

بهبود سیستم های تشخیص نفوذ با کاهش ویژگی مبتنی بر الگوریتم ژنتیک و تکنیک های داده کاوی

مهدی کشاورزی^{۱*}، حسین مومن زاده حقیقی^۲

۱- دانشجو کارشناسی ارشد mehdikeshavarzi2000@gmail.com

۲- عضو هیئت علمی momenzadeh.hossein@gmail.com

چکیده

امروزه سیستم های کامپیوتری مبتنی بر شبکه، نقش حیاتی در جامعه مدرن امروزی دارند و به همین علت ممکن است هدف دشمنی و یا نفوذ قرار گیرند. به منظور ایجاد امنیت کامل در یک سیستم کامپیوتری متصل به شبکه، استفاده از دیوار آتش و سایر مکانیزم های جلوگیری از نفوذ همیشه کافی نیست و این نیاز احساس می شود تا از سیستم های دیگری به نام سیستم های تشخیص نفوذ استفاده شود. سیستم تشخیص نفوذ را می توان مجموعه ای از ابزارها، روش ها و مدارکی در نظر گرفت که به شناسایی، تعیین و گزارش فعالیت های غیرمجاز یا تأیید نشده تحت شبکه، کمک میکند. سیستم های تشخیص نفوذ به صورت سیستم های نرم افزاری و سخت افزاری ایجاد شده و هر کدام مزایا و معایب خاص خود را دارند. به دلیل وجود مشخصه های زیاد در داده های مربوط به سیستم های تشخیص نفوذ در این تحقیق ما مشخصه های مطلوب و موثر را با استفاده از الگوریتم ژنتیک بهبود یافته انتخاب می کنیم. سپس با استفاده از تکنیک های داده کاوی استاندارد، مدلی برای طبقه بندی داده ها ارائه می دهیم. برای ارزیابی عملکرد روش پیشنهادی از پایگاه داده *NSL-KDD* که نسبت به سایر داده های تشخیص نفوذ از رکوردهای واقعی تری برخوردار است، استفاده خواهیم کرد.

واژه های کلیدی: سیستم تشخیص نفوذ، داده کاوی، انتخاب ویژگی، الگوریتم ژنتیک، پایگاه داده *NSL-KDD*

۱- مقدمه

نفوذ به مجموعه اقدامات غیرقانونی که محرمانگی و یا دسترسی به یک منبع را به خطر می اندازد [۱] اطلاق می گردد. سیستم تشخیص نفوذ را می توان مجموعه ای از ابزارها، روش ها و مدارکی در نظر گرفت که به شناسایی، تعیین و گزارش فعالیت های غیرمجاز یا تأیید نشده تحت شبکه، کمک میکند. سیستم های تشخیص نفوذ به صورت سیستم های نرم افزاری و سخت افزاری ایجاد شده و هر کدام مزایا و معایب خاص خود را دارند. سرعت و دقت از مزایای سیستم های سخت افزاری است و عدم شکست امنیتی آن ها توسط نفوذگران، قابلیت دیگر این گونه سیستم ها می باشد. اما استفاده آسان از نرم افزار، قابلیت انطباق پذیری در شرایط نرم افزاری و تفاوت سیستم های عامل مختلف، عمومیت بیشتری را به سیستم های نرم افزاری می دهد و عموماً این گونه سیستم ها انتخاب مناسب تری هستند. افزایش دسترسی به داده ها و پردازش سریعتر آنها و در عین حال افزایش حجم داده ها و نیاز به فراهم آوردن داده ها از منابع مختلف از طریق شبکه های کامپیوتری، منجر به پدید آمدن منابع تهدید آمیزی می گردد که از طریق نقاط ضعف موجود در سیستم ها، به استثمار سیستم ها و ایجاد اختلال در آن ها می پردازد. تشخیص و جلوگیری از نفوذ (*IDS*) امروزه به عنوان یکی از مکانیزم های اصلی در برآوردن امنیت شبکه ها و سیستم های کامپیوتری مطرح است. سیستم های تشخیص نفوذ

متعددی برای کشف حملات وجود دارند که چالش اصلی آنها بالا بردن کارایی می‌باشد [۲]. بیشتر سیستم‌های تشخیص نفوذ 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. استفاده از تمامی پارامترها باعث می‌شود که فرآیند تشخیص طولانی و کارایی سیستم تشخیص نفوذ تنزل یابد. در واقع چالش اساسی در سیستم تشخیص نفوذ حجم عظیم داده‌هاست [۳]. از طرفی با توجه به ترافیک بالا کاهش نرخ هشدار غلط در سیستم تشخیص نفوذ نیز از اهمیت خاصی برخوردار است. تمام سیستم‌های تشخیص نفوذ قادر به تولید هشدار در مورد وقوع نفوذ در شبکه می‌باشند. ولی به علت حجم بالای هشدارهای تولید شده توسط این سیستم‌ها و همچنین تولید هشدارهای اشتباه، این سیستم‌ها قادر به مدیریت و آنالیز هشدارهای تولید شده نمی‌باشند. امروزه بیشتر رویکردها در تشخیص نفوذ مربوط به مساله انتخاب و یا استخراج ویژگی‌های مهم متمرکز شده است. تعداد ویژگی‌هایی که بهترین دقت را در تشخیص نفوذ دارند، در مجموعه داده‌های مختلف بیشتر به صورت تجربی بدست می‌آید. اما انتخاب ویژگی‌ها ممکن است باعث از دست دادن قسمتی از داده‌ها شود.

