

## بررسی روش‌های کاهش ابعاد در سیستم‌های تشخیص نفوذ

رویا گلستانی شیشوان<sup>۱\*</sup>، امیرفرید امینیان مدرس<sup>۲</sup>.

۱- دانشجوی کارشناسی ارشد مهندسی کامپیوتر نرم‌افزار، دانشگاه صنعتی سجاد، مشهد

۲- عضو هیئت علمی، دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی سجاد، مشهد

### چکیده

برای ایجاد امنیت داده‌ها در یک سیستم کامپیوتری، علاوه بر دیوارهای آتش و دیگر تجهیزات جلوگیری از نفوذ، سیستم‌های دیگری به نام تشخیص نفوذ (IDS) مورد نیاز می‌باشند تا بتوانند در صورتی که نفوذگر از دیوار آتش و دیگر تجهیزات امنیتی عبور کرد، آن را تشخیص دهد. تشخیص نفوذ، با مشکلاتی از قبیل حجم عظیم ترافیک شبکه، توزیع داده‌های بسیار نامتعادل، سختی شناسایی مرزهای تصمیم‌گیری بین رفتار طبیعی و غیرطبیعی، مواجه است. در این مقاله، هفت روش تبدیلی کاهش ابعاد (شامل دو نوع خطی و غیرخطی) برای سیستم تشخیص نفوذ مبتنی بر ماشین بردار پشتیبان، مقایسه و ارزیابی می‌شوند. این روش‌ها شامل تحلیل مولفه‌های اصلی، تحلیل متمایز خطی (LDA)، تحلیل عاملی (FA)، تحلیل مولفه اساسی مبتنی بر هسته (Kernel PCA)، کدکننده‌های خودکار (Autoencoder)، تحلیل مولفه‌های اصلی احتمالاتی (ProbPca) و مقیاس‌گذاری چندبعدی (MDS) می‌باشند. در این مقاله از مجموعه داده‌ی NSL-KDD جهت ارزیابی سیستم تشخیص نفوذ و روش‌های کاهش ابعاد استفاده می‌شود. همچنین مدت زمان آموزش، آزمایش، میزان صحت، دقت و درصد خطا توسط سیستم تشخیص نفوذ باهم مقایسه می‌شود. نتایج مقایسه‌ها نشان می‌دهد که روش ProbPca زمان آموزش و آزمایش مناسب‌تری از روش‌های تبدیلی دیگر دارد.

کلمات کلیدی: تشخیص نفوذ، کاهش ابعاد، امنیت داده‌ها.

### ۱. مقدمه

رشد تکنولوژی اینترنت، ارتباطات داده‌ها را گسترش داده است. از طرفی تهدیدات شبکه و مسائل امنیتی، مشکل جامعیت داده‌ها و از دست رفتن اطلاعات را افزایش می‌دهد. با توجه به این موضوع بحث امنیت اهمیت بیشتری پیدا کرده است. تعریف امنیت داده‌ها شامل تعریف سه بخش: یکپارچگی، محرمانگی و دسترس‌پذیری داده‌ها است. یکپارچگی داده‌ها به این معناست که تنها افراد و سیستم‌های مجاز می‌توانند در داده تغییر ایجاد کنند. محرمانگی داده به این معناست که تنها افراد و سیستم‌های مجاز می‌توانند اطلاعات حساس یا طبقه‌بندی شده را ببینند و سایر موارد امکان هیچگونه دسترسی به

\* Email: [golestani.roya@gmail.com](mailto:golestani.roya@gmail.com)



این داده‌ها را نداشته باشند. دسترس‌پذیری شامل داده‌ها و سیستم‌ها می‌شود. اگر شبکه یا داده آن، برای کاربران مجاز در دسترس نباشد (به دلایلی مانند حمله 'DoS' یا خرابی‌های معمول شبکه) می‌تواند مسائل جدی برای سازمان و کاربران که به شبکه به عنوان یک ابزار متکی هستند به وجود آورد [۱].

نمودها به دو دسته داخلی و خارجی تقسیم می‌شوند. نمودهای خارجی به نمودهایی گفته می‌شوند که توسط افراد مجاز یا غیرمجاز، از خارج شبکه به درون شبکه صورت می‌گیرد. نمودهای داخلی توسط افراد مجاز در سیستم و شبکه داخلی، از درون خود شبکه انجام می‌پذیرد. یک سیستم تشخیص نمود تمام فعالیت‌های شبکه را برای شناسایی حملات شناخته شده و یا ناشناخته نظارت می‌کند و برای تشخیص اینکه این وقایع نشانه‌ای از یک حمله یا استفاده مشروع از سیستم هستند، تصمیم می‌گیرد [۲].

مسئله تشخیص نمود، فرایند نظارت بر وقایع رخ داده در یک شبکه و یا یک سیستم کامپیوتری در جهت کشف موارد انحراف از سیاست‌های امنیتی می‌باشد. سیستم تشخیص نمود (IDS) نرم‌افزاری است با قابلیت تشخیص، آشکارسازی و پاسخ (واکنش) به فعالیت‌های غیرمجاز یا ناهنجار که در رابطه با سیستم می‌باشد [۲].

با جمع‌آوری داده‌های ترافیک شبکه و با کمک روش‌های انتخاب ویژگی، تنها بردارهای ویژگی که حاوی اطلاعات ضروری هستند در نظر گرفته می‌شوند. داده‌های جمع‌آوری شده در این مرحله برای تعیین اینکه عادی هستند یا غیرعادی، تجزیه و تحلیل می‌شوند. در این گام از روش‌های مختلف برای تشخیص نمود استفاده می‌شود. در انتها واکنش سیستم تشخیص نمود است که هشدار لازم را به مدیر سیستم مبنی بر اتفاق افتادن یک حمله و ماهیت آن حمله می‌دهد. همچنین سیستم تشخیص نمود نیز در کنترل حملات با بستن پورت شبکه یا از بین بردن فرایندها، کمک می‌کند.

همانطور که پیش‌تر ذکر شد با افزایش مشکلات امنیتی و حملات شبکه در سال‌های اخیر، سیستم‌های تشخیص نمود، یک جزء اصلی برای امنیت شبکه به شمار می‌آیند. این سیستم‌ها با داده‌های حجیم برای تحلیل مواجه هستند. بررسی بسیاری از این داده‌ها نشان می‌دهد که بسیاری از ویژگی‌های آنها غیرمفید و بی‌تاثیر هستند. بنابراین حذف ویژگی‌های نامناسب از این مجموعه داده، یک راهکار مناسب برای کاهش مجموعه داده در سیستم‌های تشخیص نمود است [۳].

امروزه بیشتر رویکردها در تشخیص نمود مربوط به مسئله استخراج ویژگی‌های مهم متمرکز شده است. اما استخراج ویژگی‌ها باعث از دست دادن قسمتی از داده‌ها خواهد شد. بیشتر سیستم‌های تشخیص نمود کنونی از تمامی پارامترهای موجود در بسته‌های شبکه برای ارزیابی و کشف الگوهای حملات استفاده می‌نمایند، در صورتی که برخی از این پارامترها غیرمرتبط و زائد می‌باشند.

با توجه به اینکه در مجموعه داده‌های سیستم‌های تشخیص نمود با تعداد بالایی از ویژگی‌ها رو به رو هستیم، کاهش ابعاد می‌تواند به راحت‌تر شدن تحلیل‌ها، افزایش عملکرد جداکننده، حذف اطلاعات تکراری و غیرمرتبط کمک کند.

کاهش حجم داده‌ها در دو شکل امکان‌پذیر است. اول، کاهش تعداد ویژگی‌ها، که منجر به حذف ویژگی‌هایی می‌شود که تاکید کمتری بر سیستم تشخیص نمود دارند. دوم، کاهش تعداد نمونه‌ها که منجر به حذف رکوردها و نمونه‌هایی می‌شود که با حذف آنها، سیستم تشخیص نمود با دقت بهتری کار می‌کند. در این مقاله روش‌های کاهش ابعاد و تاثیر آنها در تشخیص نمود مورد بررسی و تحقیق قرار گرفته است. هدف این بررسی شناسایی مناسب‌ترین روش کاهش ابعاد در سیستم تشخیص نمود براساس ماشین بردار پشتیبان (SVM) است.

ساختار ادامه مقاله به این شرح است: بخش دوم مروری بر روش‌های شناخته شده کاهش ابعاد است. در بخش سوم بررسی تاثیر روش‌های خطی و غیرخطی کاهش ابعاد در سیستم‌های تشخیص نمود ارائه خواهد شد. بخش چهارم ارزیابی روش‌های مورد بررسی به همراه جدول نتایج و مقایسه‌ها و تحلیل آنها مطرح می‌شود و نتیجه‌گیری کلی در بخش پنجم بیان شده است.

