این تحقیق، مقدار FP آن برابر با صفر شده است. تنها از بین ۴۶۰۱ نمونه، ۳ نمونه‌ی مشروع به اشتباه اسپم در نظر گرفته شده‌اند. در آینده روش پیشنهادی به صورت وزن‌دار نیز مورد بررسی قرار خواهد گرفت. به گونه‌ای که یادگیرنده‌هایی که به ازای هر زیرمجموعه دقت بالاتری در طبقه‌بندی داده‌های آموزش دارند، وزن بیشتری در رأی‌گیری داشته باشند. همچنین روش پیشنهادی با یادگیرنده‌های متفاوت (غیر از  $k$ - نزدیک‌ترین همسایه که در این مقاله استفاده شد) مورد آزمایش قرار خواهد گرفت.

## منابع و مراجع

- [1] Tak, Gaurav Kumar, and Shashikala Tapaswi. "Knowledge Base Compound Approach towards Spam Detection." *Recent Trends in Network Security and Applications*, 490-499, (2010).
- [2] Nelson, Marten. "Spam Control: Problems and Opportunities." *Ferris Research* (2003).
- [3] Araujo, Lourdes, and Juan Martinez-Romo. "Web spam detection: new classification features based on qualified link analysis and language models." *IEEE Transactions on Information Forensics and Security* 5.3 (2010): 581-590.
- [4] Mahajan, Renuka. "Review of data mining techniques and parameters for recommendation of effective adaptive e-learning system." *Collaborative Filtering Using Data Mining and Analysis*, 1, (2016).
- [5] Bagherzadeh-Khiabani, Farideh, et al. "A tutorial on variable selection for clinical prediction models: feature selection methods in data mining could improve the results." *Journal of clinical epidemiology* 71, 76-85, (2016).
- [6] Wang, Gang, et al. "Sentiment classification: The contribution of ensemble learning." *Decision support systems* 57 (2014): 77-93.
- [۷] بستوه، پریسا، ۱۳۹۴، ارائه روشی جهت طبقه بندی نظرات افراد با استفاده از یادگیری ترکیبی، سومین کنگره سراسری فناوری‌های نوین ایران با هدف دستیابی به توسعه پایدار، تهران، موسسه آموزش عالی مهر اروند، مرکز راه‌کارهای دستیابی به توسعه پایدار .
- [8] Zhang, Xipeng, et al. "A method of SMS spam filtering based on AdaBoost algorithm." *Intelligent Control and Automation (WCICA), 2016 12th World Congress on. IEEE*, 2016.
- [9] Fisher, R. A. "On some extensions of Bayesian inference proposed by Mr Lindley." *Journal of the Royal Statistical Society. Series B (Methodological)* (1960): 299-301.
- [10] Robinson, Gary. "A statistical approach to the spam problem." *Linux journal*2003.107, 3 (2003).
- [11] Boldi, Paolo, Massimo Santini, and Sebastiano Vigna. "PageRank as a function of the damping factor." *Proceedings of the 14th international conference on World Wide Web. ACM*, 2005 .
- [12] Gordillo, José, and Eduardo Conde. "An HMM for detecting spam mail." *Expert systems with applications* 33.3 (2007): 667-682 .
- [13] Spracklin, L. M., and Lawrence V. Saxton. "Filtering spam using kolmogorov complexity estimates." *Advanced Information Networking and Applications Workshops, 2007, AINAW'07. 21st International Conference on. Vol. 1. IEEE*, 2007 .
- [14] Korelov, S. V., A. K. Kryukov, and L. U. Rotkov. "Text Messages' Digital Analysis on Spam Identification." *Russian, Proceedings of Scientific Conference on Radiophysics, Nizhni Novgorod State University, Nizhny Novgorod Oblast*. 2006
- [15] Hsiao, Wen-Feng, and Te-Min Chang. "An incremental cluster-based approach to spam filtering." *Expert Systems with Applications* 34.3, 1599-1608, (2008).
- [16] Lee, Sang Min, et al. "Spam detection using feature selection and parameters optimization." *Complex, Intelligent and Software Intensive Systems (CISIS), 2010 International Conference on. IEEE*, 2010.
- [17] Saeedian, Mehrnoush Famil, and Hamid Beigy. "Spam detection using dynamic weighted voting based on clustering." *Intelligent Information Technology Application, 2008. IITA'08. Second International Symposium on. Vol. 2. IEEE*, 2008 .
- [18] Sasaki, Minoru, and Hiroyuki Shinnou. "Spam detection using text clustering." *Cyberworlds, 2005. International Conference on. IEEE*, 2005

- [19] Cortez, Paulo, et al. "Symbiotic data mining for personalized spam filtering." Proceedings of the 2009 IEEE/WIC/ACM International Joint Conference on Web Intelligence and Intelligent Agent Technology-Volume 01. IEEE Computer Society, 2009 .
- [20] Weinstein, Lauren. "Spam wars." Communications of the ACM 46.8, 136 (2003).
- [21] Gburzynski, Pawel, and Jacek Maitan. "Fighting the spam wars: A remailer approach with restrictive aliasing." ACM Transactions on Internet Technology (TOIT) 4.1, 1-30 (2004) .
- [22] Li, Fulu, H. Mo-Han, and G. Pawel. "The community behavior of spammers." (2011).
- [23] Xu, K., et al. "Revealing social networks of spammers through spectral clustering." Communications, 2009. ICC'09. IEEE International Conference on. IEEE, 2009 .
- [24] Xu, Kevin S., Mark Kliger, and Y. Chen. "Tracking communities of spammers by evolutionary clustering." International Conference on Machine Learning Workshop on Social Analytics: Learning from Human Interactions. 2010.