بنابراین با توجه به تفاوت داده‌های مربوط به شبکه، تخمین تعداد بهینه ویژگی‌ها یکی از چالش‌های مهم می‌باشد. در این تحقیق سعی می‌شود ویژگی‌هایی را به عنوان صفات اساسی انتخاب کنیم که صحت تشخیص نفوذ را بالا ببرد و ثابا تاثیر ویژگی‌هایی که انتخاب نشده‌اند را نیز در سیستم تشخیص نفوذ لحاظ نماییم. ما با استفاده از الگوریتم‌های ژنتیک و تکنیک‌های داده کاوی مدلی برای تشخیص نفوذ به شبکه ارائه دهیم.

۲- مروری بر کارهای انجام شده

بعد از سال ۱۹۷۰ و با افزایش تعداد کامپیوترها نیاز به امنیت کامپیوتری بیش از پیش آشکار شد. در همان زمان وزارت دفاع آمریکا، متوجه افزایش استفاده از کامپیوترها در سیستم‌های نظامی و در نتیجه آن اهمیت امنیت این سیستم‌ها شد. دنینگ و نیومن تحلیل و طراحی سیستم *IDES* که یک سیستم تشخیص نفوذ بلادرنگ و خبره بود را در سال ۱۹۸۵ به انجام رسانیدند [۴]. تحقیق انجام شده به عنوان پایه ای برای بسیاری از تحقیقات در زمینه تشخیص نفوذ در آن دهه قرار گرفت. تحقیقات بسیاری در زمینه کاهش ویژگی‌های سیستم‌های تشخیص نفوذ انجام شده است، که از بین آنها، برخی روش‌ها، مبتنی بر هوش مصنوعی و برخی براساس روش‌های داده کاوی می‌باشند. خدایار و همکارانش یک روش ترکیبی یادگیری ماشین به منظور تشخیص نفوذ در شبکه ارائه دادند [۵]. روش ترکیبی ارائه شده در این تحقیق مبتنی بر مفهوم کاهش ابعاد و الگوریتم‌های درخت تصمیم و روش‌های ترکیبی بوده و از دقت ۹۱/۷۹ درصدی برخوردار است. در تحقیقی دیگر در سال ۲۰۰۸ یک روش مبتنی بر الگوریتم ژنتیک برای تشخیص نفوذ روی پایگاه داده *KDDCUP99* ارائه شد [۶]. این روش از یک رویکرد یادگیری ماشین به همراه الگوریتم ژنتیک برای استخراج ویژگی استفاده می‌کند. با تولید یک سری قوانین و بر اساس هر قاعده یک نوع خاص نفوذ شناسایی می‌شود. رویکرد یادگیری ترکیبی خوشه بندی *K-means* و طبقه بندی بیز ساده در [۷] ارائه شد. خوشه بندی تمام داده‌ها را به گروه‌های مربوطه قبل از استفاده طبقه بندی نسبت می‌دهد. در این پژوهش از داده‌های *KDDCUP99* برای آموزش استفاده شده است و نتایج عملکرد بهتری از جهت دقت و سرعت کشف نفوذ با نرخ هشدار غلط مناسب را نشان می‌دهد. سربارنا ساها و همکارانش یک سیستم تشخیص نفوذ مبتنی بر یادگیری ماشین توسعه دادند [۸]. آنها الگوریتم ژنتیک را به همراه *SVM* برای تعیین اتوماتیک مجموعه مناسب از ویژگی‌ها به کار بردند. این مسئله از کاربرد الگوریتم ژنتیک برای انتخاب ویژگی و کلاس بندی مبتنی بر *SVM* حمایت می‌کند. در این تحقیق یک دیکشنری از پیش تعریف شده از انواع حملات وجود دارد که این روش متمایزترین ویژگی‌ها برای هر نوع حمله را در دسترس قرار می‌دهد. در [۹] با ترکیب ویژگی‌های نوع *filter* و *wrapper*، یک نوع روش ترکیبی برای انتخاب ویژگی با استفاده از یک الگوریتم ژنتیک توسعه یافته شامل مکانیزم تنبیه و پاداش ارائه شد. بناچار و همکارانش نیز سیستم تشخیص نفوذی با استفاده از

الگوریتم ژنتیک توسعه دادند که در آن ایده ای برای بهبود بخش های ایجاد جمعیت اولیه و اپراتور انتخاب در الگوریتم ژنتیک ارائه شده است [۱۰]. یک روش جدید برای انتخاب ویژگی های موثر در ساخت مدل تشخیص نفوذ در تحقیق [۱۱] ارائه شد. روش انتخاب ویژگی پیشنهادی با توجه به میانگین ویژگی ها از هر کلاس نسبت به کلاس کل محاسبه می شود. در این تحقیق از روش کلاس بندی درخت تصمیم استفاده می شود. در تحقیقی دیگر از الگوریتم کلونی مورچه ها، به منظور استخراج ویژگی ها در سیستم های تشخیص نفوذ استفاده کرده است [۱۲]. این روش بدلیل استفاده از زیر مجموعه های ساده برای کلاس بندی قابلیت اجرای سریع و پیچیدگی محاسباتی کمی می باشد. کاهش تعداد ویژگی ها با استفاده از نمایش گراف و اطلاعات اکتشافی برای به روز رسانی فرمون ها باعث ارائه دقت بالاتر در تشخیص نفوذ و هشدار اشتباه پایین تر می شود. به منظور مقابله با نفوذکنندگان به سیستم ها و شبکه های کامپیوتری توسط هر دو دسته کاربران داخلی و حمله کنندگان خارجی، ما در این تحقیق به بررسی رویکردی نوینی مبتنی بر الگوریتم ژنتیک در تشخیص و ممانعت از نفوذ و آینده این تکنولوژی می پردازیم.