## ۲. مروری بر ادبیات پیشین

در این بخش ابتدا بررسی ادبیات موضوع مرتبط با کاهش ابعاد انجام می‌گیرد سپس روش‌های کاهش ابعاد در سیستم‌های تشخیص نفوذ مطرح می‌شود.

کاهش ابعاد یک راه موثر برای بهبود عملکرد الگوریتم‌های یادگیری ماشین است. کاهش ویژگی‌های بی‌ربط، باعث کاهش زمان آموزش و افزایش دقت طبقه‌بندی می‌شود. انتخاب ویژگی بعنوان یک روش کاهش ابعاد، یک زیرمجموعه ویژگی کاهش یافته را انتخاب می‌کند و همزمان می‌تواند اندازه زیرمجموعه را نیز کاهش دهد و قدرت افتراق را نیز بهبود دهد.

یکی دیگر از جنبه‌های مهم در زمینه‌ی کاهش ابعاد، ارزیابی میزان تاثیرگذاری ویژگی‌های انتخاب شده در طبقه‌بندی می‌باشد. برای ارزیابی میزان تاثیرگذاری ویژگی‌ها، ابتدا ویژگی‌ها را به چندین زیرمجموعه تقسیم کرده، سپس طبقه‌بند را با استفاده از هر یک از این زیرمجموعه ویژگی‌ها، آموزش داده و در نهایت با استفاده از داده‌های آزمون، میزان تاثیرگذاری هر یک از این زیرمجموعه‌ها بدست می‌آید.

الگوریتم‌های مختلفی برای طراحی سیستم‌های تشخیص نفوذ پیشنهاد شده‌اند که دارای قابلیت متعددی هستند. جهت بالا بردن قابلیت یادگیری و کاهش میزان محاسبات در الگوریتم‌ها از روش‌های کاهش ابعاد استفاده شده است.

یکی از تحقیقات انجام شده، استفاده از تجزیه و تحلیل مولفه اصلی (PCA<sup>۴</sup>) در سیستم تشخیص نفوذ در شبکه است. در [۴]، تاثیر PCA در سیستم تشخیص نفوذ بررسی شده است. همچنین تعداد ایده‌آل مولفه‌های اصلی مورد نیاز برای تشخیص نفوذ و تاثیر داده‌های آلوده به نویز روی PCA نیز مورد توجه بوده است. مجموعه داده‌های با اندازه‌ی اصلی  $n \times d$  به ساختاری با  $k$  مولفه‌ی اصلی معین، نگاشت شده‌اند و به مجموعه داده‌ای با اندازه‌ی  $n \times k$  تبدیل یافته‌اند، که  $n$  تعداد نمونه‌ها و  $d$  تعداد ابعاد اصلی است.  $K$ ، تعداد مولفه‌های اصلی با محدوده‌ی تغییرات از ۲ تا ۲۰ است.

نتایج انجام آزمایشات با روش PCA و با استفاده از الگوریتم‌های مختلف طبقه‌بندی مثل C4.5، Random Forest، روی دو مجموعه داده به نام‌های KDD CUP و UNB ISCX نشان می‌دهد که تعداد ۱۰ مولفه، برای طبقه‌بندی ایده‌آل است. همچنین دقت طبقه‌بندی برای ۱۰ مولفه اصلی به ترتیب در حدود ۹۹٫۷٪ و ۹۸٫۸٪ برای مجموعه داده‌های ذکر شده می‌باشد که تقریباً همان دقت به دست آمده با استفاده از ۴۱ ویژگی از ۳۱۲۷۹ نمونه برای مجموعه داده‌ی KDD و ۲۸ ویژگی از ۳۳۷۴۶ نمونه برای مجموعه داده ISCX می‌باشد.

نتایج این تحقیق نشان داد که نسبت کاهش با PCA برای مجموعه داده KDD Cup و UNB ISCX، به ترتیب ۰٫۲۴ و ۰٫۳۶ است. همچنین PCA با حضور نویز در داده‌ها، دقت طبقه‌بندی را کاهش داده است و در مقابل هنگامی که داده بدون نویز است افزایش دقت را خواهد داشت. استفاده از PCA برای طراحی یک سیستم تشخیص نفوذ ضمن اینکه پیچیدگی سیستم را کاهش خواهد داد باعث دستیابی به دقت طبقه‌بندی بالاتر نیز می‌شود.

در [۵] یک مدل تشخیص نفوذ با استفاده از ترکیب انتخاب ویژگی Chi-Square و SVM چندکلاسه ارائه شده است. بسیاری از سیستم‌های تشخیص نفوذ، تنها از یک الگوریتم طبقه‌بندی جهت دسته‌بندی ترافیک شبکه بعنوان نرمال و غیرنرمال استفاده می‌کنند. با توجه به میزان زیاد داده‌ها، این مدل از طبقه‌بندها موفق به دستیابی با نرخ تشخیص حمله بالا و کاهش نرخ هشدار غلط نمی‌شوند. با این حال، با استفاده از کاهش ابعاد داده‌ها می‌توانند به یک مجموعه‌ی بهینه از ویژگی‌ها بدون از دست دادن اطلاعات دست یابند و سپس با استفاده از روش مدل‌سازی چندکلاسه، شناسایی حملات شبکه‌ای متفاوت را طبقه‌بندی کنند.



از دیگر تحقیقات انجام شده در [۶] تجزیه و تحلیل ویژگی‌ها، ارزیابی و مقایسه الگوریتم‌های طبقه‌بندی مبتنی بر مجموعه‌ی داده نفوذ دارای نویز می‌باشد. ترافیک داده‌های شبکه در دنیای واقعی با مقدار زیادی از اطلاعات آلوده به نویز همراه است و IDS اغلب در چنین محیطی کار می‌کند. یکی از مسائل چالش برانگیز در IDS برخورد با محیط داده‌ی پر نویز جهت تشخیص حملات از فعالیت‌های شبکه می‌باشد.

در [۶]، الگوریتم‌های داده‌کاوی و یادگیری ماشین مختلف با مجموعه داده‌های NSL-KDD و KDD'99 با نویز ۱۰ درصد و ۲۰ درصد مورد ارزیابی و مقایسه قرار گرفته است. نتایج تجربی نشان می‌دهد که الگوریتم NN(SOM) در مقایسه با دیگر الگوریتم‌های مورد مطالعه از لحاظ مقاوم بودن به محیط نویزی به مراتب بهتر است. با این حال، JRip و J48 از خانواده الگوریتم‌های درخت، عملکرد بهتری نسبت به دیگر الگوریتم‌ها داشتند.

میزان وابستگی ویژگی‌ها در مجموعه داده‌ها برای یک طبقه‌بند خاص توسط روش طبقه‌بندی مبتنی بر عملکرد (PMR) تجزیه و تحلیل شده است [۶]. ارزیابی نتایج آماری ثابت کرده که هر طبقه‌بند دارای یک ترکیب منحصر به فرد از یک زیرمجموعه ویژگی با نتایج عملکرد مطلوب است و زیرمجموعه ویژگی انتخاب شده توسط هر الگوریتم طبقه‌بندی، متفاوت از دیگری است. نتایج تجربی نشان داده است که ارزیابی سیستم تشخیص نفوذ براساس مجموعه داده‌ی NSL-KDD دارای نتایج واقعی تری در مقایسه با مجموعه داده‌ی اصلی KDD'99 است. نتایج تجربی نشان داده است الگوریتمی که به خوبی روی مجموعه‌ی داده KDD'99 (۱۰٪ داده نویزی) اجرا می‌شود روی مجموعه داده‌ی مشابه NSI-KDD (۲۰٪ داده نویزی) نتیجه نمی‌دهد، لذا این امر ثابت می‌کند مجموعه داده NSL-KDD نشان‌دهنده محیطی واقعی‌تر برای ارزیابی الگوریتم‌های طبقه‌بندی نسبت به مجموعه داده‌ی KDD'99 است [۶].

در [۷]، یک سیستم تشخیص نفوذ برای تشخیص حملات به طور موثر مطرح شده است. برای این منظور، یک الگوریتم انتخاب ویژگی جدید به نام الگوریتم انتخاب ویژگی بهینه مبتنی بر نسبت اطلاعات ارائه شده است. این الگوریتم انتخاب ویژگی، فقط تعدادی از ویژگی‌های مطلوب و مهم را از مجموعه داده‌ی KDD Cup انتخاب می‌کند. علاوه بر این، دو روش طبقه‌بندی ماشین بردار پشتیبان و طبقه‌بند مبتنی بر قانون، جهت تاثیر بر طبقه‌بندی داده‌ها و رسیدن به دقت بیشتر استفاده شده است.

در [۸] یکی از چالش‌های مهم در تحقیقات تشخیص نفوذ، طراحی یک سیستم تشخیص نفوذ دقیق از نظر میزان تشخیص بالا، دقت بالا و نرخ هشدار غلط کم است. در این مقاله، یک ساختار کلی از یک رویکرد یادگیری ترکیبی ارائه شده است. سپس روش پیشنهادی با استفاده از خوشه‌بندی K-means و طبقه‌بندهای چندتایی پیاده‌سازی شده است. داده‌ها با روش مبتنی بر الگوریتم خوشه‌بندی K-means تقسیم‌بندی شده است. سپس، هر قسمت با استفاده از یک طبقه‌بند متمایز تقسیم‌بندی شده است. شبکه بیزین، ماشین بردار پشتیبان و الگوریتم‌های طبقه‌بندی OneR بعنوان طبقه‌بند استفاده شده است. روش ترکیبی ارائه شده نتایج بهتری نسبت به تک طبقه‌بند از نظر سرعت کشف و شناسایی، دقت و نرخ هشدار غلط دارا بوده است. نرخ تشخیص روش ترکیبی ارائه شده ۹۹٫۵۰٪ است.

در تحقیقات دیگر، یکی از چالش‌های مهم امنیت شبکه، تشخیص ناهنجاری می‌باشد. در [۹] روش جدیدی به نام G-LDA ارائه شد که ترکیب یکپارچه‌سازی تخصیص پنهان دیریکله و الگوریتم ژنتیک با هدف شناسایی ناهنجاری‌ها در ترافیک شبکه بود. علاوه بر این، انتخاب ویژگی نقش مهمی را در شناسایی زیرمجموعه‌ی بهینه از ویژگی‌ها برای تعیین بسته‌های ناهنجار ایفا می‌کند. تخصیص پنهان دیریکله، شناسایی مجموعه‌ی بهینه از ویژگی‌ها را برای طبقه‌بندی انجام می‌دهد و از الگوریتم ژنتیک برای محاسبه نمره اقلام داده‌ای و تولید یک جمعیت از زیرمجموعه‌های کاندید استفاده شده است. با استفاده از تابع ارزیابی، میزان شایستگی عناصر جمعیت فعلی مشخص شده و در نهایت بعد از فیلتر شدن عناصر بهتر برای جمعیت نسل بعد انتخاب می‌شوند. این روش بر روی مجموعه داده‌ی KDD-Cup'99 انجام شده و نتایج



تجربی نشان می‌دهد که روش ترکیبی دقت بهتری را برای تشخیص حملات شناخته شده و ناشناخته بدست آورده است. همچنین نرخ مثبت کاذب کمتری نیز گزارش شده است.

از دیگر تحقیقات انجام شده، روش مرکز خوشه و نزدیکترین همسایه است [۱۰]. هدف در سیستم‌های تشخیص نفوذ (IDS) تشخیص انواع مختلفی از ترافیک‌های مخرب شبکه است که توسط یک فایروال عادی و معمولی تشخیص داده نمی‌شود. بسیاری از سیستم‌های تشخیص نفوذ براساس تکنیک‌های یادگیری ماشین توسعه داده شده‌اند. در [۱۰] یک روش نمایش ویژگی جدید به نام روش مرکز خوشه و نزدیکترین همسایه ارائه شده است. روش نمایش ویژگی یک الگوی طبقه‌بند مهم است که کلاس‌بندی صحیح را آسان می‌کند، با این حال، مطالعات کمی روی چگونگی استخراج بیشتر ویژگی‌های مهم برای اتصالات نرمال و تشخیص موثر حملات تمرکز کرده‌اند. در روش مرکز خوشه و نزدیکترین همسایه (CANN) دو فاصله، اندازه‌گیری می‌شود: فاصله اول مبتنی بر فاصله بین هر نمونه داده و مرکز خوشه‌اش و فاصله دوم بین داده‌ها و نزدیکترین همسایه‌اش در همان خوشه است. نتایج نشان می‌دهد که دقت طبقه‌بندی CANN نسبت به KNN و SVM در مجموعه داده‌ی KDD-Cup'99 بیشتر است. همچنین کارایی محاسباتی بالایی را در زمان آموزش و آزمون طبقه‌بند فراهم می‌آورد.

در تحقیقات دیگر، برای به حداکثر رساندن اثربخشی هر یک از الگوریتم‌های استخراج ویژگی و ایجاد یک سیستم تشخیص نفوذ کارآمد، یک مجموعه‌ای از الگوریتم‌های استخراج ویژگی تحلیل الگوی متمایز خطی (LDA) و PCA پیاده‌سازی شده است [۱۱]. این روش منجر به نتایج خوبی شده است و دقت بیشتری را در مقایسه با یک روش استخراج ویژگی خاص نشان داده است. هدف، بررسی امکان استفاده از مجموعه ویژگی‌های تولید شده توسط دو روش انعطاف‌پذیر استخراج ویژگی برای تشخیص نفوذ است. بنابراین، حضور گروهی از ویژگی‌های مشابه در هر دو مجموعه ویژگی‌های تولید شده توسط دو روش استخراج ویژگی، اهمیت ویژگی‌های مشترک در مجموعه داده‌ی شبکه را تایید می‌کند. با حداقل نگه داشتن تعداد ویژگی‌ها و بدست آمدن یک مجموعه ایده‌آل از ویژگی‌ها، با پیروی از روش‌های ارائه شده در [۱۱]، دستیابی به دقت بالا در سیستم تشخیص نفوذ محقق شده است. تعداد کم ویژگی‌ها به این معنی است که طبقه‌بند برای آموزش نیاز به اطلاعات کمتری دارد. نتایج تجربی نشان می‌دهد که کاهش ابعاد می‌تواند میزان تشخیص را افزایش دهد، همچنین مجموع روش‌های استخراج ویژگی، به طور مستقیم عملکرد بهتری را در تشخیص میزان نفوذ نشان داده است.

در پژوهش دیگر [۱۲]، به صورت مصنوعی از روش تحلیل مولفه‌ی اصلی مبتنی بر هسته (KPCA) و ماشین احتمال Minimax برای شناسایی نفوذ استفاده شده است. ویژگی‌های نفوذ توسط KPCA استخراج شده و با ماشین احتمال Minimax طبقه‌بندی شده است. داده‌های آزمایشی از مجموعه داده‌ی KDD'99 می‌باشند. در این مقاله، جهت استخراج موثر ویژگی‌ها، یک روش جدید برای تشخیص نفوذ مبتنی بر طبقه‌بند ماشین احتمال Minimax با قابلیت استخراج ویژگی توسط KPCA پیشنهاد شده است. شبیه‌سازی نشان می‌دهد که ماشین احتمال Minimax با استخراج ویژگی KPCA می‌تواند عملکرد تعمیم بهتری را از ماشین بردار پشتیبان و ماشین احتمال Minimax بدون استخراج ویژگی بدست آورد. این آزمایش همچنین نشان داد که روش ارائه شده نیاز به زمان آموزش کمتری دارد. بنابراین، الگوریتم طبقه‌بند ماشین احتمال Minimax<sup>۱</sup> با استخراج ویژگی‌ها توسط KPCA برای تشخیص نفوذ بسیار کارآمدتر می‌باشد.

### ۳. بررسی تاثیر روش‌های کاهش ابعاد

این بخش نتایج پیاده‌سازی و آزمایشات روش‌های کاهش ابعاد می‌باشد. در بخش اول مجموعه داده به کار برده شده به همراه اطلاعات مربوط به آن آورده شده است. بخش دوم معیارهای ارزیابی و در بخش سوم انتخاب پارامترهای طبقه‌بندی ارائه شده است. در نهایت نتایج آزمایشات و مقایسه‌ها در بخش سوم صورت گرفته است.

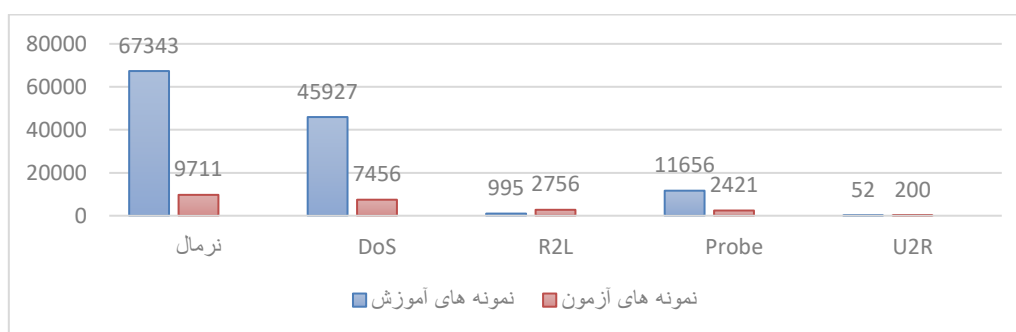
#### ۳-۱- مجموعه داده

جهت بررسی ارزیابی روش پیشنهادی از مجموعه داده ترافیک شبکه استفاده شده است. این مجموعه داده توسط واحد فناوری و سیستم‌های سایبری آزمایشگاه لینکلن MIT جمع‌آوری شده است. مجموعه داده NSL-KDD (مجموعه داده بهبود یافته KDDcup, 1999) با مشخصات مندرج در جدول (۱) در این مقاله استفاده شده است. برای هر اتصال در این مجموعه داده‌ها، ۴۱ ویژگی گسسته و پیوسته تعریف گردیده است و هر اتصال دارای برجستگی است که مشخص می‌کند اتصال مذکور نرمال است یا یکی از انواع حملات. انتخاب مجموعه داده KDD-NSL نسبت به مجموعه داده اصلی KDD این مزیت را دارد که شامل مجموعه رکوردهای تکراری در مجموعه داده یادگیری نمی‌باشد به طوری که طبقه‌بند نسبت به رکوردهای مکرر گرایش پیدا نمی‌کند. این مجموعه داده شامل حملات نوع DoS, R2L, U2R, و Probing است.