۳- سیستم تشخیص نفوذ پیشنهادی

با توجه به اینکه در پردازش داده ها با اعداد سروکار دارم و برخی از مشخصه های این مجموعه داده دارای مقادیر رشته ای هستند، لذا تغییر آنها به مقادیر عددی برای انجام پردازش الزامی می باشد. برای تبدیل نوع های رشته ای به عدد نیز به این صورت عمل می کنیم که به ازای هر نوع یک عدد منحصر به فرد در نظر می گیریم. گام اول در ایجاد هر مدلی مبتنی بر تکنیک های داده کاوی، مرحله پیش پردازش داده ها می باشد. پیش پردازش جهت آماده سازی داده ها جهت پردازش و همچنین بهبود کیفیت داده های واقعی انجام می شود. این مرحله شامل نرمال سازی و در هم ریختن داده ها می باشد. در مرحله بعد با استفاده از الگوریتم های ژنتیک موازی ویژگی های مطلوب استخراج می شود. در برخی از مجموعه داده ها تعداد ویژگی ها زیاد می باشند و ممکن است برخی از این ویژگی ها در طبقه بندی داده ها نقشی نداشته باشند. بنابراین نیاز است که زیر مجموعه ای از بهترین ویژگی ها انتخاب شوند. روشی که برای انتخاب ویژگی ها در این پژوهش در نظر گرفته شده، براساس اندازه گیری توزیع ویژگی ها عمل می کند. به این صورت که در تکرارهای متوالی ویژگی های هنجار استفاده شده در کروموزوم های الگوریتم ژنتیک شناسایی شده و از آنها در ساخت راه حل های جدید استفاده می شود. روش ارائه شده می تواند روی مجموعه داده هایی با ابعاد متفاوت کار کند. بعد از انتخاب ویژگی ها برای ارزیابی دقت کلاس بندی از دو طبقه بند k نزدیک ترین همسایه (KNN) و درخت تصمیم (DT) استفاده می کنیم. در مرحله آخر، طبقه بندی داده ها با ویژگی های استخراج شده محاسبه می شود. منظور از تعداد ویژگی ها مشخص کردن تعداد ویژگی های انتخابی از کل ویژگی ها در مجموعه داده می باشد. بعد از تعیین ویژگی های مطلوب نهایی، دقت طبقه بندی ویژگی های انتخاب شده را روی بهترین روش از بین KNN و DT نشان می دهیم.

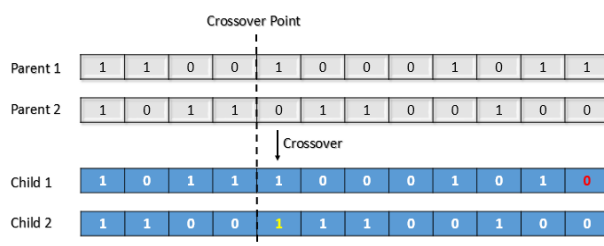
۳-۱- الگوریتم ژنتیک بهبود یافته برای انتخاب ویژگی ها

اولین مرحله در الگوریتم ژنتیک مشخص کردن ساختار کروموزوم ها می باشد. در واقع کد گذاری و نمایش کروموزوم ها به نوع مسئله و کاربرد آن بستگی دارد. در این مسئله شماره هر ویژگی را به عنوان معیار انتخاب آن ویژگی در نظر می گیریم. در این صورت ویژگی هایی که شماره آن در کروموزوم وجود داشته باشد در طبقه بندی داده ها شرکت خواهند داشت و بقیه در نظر گرفته نمی شوند. طول کروموزوم با توجه به تعداد کل ویژگی ها تعیین می شود و تعداد ژن های یک کروموزوم با توجه به تعداد ویژگی های مطلوب (DNF) مشخص می شود. در شکل ۱ نمونه ای از کروموزوم پیشنهادی را مشاهده می کنید.

۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰	۱۱	۱۲
۰	۱	۱	۰	۱	۱	۰	۰	۱	۰	۰	۱