جدول ۱- اطلاعات مجموعه داده NSL-KDD [۱۳]

تعداد داده‌ها	مجموعه داده‌های آموزشی (KDDtrain+)	مجموعه داده‌های آزمایشی (KDDTest+)
نرمال	67343	9711
حمله	58630	12833
مجموع	125973	22544

به‌طور کلی، حمله‌ها به چهار دسته اصلی تقسیم می‌شوند که عبارتند از: U2R, R2L, DoS, و Probe. در شکل ۱ تعداد نمونه‌های آموزش و آزمون در گروه‌های داده‌های نرمال و حمله نشان داده شده است. در این شکل محور افقی نوع داده‌ها و محور عمودی تعداد نمونه‌ها را نشان می‌دهد که در آن داده‌های آزمایش با رنگ آبی و داده‌های آزمون با رنگ قرمز نمایش داده شده‌اند.



شکل ۱- تعداد نمونه‌های آموزش و آزمون در گروه‌های مختلف

#### ۳-۲- معیارهای ارزیابی



ارزیابی داده‌های کاهش یافته حاصل از روش‌های مطرح شده در بخش قبلی توسط طبقه‌بند ماشین بردار پشتیبان صورت می‌گیرد. در این بخش داده‌های کاهش یافته به‌عنوان داده‌های آموزش و آزمون جهت طبقه‌بندی در نظر گرفته می‌شوند که برای این تقسیم‌بندی داده‌ها، از اعتبارسنجی ضربدری ۱۰ تایی استفاده شده است. معیارهای ارزیابی در نظر گرفته شده برای بررسی داده‌های تست شامل دقت، صحت و  $G\text{-mean}$  می‌باشد. معیار  $G\text{-mean}$  از جمله مهمترین معیارهای ارزیابی داده‌های نامتوازن می‌باشد. با توجه به تفاوت زیاد تعداد نمونه‌های کلاس نرمال با تعداد نمونه‌های هر یک از کلاس‌های حمله، می‌توان این مجموعه داده را به‌صورت نامتوازن در نظر گرفت، لذا معیار  $G\text{-mean}$  به‌عنوان یک معیار ارزیابی در مساله مورد بررسی، می‌تواند معیاری کاربردی باشد. ابتدا معیارهای ارزیابی را برای کل مجموعه داده با ۴۱ ویژگی، بدون کاهش ابعاد ارزیابی می‌کنیم. نتایج این بررسی در جدول (۲) آورده شده است. صحت تشخیص کلاس عبارتست از تعداد حملات مربوط به کلاس که درست تشخیص داده شده تقسیم بر کل حملات مربوط به آن کلاس، خطای مثبت عبارتست از رفتارهای نرمال که حمله تشخیص داده شده‌اند تقسیم بر کل رکوردهای نرمال.

### ۳-۳- انتخاب پارامترهای دسته‌بندی

آزمایشات صورت گرفته بر روی روش‌ها توسط نرم‌افزار متلب ۲۰۱۶ بر روی پردازنده ۴،۲ GHZ و ۳۲ گیگا بایت حافظه اصلی انجام شده است. در تمام آزمایشات، از روش اعتبارسنجی ضربدری ۱۰ تایی به‌منظور ارزیابی استفاده شده است. در بررسی های انجام شده از هسته تابع پایه‌ای شعاعی استفاده شده است. از طرف دیگر دقت طبقه‌بندی داده‌های کاهش داده شده توسط طبقه‌بند ماشین بردار پشتیبان ارزیابی شده است، به این صورت که این طبقه‌بندها با داده‌های کاهش داده شده آموزش می‌بینند سپس با داده‌های آزمون این طبقه‌بندها ارزیابی می‌شوند.

### ۴. ارزیابی و مقایسه

در این بخش نتایج آزمایشات گزارش شده است. ابتدا، نتایج آزمایش‌ها بدون کاهش ابعاد در مجموعه داده مورد نظر یعنی با ۴۱ ویژگی ارائه شده است، سپس نتایج حاصل از طبقه‌بندی SVM با ورودی‌های تعداد ویژگی‌های متفاوت با ارائه جدول و رسم نمودار تحلیل و بررسی شده است. همچنین نتایج به‌دست آمده از بهترین روش با روش ارائه شده در [۱۴] مقایسه شده است. در جدول (۲) نتایج مربوط به طبقه‌بندی کننده‌های SVM را با ۴۱ ویژگی نشان می‌دهد.

جدول ۲- نتایج حاصل از طبقه بندی کننده های SVM با ورودی‌های شامل ۴۱ ویژگی (بدون کاهش ابعاد)

مدت زمان (آزمایش ثانیه)	مدت زمان آموزش (ثانیه)	G- mean	میانگین مربعات خطا	خطا	دقت	صحت	کلاس
۳۲،۵۳	۱۰۵۲	۹۹،۵۰	۰،۵۰	۰،۵۱	۹۸،۹۶	۹۹،۴۸	Normal
۳۱،۶۳	۹۲۰،۳۹	۹۸،۲۲	۰،۴۱	۰،۴۲	۹۹،۳۹	۹۹،۵۷	Probe
۳۲،۶۵	۱۳۵۵	۹۹،۷۱	۰،۲۵	۰،۲۵	۹۹،۷۹	۹۹،۷۴	DOS

با توجه به جدول (۲)، SVM با ۴۱ ویژگی بالاترین دقت طبقه‌بندی را برای کلاس DOS دارد و با توجه به مقادیر دقت و  $G\text{-mean}$  به‌دست آمده برای این کلاس می‌توان نتیجه گرفت که داده‌های این کلاس با سایر کلاس‌ها تفکیک‌پذیری بیشتری دارد، لذا پیش‌بینی می‌شود که مقادیر دقت و  $G\text{-mean}$  این کلاس بعد از کاهش نیز زیاد باشد. در جدول (۲)



میانگین زمان آموزش برابر است با ۱۱۰۹,۱۳ ثانیه و میانگین زمان آزمایش برابر با ۳۲,۲۷ ثانیه است. از آنجایی که یکی از اهداف کاهش، بهبود سرعت پردازش داده‌هاست، در جدول (۳) میانگین زمان آموزش و آزمون برای بیشترین (۹ بعد) و کمترین (۲۳ بعد) میزان کاهش آورده شده است تا با زمان آموزش و آزمون داده‌های اصلی مقایسه شود. با توجه به جدول (۳)، کاملاً بدیهی است که زمان پردازش داده‌های کاهش یافته کمتر از داده‌های اصلی باشد. در برخی از روش‌ها مثل kernelPCA, LDA, ProbPCA و Autoencoder کاهش زمان پردازش چشمگیری مشاهده می‌شود. لذا با اثبات اینکه کاهش داده‌ها ارتباط مستقیمی با زمان پردازش آنها دارد، در ادامه‌ی آزمایش‌های بعدی به بررسی دقت طبقه‌بندی داده‌های کاهش یافته می‌پردازیم. همانطور که در جدول (۲) ملاحظه شد، SVM با ۴۱ ویژگی، معیار صحت برابر با روش PCA را برای کلاس DOS دارد و بالاترین معیار صحت برای کلاس نرمال مربوط به روش‌های PCA, ProbPCA و MDS<sup>۱۱</sup> می‌باشد. در شناسایی DOS بهترین نتیجه با ویژگی‌های حاصل از FA<sup>۱۲</sup> بدست آمده است.

جدول ۳- مقایسه زمان آموزش و آزمایش SVM روش‌های کاهش ابعاد با بیشترین و کمترین میزان کاهش

نرخ شناسایی FA با	نرخ شناسایی با Auto Encoder	نرخ شناسایی با MDS	نرخ شناسایی با ProbPCA	نرخ شناسایی با KPCA	نرخ شناسایی با LDA	نرخ شناسایی با PCA	
۱۹۶۴	۳۱۵,۶۸	۱۳۰۰	۳۵۷,۵۳	۶۷۱,۰۷	۳۲۱,۷۱	۱۲۶۴	میانگین زمان آموزش (ثانیه)
۳۱,۶۳	۱۰,۷۷	۲۲,۲۳	۶,۴۵	۲۴,۵۰	۳۵,۲۱	۲۲,۲۹	میانگین زمان آزمایش (ثانیه)
۷۱۸,۸۸۸	۷۳	۶۱۵,۲۵	۹۲,۹۱۲	۹۶,۶۹	۴۳,۳۲	۶۱۶,۸۲۲	میانگین زمان آموزش (ثانیه)
۱۳,۸۶۴	۶,۰۴	۱۳,۲۲۲	۱,۴۷	۶,۱۹	۴,۳۰	۱۳,۱۴	میانگین زمان آزمایش (ثانیه)

همانطور که گفته شد آزمایش‌ها بر روی داده‌های کاهش یافته با ابعاد ۹، ۱۴، ۱۷، ۲۰ و ۲۳ بعد، برای سه کلاس شامل یک کلاس نرمال و دو کلاس حمله، برای معیارهای ارزیابی شامل دقت، صحت، میانگین مربعات خطا، زمان اجرا و G-mean صورت گرفته است. گزارش تمامی مقادیر به صورت جدول به دلیل حجم بالای نتایج امکان‌پذیر نیست، لذا نتایج برای ۲۳ بعد به صورت جدول، کامل ارائه می‌گردد و برای سایر ابعاد به ارائه نمودار بسنده شده است.

جدول ۴- نتایج حاصل از طبقه بندی کننده های SVM با ورودیهای شامل ۲۳ ویژگی با PCA

مدت زمان آزمایش (ثانیه)	مدت زمان آموزش (ثانیه)	G-mean	میانگین مربعات خطا	خطا	دقت	صحت	کلاس
----------------------------	---------------------------	--------	-----------------------	-----	-----	-----	------



۲۲,۳۳	۱۱۷۸	۹۹,۵۴	۰,۴۶	۰,۴۷	۹۹,۰۵	۹۹,۵۲	Normal
۲۲,۰۶	۱۰۶۵	۹۸,۵۱	۰,۳۸	۰,۳۸	۹۹,۰۹	۹۹,۶۱	Probe
۲۲,۵۰	۱۵۴۸	۹۹,۷۰	۰,۲۶	۰,۲۶	۹۹,۷۹	۹۹,۷۳	DOS

جدول ۵- نتایج حاصل از طبقه بندی کننده های SVM با ورودیهای شامل ۲۳ ویژگی با KernelPCA

مدت زمان آزمایش (ثانیه)	مدت زمان آموزش (ثانیه)	G-mean	میانگین مربعات خطا	خطا	دقت	صحت	کلاس
۳۶,۸۴	۹۵۱	۹,۷۳	۴۴,۰۲	۴۳,۹۵	۵۵,۹۷	۵۶,۰۴	Normal
۱,۷۱	۲۵,۲۲	۳۲,۶۱	۲,۱۶	۲,۱۶	۱۰۰	۹۷,۸۳	Probe
۳۴,۹۶	۱۰۳۷	۶,۳۶	۴۱,۷۳	۴۱,۶۶	۱۰۰	۵۸,۳۳	DOS

جدول ۶- نتایج حاصل از طبقه بندی کننده های SVM با ورودیهای شامل ۲۳ ویژگی با LDA

مدت زمان آزمایش (ثانیه)	مدت زمان آموزش (ثانیه)	G-mean	میانگین مربعات خطا	خطا	دقت	صحت	کلاس
۵۲,۳۰	۴۷۱,۸۳	۰	۴۷,۶۵	۴۸,۲۷	۰	۵۱,۷۲	Normal
۹,۸۴	۱۰۳,۶۱	۰	۱۰,۲۸	۱۰,۴۱	۰	۸۹,۵۸	Probe
۴۳,۵۰	۳۸۹,۷۰	۰	۳۹,۹۷	۴۰,۴۸	۰	۵۹,۵۱	DOS

جدول ۷- نتایج حاصل از طبقه بندی کننده های SVM با ورودیهای شامل ۲۳ ویژگی با ProbPCA

مدت زمان آزمایش (ثانیه)	مدت زمان آموزش (ثانیه)	G-mean	میانگین مربعات خطا	خطا	دقت	صحت	کلاس
۶,۸۳	۳۴۷,۷۷	۹۹,۸۷	۰,۱۳	۰,۱۳	۹۹,۸۶	۹۹,۸ ۶	Normal
۶,۸۴	۳۷۵,۳۲	۹۹,۸۵	۰,۰۳	۰,۰۳	۱۰۰	۹۹,۹ ۶	Probe
۵,۷۰	۳۴۹,۵۱	۹۹,۹۹	۰,۰۱	۰,۰۱	۹۹,۹۷	۹۹,۹ ۸	DOS

جدول ۸- نتایج حاصل از طبقه بندی کننده های SVM با ورودیهای شامل ۲۳ ویژگی با MDS



مدت زمان آزمایش (ثانیه)	مدت زمان آموزش (ثانیه)	G-mean	میانگین مربعات خطا	خطا	دقت	صحت	کلاس
۲۲,۳۲	۱۱۶۲	۹۹,۴۳	۰,۵۷	۰,۵۸	۹۸,۸۰	۹۹,۴۱	Normal
۲۱,۷۵	۱۱۸۸	۹۸,۱۷	۰,۴۳	۰,۴۴	۹۹,۲۹	۹۹,۵۵	Probe
۲۲,۶۴	۱۵۵۲	۹۹,۷۵	۰,۲۲	۰,۲۲	۹۹,۷۹	۹۹,۷۷	DOS

جدول ۹- نتایج حاصل از طبقه بندی کننده‌های SVM با ورودی‌های شامل ۲۳ ویژگی با FactorAnalysis

مدت زمان آزمایش (ثانیه)	مدت زمان آموزش (ثانیه)	G-mean	میانگین مربعات خطا	خطا	دقت	صحت	کلاس
۳۱,۸۱	۱۸۸۸	۹۶,۹۰	۳,۱۰	۳,۱۵	۹۳,۸۷	۹۶,۸۵	Normal
۳۱,۲۴	۲۰۵۵	۹۳,۶۱	۱,۳۱	۱,۳۳	۹۹,۴۴	۹۸,۶۶	Probe
۳۱,۸۶	۱۹۵۰	۹۶,۳۷	۲,۸۹	۲,۹۳	۹۹,۶۴	۹۷,۰۶	DOS

جدول ۱۰- نتایج حاصل از طبقه بندی کننده‌های SVM با ورودی‌های شامل ۲۳ ویژگی با Autoencoder

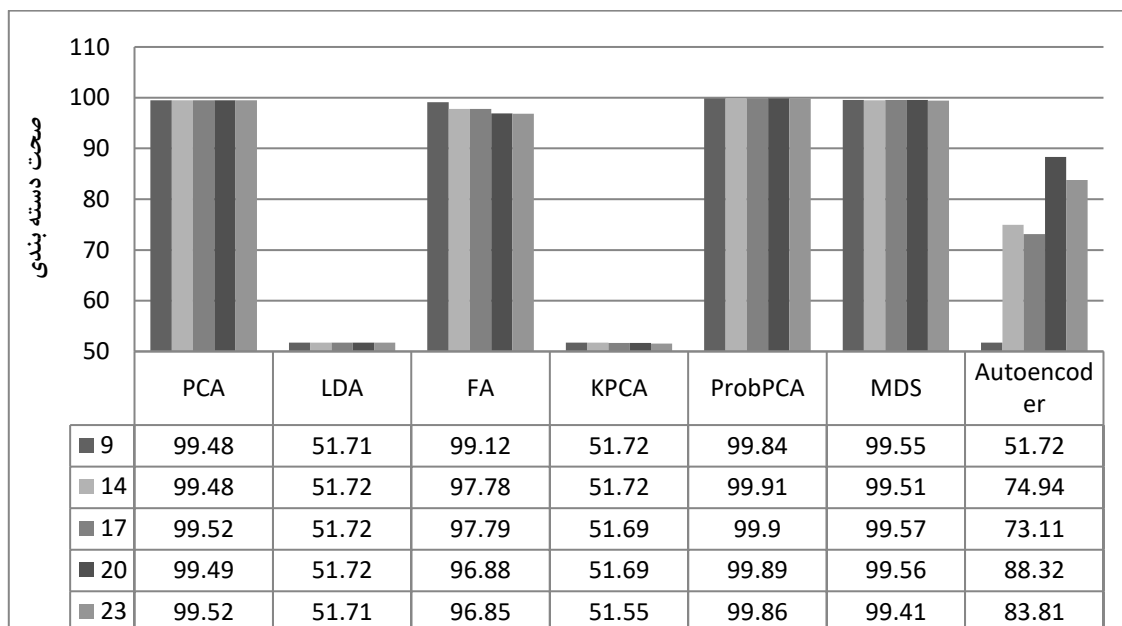
مدت زمان آزمایش (ثانیه)	مدت زمان آموزش (ثانیه)	G-mean	میانگین مربعات خطا	خطا	دقت	صحت	کلاس
۱۵,۶۹	۲۲۵,۲۹	۸۳,۵۸	۱۵,۹۷	۱۶,۱۸	۸۵,۸۲	۸۳,۸۱	Normal
۶,۳۶	۳۸۲,۶۷	۰	۱۰,۲۸	۱۰,۴۱	۰	۸۹,۵۸	Probe
۱۰,۲۷	۳۳۹,۰۹	۹۵,۸۰	۴,۷۸	۴,۸۴	۸۹,۵۶	۹۵,۱۵	DOS