شکل ۱: ساختار کروموزوم پیشنهادی

الگوریتم ژنتیک کار خود را با جمعیت اولیه از کروموزوم ها آغاز می کند. جمعیت اولیه به دو صورت تصادفی و هیوریستیکی از فضای جستجو تولید می شود. ما در این تحقیق کروموزوم ها را با هر دو روش از فضای جستجو انتخاب می کنیم. الگوریتم ژنتیک در ابتدا کار خود را با تولید جمعیت اولیه تصادفی شروع می کند. سپس با توجه به معیارهایی (در بخش های بعدی معرفی خواهد شد) ویژگی های هنجار و ناهنجار را از فضای جستجو استخراج کرده و از این پس تولید جمعیت جدید با توجه به این ویژگی ها انجام می پذیرد. در اینجا ویژگی های هنجار، برداری از ویژگی ها می باشد که در تولید جمعیت اولیه باید از آنها استفاده شود و ویژگی های ناهنجار برداری از ویژگی ها است که استفاده از آنها در تولید جمعیت های پیشین نتایج مطلوبی نداشته و از آنها در تولید جمعیت جدید استفاده نخواهد شد. با توجه به اینکه در روش پیشنهادی الگوریتم ژنتیک به صورت متوالی اجرا می شود، لذا در هر تکرار نیاز به تولید جمعیت اولیه از ویژگی ها می باشد، که در تکرار اول این جمعیت به صورت تصادفی و در سایر تکرارها به صورت هیوریستیکی از فضای جستجو تولید می شود. برای نمایش تعداد اعضای جمعیت از پارامتر NP استفاده خواهیم کرد. تابع برازندگی ($fitness$) در مسئله انتخاب ویژگی ها، نرخ خطا کلاس بندی داده ها می باشد. بنابراین نمونه های آموزشی را با توجه به هر کروموزوم (ویژگی های انتخاب شده) کلاس بندی کرده و خطا محاسبه شده را به عنوان معیار خوبی آن کروموزوم در نظر می گیریم. بنابراین مسئله ما مینیمم کردن تابع هدف می باشد. در این تحقیق از دو طبقه بند معروف $KNN [۱۳]$ و $DT [۱۴]$ جهت کلاس بندی داده ها و محاسبه معیار برازندگی کروموزوم ها استفاده شده است. دلیل استفاده از این دو طبقه بند سرعت بالای آنها در طبقه بندی پایگاه داده بزرگ می باشد. با توجه به اینکه در روش انتخاب ویژگی، نیاز به محاسبه متوالی برازندگی ویژگی های انتخاب شده وجود دارد، لذا استفاده از روش هایی که با سرعت بالا (نسبت به دقت) بتوانند طبقه بندی داده ها را انجام دهند، در اولویت است. به منظور انتخاب کروموزوم ها جهت تولید مثل از روش انتخاب تورنمنت استفاده می کنیم. مهمترین عملگر در الگوریتم ژنتیک، عملگر ترکیب (آمیزش) است. جفت هایی که در قسمت انتخاب، به عنوان والد در نظر گرفته شدند، در این قسمت ژن هایشان را با هم مبادله می کنند و اعضای جدیدی بوجود می آورند. شکل های مختلفی از عملگر ترکیب وجود دارد که از مهمترین آنها می توان به ترکیب تک نقطه ای، ترکیب دو نقطه ای، ترکیب n نقطه ای و ترکیب یکنواخت اشاره کرد. لازم به ذکر است که آمیزش معمولاً بر روی همه زوج کروموزوم های انتخاب شده برای جفت گیری به کار برده نمی شود. معمولاً احتمال آمیزش برای هر زوج کروموزوم بین $۰/۶$ تا $۰/۹۵$ در نظر گرفته می شود که به این عدد نرخ آمیزش (Pc) گفته می شود. در صورتی که بر روی یک زوج کروموزوم عمل آمیزش صورت نگیرد، فرزندان با تکرار نمودن والدین تولید می شوند. در این تحقیق از اپراتور ترکیب تک نقطه ای استفاده شده که نمونه ای از آن را در شکل ۲ مشاهده می کنید.



شکل ۲: اپراتور ترکیب تک نقطه ای بر روی کروموزوم پیشنهادی

همانطور که مشاهده می‌کنید دو مورد از ژن‌های فرزندان تغییر کرده است. در این مثال تعداد کل ویژگی‌ها (طول کروموزوم) ۱۲ و تعداد ویژگی‌های انتخاب شده (تعداد یک‌ها در کروموزوم) ۶ می‌باشد. تعداد ویژگی‌های انتخاب شده در تمامی کروموزوم‌ها برابر است. همانطور که در مثال بالا مشاهده می‌کنید، با نقطه برش انتخاب شده و با اعمال اپراتور ترکیب تک نقطه‌ای تعداد ویژگی‌ها در فرزند اول ۷ و در فرزند دوم ۵ می‌باشد. با توجه به موارد گفته شده باید تعداد ویژگی‌های انتخابی (ژن‌های یک کروموزوم) در فرزندان تولید شده برابر ۶ شود. لذا در فرزند اول که تعداد ویژگی‌های انتخابی بیشتر است یک مورد را به تصادف حذف می‌کنیم. در فرزند دوم نیز یک ویژگی استفاده نشده از والد اول به تصادف انتخاب کرده و در ژن متناظر قرار می‌دهیم. اپراتور ترکیب در روش پیشنهادی دو فرزند تولید می‌کند.

جهش یک فرایند تصادفی است که در آن محتوای یک ژن جهت تولید یک ساختار ژنتیک جدید جایگزین می‌شود. نقش جهش در الگوریتم ژنتیک بازگرداندن مواد ژنتیکی گم شده و یا پیدا نشده داخل جمعیت است. مزیت عملگر جهش اینست که توان دسترسی به همه فضای جستجو را به ما می‌دهد. در این قسمت از دو فرزند تولید شده توسط اپراتور ترکیب، یک مورد به تصادف انتخاب شده و جهش می‌یابد. در این تحقیق از اپراتور جهش تغییر بیت استفاده می‌کنیم. نمونه‌ای از این اپراتور در شکل ۳ نشان داده شده است. ایده اصلی این روش استفاده از ویژگی‌هایی می‌باشد که از آنها در راه حل پیشنهادی استفاده نشده است. در این اپراتور ابتدا یک ویژگی از کروموزوم (ژن با محتوای یک) انتخاب شده و حذف می‌گردد، سپس ویژگی جدید به کروموزوم اضافه می‌شود. در اپراتور ترکیب فقط ویژگی‌های بکار رفته در دو والد انتخابی در فرزند استفاده می‌شود. لذا بکارگیری ویژگی‌های مشاهده نشده در بهبود راه حل پیشنهادی لازم و ضروری می‌باشد. به یاد داشته باشید که ساخت جمعیت اولیه با توجه به ویژگی‌های هنجار و ناهنجار استخراج شده انجام می‌گرفت. لذا در این اپراتور نباید ویژگی‌هایی اضافه یا حذف شود که معیار ساخت جمعیت اولیه را نقض کند. بنابراین از بین ویژگی‌های بکار رفته در کروموزوم فقط ویژگی‌هایی می‌توانند برای حذف انتخاب شوند که در لیست ویژگی‌های هنجار نباشند. همچنین برای اضافه کردن ویژگی جدید، نباید از لیست ویژگی‌های ناهنجار موردی انتخاب شود.