همانطور که ملاحظه می‌شود نتایج حاصل از طبقه‌بندی SVM با ورودی‌های شامل ۲۳ ویژگی با استفاده از روش‌های کاهش ابعاد مختلف و محاسبه معیارهای ارزیابی در جداول (۴) تا (۱۰) گردآوری شده است. از مقایسه سه روش خطی PCA، LDA و FA نسبت به ۴۱ ویژگی چنین برمی‌آید که روش PCA بالاترین دقت طبقه‌بندی را در کلاس Dos دارا می‌باشد. با استفاده از این روش، با وجود اینکه ابعاد مسئله و تعداد ویژگی‌ها کاهش یافته است، دقت تشخیص و طبقه‌بندی کل سیستم نیز کاهش معناداری نیافته است. بنابراین، چنین نتیجه‌گیری می‌شود که روش PCA در کاهش ابعاد می‌تواند یکی از انتخاب‌های مناسب در جهت کاهش تعداد ویژگی‌ها باشد. از طرفی، همانطور که مشاهده می‌شود طبقه‌بندی SVM با ورودی ۲۳ ویژگی در روش LDA قادر به طبقه‌بندی نبوده است و تمامی نمونه‌ها برچسب کلاس

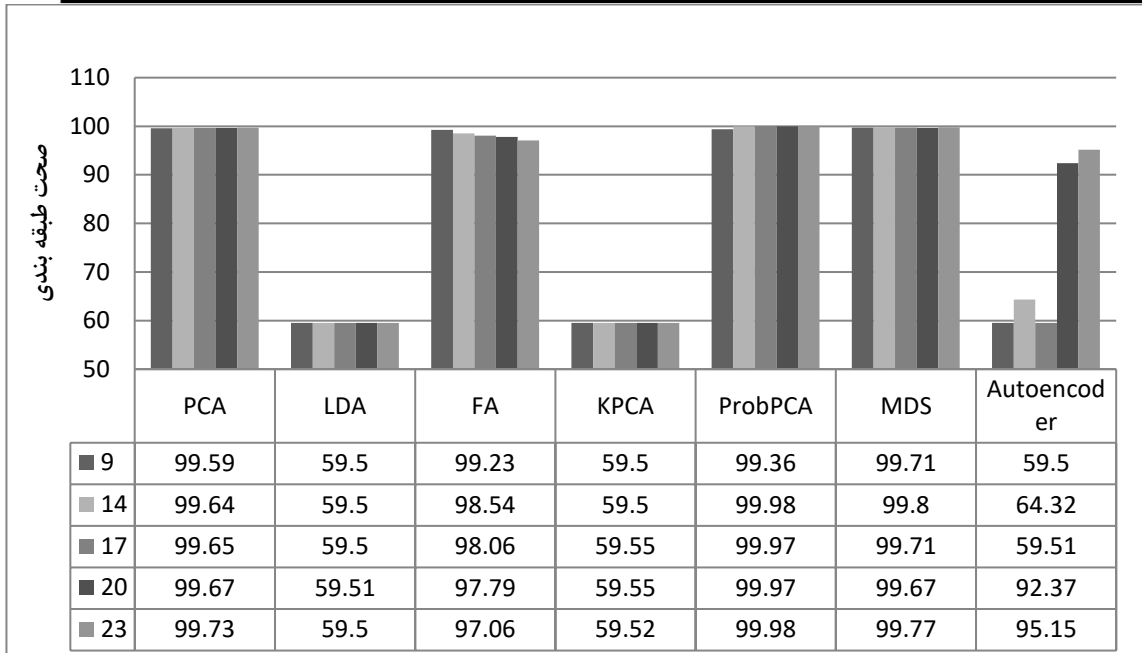
نرمال را دریافت کرده‌اند و به همین دلیل است که دارای معیار صحت کم و معیارهای دقت (Precision) و GMean صفر است، در واقع هیچ حمله‌ای تشخیص داده نشده است ( $TP^{\text{ن}}=0$ ).

در بررسی روش‌های غیرخطی، با توجه به کاهش شدید معیارهای ارزیابی نظیر دقت و صحت در روش Autoencoder، چنین استنتاج می‌گردد که کاهش ابعاد از طریق این روش کارایی مناسبی ندارد. بنابراین، هرچند کاهش ابعاد با این روش می‌تواند در کاهش حجم محاسبات و افزایش راندمان زمانی سیستم موثر باشد، اما به علت افت شدید دقت سیستم، استفاده از این روش مجاز نیست. همچنین با توجه به نتایج بدست آمده از آزمایشات، ProbPCA با بالاترین دقت و به دلیل احتمالاتی بودن این روش دارای کمترین زمان اجرا بوده و بهترین نتایج حاصل شده است و از آنجائیکه روش ProbPCA سرعت همگرایی کمی دارد بنابراین دقت آن بیشتر است. بعد از ProbPCA، روش مقیاس‌گذاری چندبعدی یا MDS بهترین نتایج را دارا می‌باشد. بالاترین دقت طبقه‌بندی کلاس نرمال در بین دیگر روش‌ها مربوط به روش غیرخطی ProbPCA می‌باشد.

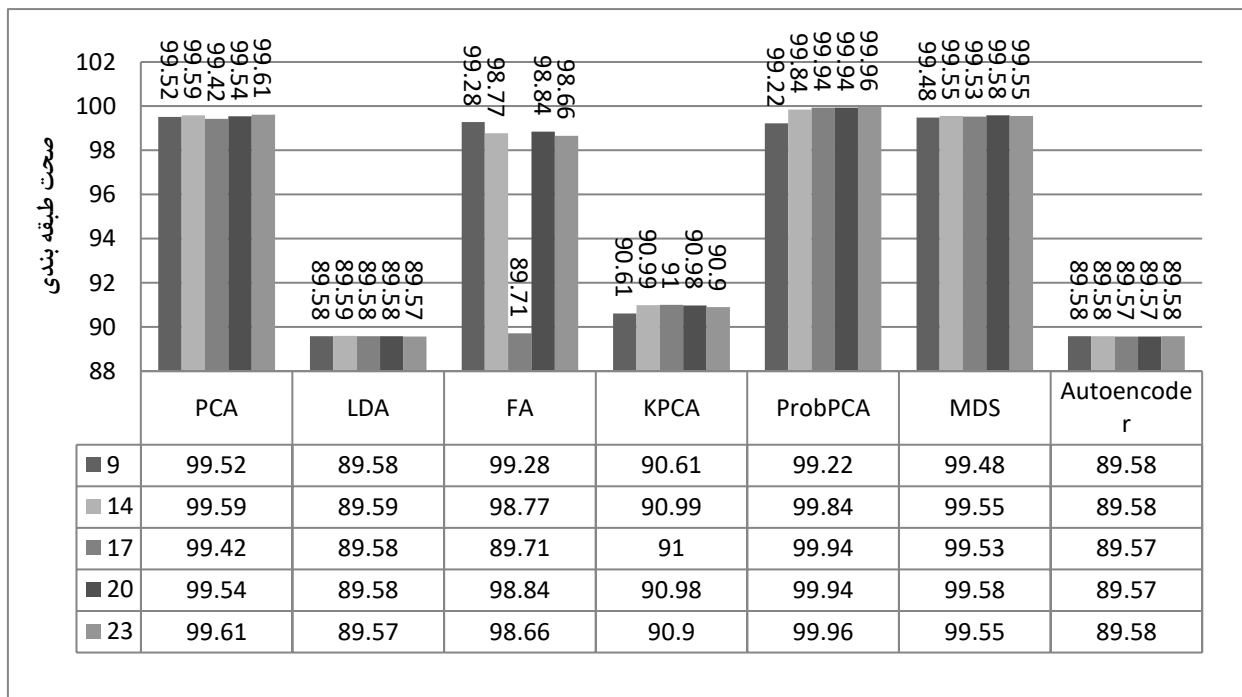
در ادامه معیارهای به‌دست آمده برای سایر ابعاد کاهش یافته بررسی می‌شود. مقادیر به‌دست آمده در آزمایش‌ها برای دو معیار دقت و GMean و همچنین دو معیار خطا و میانگین مربعات خطا بسیار به هم نزدیک هستند. لذا در ادامه معیارهای ارائه شده توسط نمودارهای شکل‌های (۲) تا (۷) برای کلاس‌های مختلف داده‌ها شامل صحت و GMean می‌باشد.