شکل ۳: اپراتور جهش تغییر بیت بر روی یکی از فرزندان تولید شده

در این تحقیق به اندازه تعداد کروموزوم‌های جمعیت اولیه (NP) فرزند تولید می‌شود. سپس از میان جمعیت نسل قبل و جمعیت نسل فعلی تعدادی انتخاب شده و به نسل بعد می‌روند. برای انتخاب جمعیت نسل بعد، ابتدا کروموزوم‌های جمعیت نسل قبل و جمعیت نسل جاری بر اساس معیار برازندگی به صورت نزولی مرتب می‌شوند. سپس $1/4$ ابتدایی لیست (بهترین‌ها) به صورت مستقیم به نسل بعد می‌روند. کروموزوم‌های $1/4$ انتهایی لیست (بدترین‌ها) حذف می‌گردند و در نهایت سایر اعضای جمعیت بین کروموزوم‌های باقیمانده لیست، به صورت تصادفی انتخاب می‌شوند. یکی از پدیده‌های جالب در الگوریتم‌های ژنتیک این است که نسل‌های میانی کروموزوم‌هایی تولید می‌کنند که از نظر تابع ارزش و خوب بودن بسیار مناسب هستند. این کروموزوم‌ها ممکن است در اثر انجام اپراتورهای ترکیب و جهش از بین رفته و دیگر تولید نشوند. یک روش این است که اینگونه موارد را شناسایی کرده و در نسل‌های بعدی نیز از آن‌ها استفاده کنیم. به این تکنیک نخبه‌گرایی می‌گویند که عملاً تأثیر بسزایی در رسیدن به جواب مسئله دارد، چون این روش باعث می‌شود کروموزوم‌های خوب مستقیماً و بدون واسطه به نسل‌های بعدی بروند. در این تحقیق در هر

نسل تنها یک کروموزوم که بهترین راه حل را از نظر طبقه بندی پیشنهاد داده، مستقیماً به نسل بعدی منتقل می شود. نکته ای که در طراحی هر الگوریتمی بایستی به آن توجه شود، شرط خاتمه الگوریتم می باشد. در این تحقیق از تعداد مشخصی نسل برای شرط خاتمه الگوریتم ژنتیک استفاده میکنیم.

۳-۲- استفاده از فاکتور توزیع ویژگی ها در تشخیص ویژگی های هنجیده و ناهنجیده

در پایان الگوریتم ژنتیک، جمعیتی از راه حل ها بدست می آید. در اکثر روش های انتخاب ویژگی مبتنی بر الگوریتم ژنتیک، ویژگی ها بکار رفته در بهترین راه حل را به عنوان ویژگی های انتخاب شده در نظر گرفته و داده های آموزشی را بر مبنای این ویژگی ها کلاس بندی می کنند. ساختار الگوریتم ژنتیک بر مبنای جستجوی تصادفی می باشد، به همین دلیل همیشه راه حل بهینه و یکسانی را تولید نمی کند. به عنوان مثال اگر الگوریتم ژنتیک پیشنهادی را روی مجموعه داده تشخیص نفوذ با تعداد ثابتی از ویژگی ها در چندین اجرای مجزا اعمال کنیم، همیشه ویژگی های انتخابی در بهترین راه حل یکسان نیستند. با این شرایط پیدا کردن ویژگی های مطلوبی که بهترین عملکرد را در طبقه بندی داده ها دارند، امکان پذیر نخواهد بود. لذا در این تحقیق رویکردی را پیشنهاد می دهیم که تا حد زیادی به انتخاب بهترین ویژگی ها به ما کمک خواهد کرد. هدف ما در این بخش شناسایی ویژگی های هنجیده و ناهنجیده، با توجه به راه حل های بدست آمده از الگوریتم ژنتیک می باشد. برای شناسایی ویژگی های هنجیده و ناهنجیده از معیار فاکتور توزیع ویژگی ها (FD) در جمعیت استفاده می کنیم. توزیع هر ویژگی در جمعیت نشان دهنده میزان تکرار آن ویژگی می باشد. مفهوم توزیع ویژگی های هنجیده در جمعیت در واقع میزان تکرار آن ویژگی در بخش جمعیت هنجیده می باشد. به عنوان مثلاً جمعیت هنجیده راه حل هایی می باشند که برازندگی آنها از میانگین برازندگی کل جمعیت بیشتر باشد. همچنین مفهوم توزیع ویژگی های ناهنجیده در جمعیت، میزان تکرار آن ویژگی در قسمت ناهنجیده جمعیت می باشد. به عنوان مثلاً جمعیت ناهنجیده راه حل هایی می باشند که برازندگی آنها از میانگین برازندگی کل جمعیت کمتر باشد. فاکتور توزیع یک ویژگی در جمعیت هنجیده، نسبت تکرار آن ویژگی به کل جمعیت هنجیده می باشد و فاکتور توزیع یک ویژگی در جمعیت ناهنجیده، نسبت تکرار آن ویژگی به کل جمعیت ناهنجیده می باشد. شکل ۴ مفهوم فاکتور توزیع ویژگی ها را به خوبی نشان داده است.