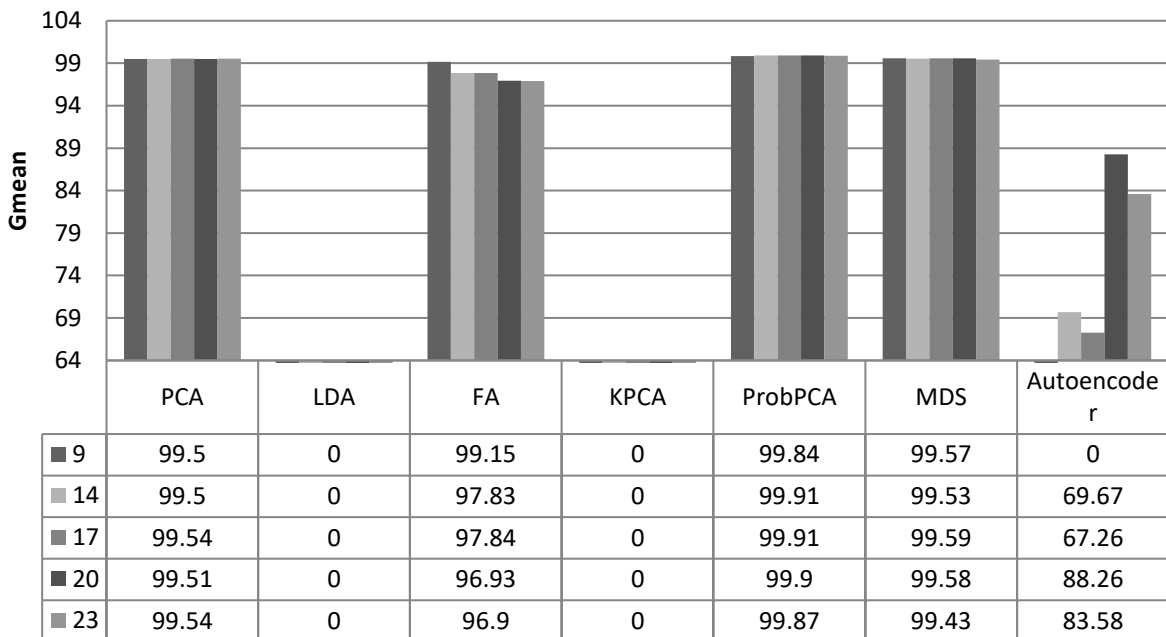
شکل ۲- صحت طبقه‌بندی آزمایش روش‌های مختلف در کلاس Normal



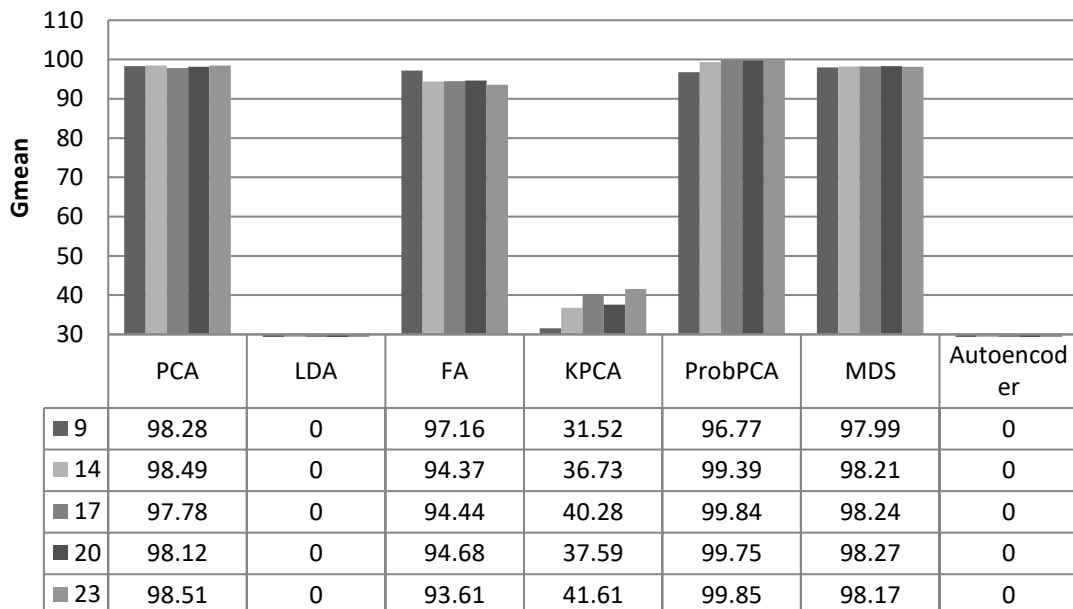
شکل ۳- صحت طبقه‌بندی آزمایش روش‌های مختلف در کلاس حمله DOS



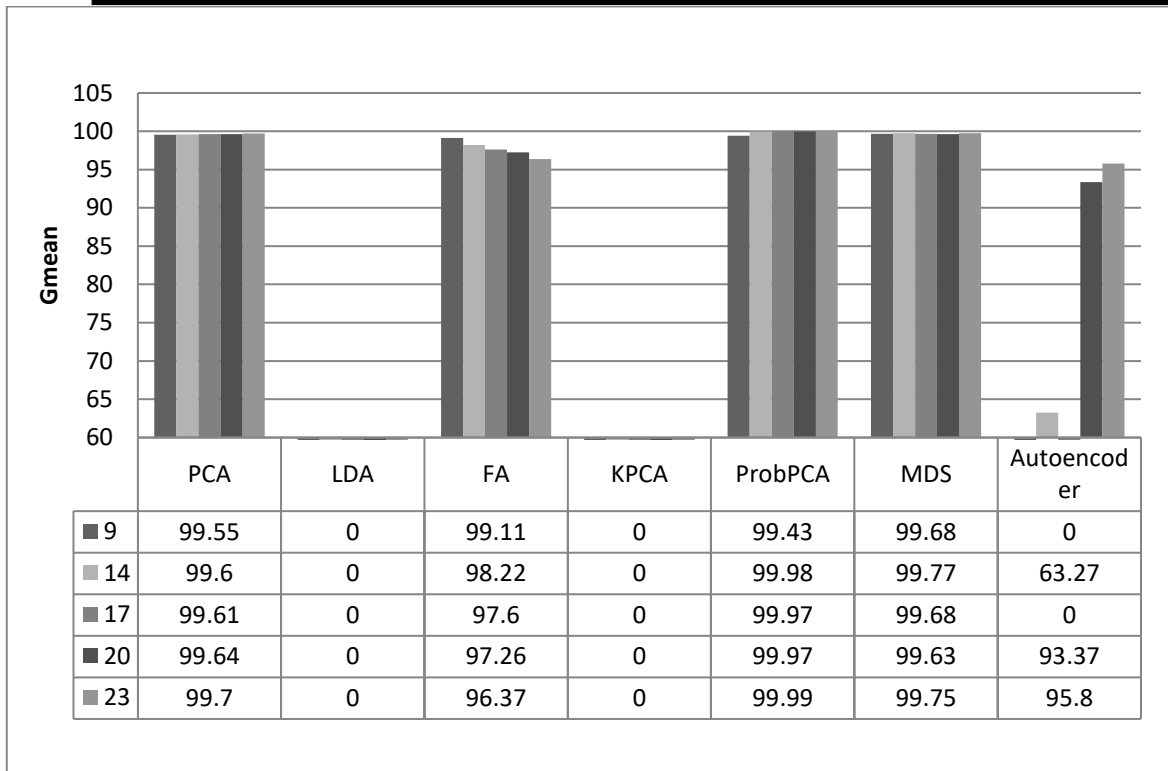
شکل ۴- صحت طبقه‌بندی آزمایش روش‌های مختلف در کلاس حمله Probe



شکل ۵- ارزیابی معیار GMean برای روش‌های مختلف در کلاس حمله Normal



شکل ۶- ارزیابی معیار GMean برای روش‌های مختلف در کلاس حمله Probe



شکل ۷- ارزیابی معیار GMean برای روش‌های مختلف در کلاس حمله DOS

با توجه به نتایج بدست آمده، تغییر میزان ابعاد در روش‌های LDA و MDS بر روی معیارهای ارزیابی تاثیر چشمگیری نداشته است. اما در روش‌هایی مانند PCA، Autoencoder و ProbPCA با تغییر ابعاد می‌توان به جواب بهینه دست یافت که در PCA و Autoencoder با افزایش ابعاد عملکرد بهتری دارد. در روش ProbPCA با کاهش ابعاد برای ۲۰ و ۱۷ بعد، برای کلاس حمله Probe شاهد کاهش صحت نیستیم اما برای تعداد ۱۴ معیار صحت دچار افت شده است. از سوی دیگر در FA بهترین صحت برای ۹ ویژگی است و هرچه ابعاد افزایش می‌یابد، صحت به مقدار ناچیزی کاهش یافته است. به‌طور کلی، با بررسی معیار صحت برای روش‌های مختلف می‌توان به این نتیجه رسید که تمام روش‌ها به جز روش‌های LDA و Autoencoder (برای کلاس‌های نرمال، Probe و DOS)، صحت بسیار بالایی دارند.