Solution	f1	f2	f3	f4	f5	f6	f7	f8	Fitness
1	0	0	1	1	0	1	0	1	87
2	1	1	0	1	1	0	0	0	85
3	1	0	1	1	0	1	0	0	78
4	1	1	0	1	0	0	0	1	70
5	0	0	1	0	0	1	1	1	50
6	1	0	1	0	1	0	1	0	45
7	0	1	0	1	0	0	1	1	40
8	1	0	1	0	1	1	0	0	39
9	0	0	0	1	1	1	0	1	37
10	0	0	0	1	1	1	1	0	32

توزیع ویژگی ها در جمعیت بد

$$FD_1^{best} = \frac{3}{4} \quad FD_2^{best} = \frac{2}{4} \quad FD_1^{bad} = \frac{2}{6} \quad FD_2^{bad} = \frac{1}{6}$$

$$FD_3^{best} = \frac{2}{4} \quad FD_4^{best} = \frac{4}{4} \quad FD_3^{bad} = \frac{3}{6} \quad FD_4^{bad} = \frac{3}{6}$$

$$FD_5^{best} = \frac{1}{4} \quad FD_6^{best} = \frac{2}{4} \quad FD_5^{bad} = \frac{4}{6} \quad FD_6^{bad} = \frac{4}{6}$$

$$FD_7^{best} = \frac{0}{4} \quad FD_8^{best} = \frac{2}{4} \quad FD_7^{bad} = \frac{4}{6} \quad FD_8^{bad} = \frac{3}{6}$$

شکل ۴: مثالی از فاکتور توزیع ویژگی ها

همانطور که در مثال بالا مشاهده می کنید، با توجه به میانگین برازندگی جمعیت ۴ راه حل در جمعیت هنجیده و ۶ راه حل در جمعیت ناهنجیده فرض شده است. میزان تکرار ویژگی اول ($f1$) در جمعیت هنجیده ۳ و در جمعیت ناهنجیده ۲ می باشد. بنابراین فاکتور توزیع این ویژگی برای جمعیت هنجیده $FD_1^{best} = \frac{3}{4}$ و برای جمعیت ناهنجیده $FD_1^{bad} = \frac{2}{6}$ می باشد. در اینجا با توجه

به فاکتور توزیع ویژگی‌ها لیست ویژگی‌های هنجیده و ناهنجیده معرفی می‌شوند. به این صورت که بعد از اتمام الگوریتم ژنتیک، فاکتور توزیع را برای همه ویژگی‌ها محاسبه می‌کنیم. سپس ویژگی‌هایی که در جمعیت هنجیده فاکتور توزیع آنها از مقدار ثابتی مثل α بیشتر باشد به لیست ویژگی‌های هنجیده اضافه می‌شوند. همچنین ویژگی‌هایی که در جمعیت ناهنجیده فاکتور توزیع آنها از مقدار ثابتی مثل β کمتر باشد به لیست ویژگی‌های ناهنجیده اضافه می‌شوند. پارامتر α و β در واقع میزان تشابه راه حل‌ها و یا سخت‌گیری در انتخاب آنها را برای انتخاب یک ویژگی در جمعیت هنجیده و ناهنجیده کنترل می‌کنند. توجه داشته باشید که تعداد ویژگی‌های مطلوب در ابتدا به صورت ثابت تعریف شده و به عنوان ورودی به الگوریتم داده می‌شود. لذا برای تکمیل لیست ویژگی‌های هنجیده نیاز به تکرار مجدد الگوریتم ژنتیک و یافتن راه حل‌های جدید می‌باشد. بنابراین الگوریتم ژنتیک به صورت متوالی برای یافتن ویژگی‌های هنجیده تکرار شده تا زمانی که تعداد ویژگی‌های هنجیده استخراج شده به تعداد ویژگی‌ها مطلوب درخواستی برسد.

همانطور که در بخش‌های قبل نیز گفته شد، برای کمک به الگوریتم ژنتیک برای یافتن راه حل‌های بهینه‌تر، از لیست ویژگی‌های هنجیده و ناهنجیده در تولید جمعیت اولیه استفاده می‌کنیم. به طوری که تمامی جمعیت تولیدی شامل ویژگی‌های هنجیده استخراج شده باشند. همچنین از ویژگی‌های ناهنجیده استخراج شده در تولید جمعیت جدید استفاده نشود. این روش با ثابت قرار دادن تعدادی از ویژگی‌ها در هر تکرار، فضای جستجو را به میزان قابل توجهی کاهش می‌دهد. بعد از یافتن ویژگی‌های مطلوب، نرخ خطا کلاس بندی با دو کلاسیفایر معرفی شده محاسبه می‌شود.

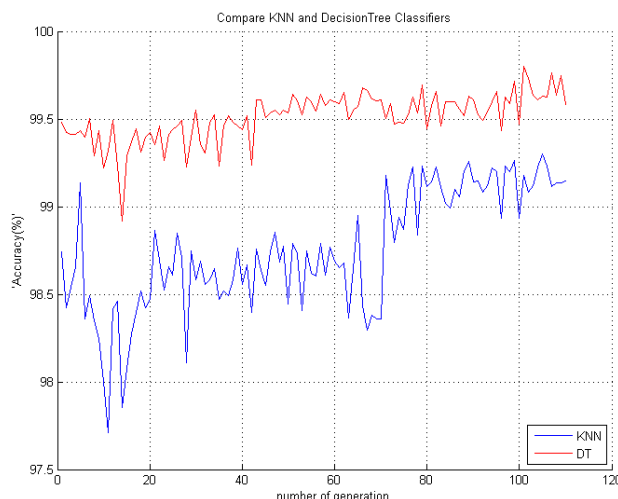
۳-۳- کنترل تطبیقی پارامترها

کنترل تطبیقی پارامترها در واقع روشی در تئوری کنترل به منظور تطبیق سیستم کنترلی با پارامترهایی در سیستم که مقدار آنها در طول فرآیند کاری سیستم تغییر می‌کند، است. مبنای کنترل تطبیقی براساس تخمین پارامترها می‌باشد. روش‌های مختلف تخمین شامل روش بازگشتی حداقل مربعات و روش گرادیان کاهشی می‌باشد. هر دوی این روش‌ها قوانین تطبیقی را به صورت آنلاین در سیستم تنظیم می‌کنند. در این تحقیق مقادیر دو پارامتر نرخ جهش (Pm) و میزان تشابه (α و β) در طی اجرای الگوریتم تغییر می‌کند. پارامتر نرخ جهش در ابتدای کار مقدار نسبتاً بالایی دارد و در طی روند اجرای الگوریتم به صورت متوالی از مقدار آن کاسته می‌شوند. پارامتر میزان تشابه نیز در ابتدا درصد بالایی از فضای انتخابی را شامل می‌شود. ولی هر قدر به جلو می‌رویم به دلیل تفاوت میان ویژگی‌های انتخابی از مقدار آن کاسته می‌شوند. این پارامتر در تعداد تکراری مشخص، که بهبودی در شناسایی ویژگی‌های هنجیده نداشته باشیم، به میزانی کاسته می‌شود. این روش تا حدودی مشکل همگرایی زودرس به نواحی بهینه محلی در الگوریتم ژنتیک با عملگرهای ژنتیکی نرخ ثابت را هم رفع می‌کند.

۴- آزمایشات و نتایج

در این تحقیق برای ارزیابی روش پیشنهادی از مجموعه داده **NSL-KDD** استفاده شده است. برای دسترسی کامل به این مجموعه داده می‌توانید از [۱۵] استفاده کنید. این مجموعه داده بهبود یافته مجموعه داده **KDDCUP'99** است و برای حل برخی از مشکلات ذاتی آن پیشنهاد شده است [۱۶]. مجموعه داده **NSL-KDD** از ۴۱ ویژگی و ۵ کلاس تشکیل شده است. رکوردها شامل یک کلاس نرمال و ۴ نوع کلاس حمله **U2R** **R2L** **Dos** و **Probing** هستند. برای پیاده‌سازی و آنالیز مجموعه داده‌ها از نرم افزار متلب ورژن ۲۰۱۳ استفاده شده است. نتایج بدست آمده از آزمایشات به منظور افزایش دقت ارزیابی، میانگین ۳۰ مرتبه تکرار تست می‌باشد. در پیاده‌سازی انجام شده اندازه جمعیت ۲۵، تعداد نسل ۳۰، نرخ ترکیب ۰/۸۵ و نرخ جهش ۰/۱۵ در نظر گرفته شده است. سخت‌گیری الگوریتم در انتخاب ویژگی‌های هنجیده $\alpha = 0.95$ و ناهنجیده $\beta = 0.90$ می‌باشد. با اعمال الگوریتم پیشنهادی

تعداد ۲۳ ویژگی به صورت خودکار از فضای جستجو انتخاب شد. ما کارایی دو طبقه بند معرفی شده را از لحاظ دقت روی هر کروموزوم تولید شده بررسی کرده و نتایج را در شکل ۵ نشان داده ایم. با توجه به نتایج الگوریتم انتخاب ویژگی، روش طبقه بندی DT کارایی بهتری نسبت به KNN داشته و به همین منظور نتایج طبقه بندی بر مبنای این روش است.



شکل ۵: مقایسه روش های طبقه بندی KNN و DT از لحاظ دقت تشخیص

۴-۱- معیارهای ارزیابی

ارزیابی یک مدل دسته بندی می تواند بر اساس نمونه های آموزشی و آزمایشی صورت گیرد. برای ارزیابی باید برچسبی که مدل دسته بندی به آن دسته حمله نسبت داده شده، مقایسه شود. وقوع حالات مختلف برای دسته ها با توجه به مجموعه داده های ورودی برای دسته بندی با مقادیر FN (False Negative)، TN (True Negative)، FP (False Positive) و TP (True Positive) برای سیستم تشخیص نفوذ با دو دسته در جدول ۱ نشان داده شده است. این جدول به ماتریس گیجی (Confusion) معروف است.

جدول ۱: ماتریس درهم ریختگی برای داده های سیستم تشخیص نفوذ

<i>Actual Records</i>	<i>Predicted Normal</i>	<i>Predicted Attack</i>
<i>Normal</i>	<i>TN</i>	<i>FP</i>
<i>Intrusions</i>	<i>FN</i>	<i>TP</i>

معیار TN : درصد رکوردهای معتبر که به درستی طبقه بندی شده است.

معیار TP : درصد رکورد حمله که به درستی طبقه بندی شده است.

معیار FN : درصد رکوردهای که به اشتباه به عنوان فعالیت درست قرار گرفتند در حالی که در واقع آنها حمله هستند.