از نظر معیار Gmean، داده‌های کاهش یافته با روش LDA برای تمام ابعاد ضعیف است. می‌توان این‌گونه نتیجه گرفت که روش LDA سبب نامتوازن شدن داده‌ها می‌شود. از طرفی داده‌های کاهش یافته تمام ابعاد با روش‌های PCA، MDS، FA و ProbPCA برای تمام کلاس‌ها از معیار Gmean بالایی برخوردار است. به‌طور کل داده‌های کاهش یافته با روش Autoencoder نتایج قابل قبولی برای Gmean ارائه نداده است. همچنین با توجه به جداول و نمودارها می‌توان نتیجه‌گیری کرد که روش ProbPCA از روش‌های دیگر کاهش بعد مناسب‌تر است.

در جدول (۱۱) نتایج بدست آمده از روش کاهش بعد انتخابی که در مرجع [۱۴] پیاده‌سازی شده و همچنین نتایج حاصل از روش تبدیلی ProbePCA که بهترین روش تبدیل است، آورده شده است. در این جدول، در روش تبدیلی ProbePCA معیار صحت و نرخ تشخیص بهتر از روش انتخابی مرجع [۱۴] است ولی زمان آزمایش بدتر است که علت آن مدت زمانی است که سیستم عملیات محاسباتی برای تبدیل ۴۱ ویژگی به ۲۱ ویژگی کاهش بعد با روش ProbPCA را انجام می‌دهد.



جدول ۱۱- مقایسه کاهش بعد انتخابی در مرجع [۱۴] و روش کاهش بعد تبدیلی در این مقاله

ردیف	Dos		Probe		Normal		صحت	تعداد نمونه‌ها	روش کاهش بعد
	TP rate	FP rate	TP rate	FP rate	TP rate	FP rate			
۶،۱۳	۹۹،۹۴	۰،۰۱	۹۹،۰۱	۰،۰۵	۹۹،۹۵	۰،۲۷	۹۹،۹۱	۲۱	روش ProbPCA پیاپی‌سازی شده در این مقاله
۲،۸۴	۹۹،۱۴	۰،۱۴۹	۹۰،۳۵	۰،۰۴۰	۹۹،۰۷	۰،۱۷۴	۹۷،۶۷	۲۱	روش انتخابی در مرجع [۱۴]

##### ۵. نتیجه‌گیری

در این پژوهش، کاربردی‌ترین روش‌های خطی و غیرخطی کاهش ابعاد داده‌ها، بر روی مجموعه داده NSL-KDD انجام شد، سپس داده‌های کاهش شده حاصل از این روش‌ها، توسط طبقه‌بند ماشین بردار پشتیبان طبقه‌بندی و نتایج حاصل از آن‌ها مورد ارزیابی قرار گرفت. همچنین تاثیر حذف داده‌های پرت نیز بررسی شد، به این صورت که جهت دقت طبقه‌بندی داده‌ها ابتدا داده‌های پرت از آن‌ها جدا شده و سپس روش‌های کاهش ابعاد خطی و غیرخطی روی آن‌ها اعمال شده است.

از آنجائیکه بهترین درصدهای شناسایی نفوذ بوسیله سیستم تشخیص نفوذ با ویژگی‌های ProbPCA حاصل شده است، بنابراین می‌توان نتیجه‌گیری کرد که اطلاعات مفید بیشتری در ویژگی‌های بدست آمده از ProbPCA نسبت به روش‌های دیگر وجود دارد. ضمناً با ProbPCA زمان آموزش و آزمایش مناسب‌تری از روش‌های تبدیلی دیگر حاصل می‌گردد. بنابراین می‌توان نتیجه گرفت که ProbPCA مناسب‌تر از روش‌های دیگر می‌باشد زیرا هم دارای دقت بالاتر و هم زمان محاسباتی قابل قبول‌تر نسبت به سایر روش‌ها است.

در کارهای آینده می‌توان داده‌های کاهش ابعاد یافته در این مقاله (با کمک روش ProbPCA) را با دیگر الگوریتم‌های طبقه‌بندی و یا توسعه‌های بردار پشتیبان نیز مورد آزمایش قرار داد و نتایج طبقه‌بندی را مورد تحلیل و بررسی قرار داد.

##### ۶. مراجع

- [1] J. Andress, *THE BASICS OF INFORMATION SECURITY*: Elsevier, 2011.
- [2] H. Debar, M. Dacier, and A. Wespi, "Towards a taxonomy of intrusion-detection systems," *Computer Networks*, vol. 31, pp. 805-822, 1999.
- [3] S. X. Wu and W. Banzhaf, "The use of computational intelligence in intrusion detection systems: A review," *Applied Soft Computing*, vol. 10, pp. 1-35, 2010.
- [4] K. K. Vasan and B. Surendiran, "Dimensionality reduction using Principal Component Analysis for network intrusion detection," *Perspectives in Science*, vol. 8, pp. 510-512, 2016.
- [5] I. S. Thaseen and C. A. Kumar, "Intrusion Detection Model using fusion of chi-square feature selection and multi class SVM," *Journal of King Saud University-Computer and Information Sciences*, 2016.



- [6] J. Hussain and S. Lalmuanawma, "Feature Analysis, Evaluation and Comparisons of Classification Algorithms Based on Noisy Intrusion Dataset," *Procedia Computer Science*, vol. 92, pp. 188-198, 2016.
- [7] S. Balakrishnan, K. Venkatalakshmi, and A. Kannan, "Intrusion detection system using Feature selection and Classification technique," *International Journal of Computer Science and Application*, 2014.
- [8] S. V. Farrahi and M. Ahmadzadeh, "KCMC: A Hybrid Learning Approach for Network Intrusion Detection using K-means Clustering and Multiple Classifiers," *International Journal of Computer Applications*, vol. 124, 2015.
- [9] B. Kasliwal, S. Bhatia, S. Saini, I. S. Thaseen, and C. A. Kumar, "A hybrid anomaly detection model using G-LDA," in *Advance Computing Conference (IACC), 2014 IEEE International*, 2014, pp. 288-293.
- [10] W.-C. Lin, S.-W. Ke, and C.-F. Tsai, "CANN: An intrusion detection system based on combining cluster centers and nearest neighbors," *Knowledge-based systems*, vol. 78, pp. 13-21, 2015.
- [11] A. A. Aburomman and M. B. I. Reaz, "Ensemble of binary SVM classifiers based on PCA and LDA feature extraction for intrusion detection," in *Advanced Information Management, Communicates, Electronic and Automation Control Conference (IMCEC), 2016 IEEE*, 2016, pp. 636-640.
- [12] Z. Chen, H. Ren, and X. Du, "Minimax probability machine classifier with feature extraction by kernel PCA for intrusion detection," in *Wireless Communications, Networking and Mobile Computing, 2008. WiCOM'08. 4th International Conference on*, 2008, pp. 1-4.
- [13] E. B. M. Tavallae, W. Lu, and A. Ghorbani, "A Detailed Analysis of the KDD CUP 99 Data Set," *CISDA*, 2009.
- [14] R. Singh, H. Kumar, and R. Singla, "An intrusion detection system using network traffic profiling and online sequential extreme learning machine," *Expert Systems with Applications*, vol. 42, pp. 8609-8624, 2015.

<sup>1</sup> Denial-Of-service

<sup>2</sup> Intrusion Detection System

<sup>3</sup> Support Vector Machine

<sup>4</sup> Principal Component Analysis

<sup>5</sup> Performance-based Method of Ranking

<sup>6</sup> cluster center and nearest neighbor

<sup>7</sup> Linear Discriminant Analysis

<sup>8</sup> Kernel Principal Component Analysis

<sup>9</sup> Minimax Probability Machine Classifier

<sup>10</sup> Probabilistic Principal Component Analysis

<sup>11</sup> Multidimensional Scaling

<sup>12</sup> Factor Analysis

<sup>13</sup> True Positive