معیار FP : درصد رکوردهای که به اشتباه به عنوان حمله قرار گرفتند در حالی که در واقع آنها فعالیت معتبر هستند.

مهمترین معیار برای تعیین کارایی یک الگوریتم دسته بندی، معیار $Accuracy$ می باشد. این معیار دقت کل یک دسته بند را محاسبه می کند. رابطه زیر نحوه محاسبه این معیار را نشان می دهد.

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

معیار **Precision** درصدی را نشان می‌دهد که از میان تمامی دسته‌ها که توسط دسته‌بند به آن دسته نسبت داده شده‌اند، درست دسته‌بندی شده‌اند. به عبارتی دقت دسته‌بندی دسته i را با توجه به کل مواردی نشان می‌دهد که برچسب i برای نمونه مورد بررسی توسط دسته‌بند پیشنهاد شده است. نحوه محاسبه این معیار در رابطه زیر نشان داده شده است.

$$Precision_i = \frac{TP_i}{TP_i + FP_i} \quad (2)$$

معیار **Recall** برای یک دسته، که از میان تمامی دسته‌های حملات متعلق به آن دسته، به درستی دسته‌بندی شده است. به عبارتی دقت دسته‌بندی دسته i را با توجه به کل نمونه‌های با برچسب i نشان می‌دهد. این معیار به صورت زیر محاسبه می‌شود.

$$Recall_i = \frac{TP_i}{TP_i + FN_i} \quad (3)$$

نکته قابل توجه این است که معیار **Recall** کارایی دسته‌بند را با توجه به تعداد رخداد دسته i نشان می‌دهد، درحالی‌که معیار **Precision** اساساً مبتنی بر دقت پیشبینی دسته می‌باشد و بیانگر آن است که به چه میزان می‌توانیم به خروجی دسته‌بند اعتماد کنیم.

معیار **F-measure** از ترکیبی معیارهای **Precision** و **Recall** بدست می‌آید و در مواردی استفاده می‌شود که نتوان اهمیت ویژگی‌های را برای هر یک از دو معیار **Precision** و **Recall** نسبت به یکدیگر قائل شد. رابطه زیر نحوه محاسبه این معیار را نشان می‌دهد.

$$F\text{-measure} = \frac{2 * Precision_i * Recall_i}{Precision_i + Recall_i} \quad (4)$$

کارایی سیستم‌های تشخیص نفوذ را میتوان با معیارهای معرفی شده مورد ارزیابی قرار داد.

۴-۲- ارزیابی روش پیشنهادی

بعد از پیدا کردن ویژگی‌های مطلوب توسط الگوریتم ژنتیک، اکنون نتایج طبقه‌بندی مجموعه داده **NSL-KDD** را مطابق با ویژگی‌های انتخاب شده بررسی می‌کنیم. جدول ۲ نتایج پیش‌بینی روش پیشنهادی را برای بخش‌های آموزش و تست به صورت جداگانه و در حالت دو کلاسه نشان می‌دهد. میزان دقت هر یک از ۴ دسته‌ی حملات **U2R**، **Probe**، **DOS** و **R2L** به همراه دسته‌نرمال محاسبه شده است. نتایج هر دسته را می‌توانید در جدول ۳ مشاهده کنید. ماتریس درهم‌ریختگی داده‌های سیستم تشخیص نفوذ برای هر یک از ۴ دسته‌ی حملات به همراه دسته‌نرمال محاسبه شده و در جدول ۴ نشان داده شده است. در این جدول تعداد رکوردها برای هر نوع حمله به همراه تعداد پیش‌بینی‌ها ذکر شده است. در جدول ۵ بهترین نتایج طبقه‌بندی روش پیشنهادی را با معیارهای متفاوت مشاهده می‌کنید. نتایج بر مبنای هر کلاس در مقابل سایر کلاس‌ها محاسبه شده است. در ادامه برای ارزیابی هر چه بیشتر رویکرد فوق، عملکرد سیستم پیشنهادی را با سایر روش‌های تشخیص نفوذ که نتایج آزمایشات خود را بر روی داده‌های **NSL-KDD** ارائه کرده‌اند، مقایسه می‌کنیم. نتایج جدول ۶ نشان می‌دهد که روش پیشنهادی در مقایسه با سایر روش‌های تشخیص نفوذ و به ازای برخی از حملات دقت بیشتری داشته و در سایر موارد دقت مناسبی را ارائه می‌دهد.

جدول ۲: نتایج روش پیشنهادی در حالت دو کلاسه

Dataset	Predicted Normal	Predicted Attack	Predicted
Train	99.94	99.89	99.92
Test	99.87	99.78	99.81

جدول ۳: درصد تشخیص روش پیشنهادی به تفکیک نوع حمله

Dataset	Normal	DOS	U2R	R2L	Probe
Train	99.94	99.96	84.62	98.59	99.80
Test	99.87	99.97	97.52	99.42	99.71

جدول ۴: ماتریس درهم ریختگی به تفکیک نوع حمله

Actual Records			Predicted				
Records Type	Dataset	Number	Normal	DOS	U2R	R2L	Probe
Normal	Train	67343	67303	6	7	11	16
	Test	9710	9683	3	2	7	1
DOS	Train	45927	8	45909	0	0	10
	Test	7458	1	7454	1	0	0
U2R	Train	52	7	0	44	1	0
	Test	200	2	0	197	3	0
R2L	Train	995	13	0	1	981	0
	Test	2754	9	3	2	2753	2
Probe	Train	11656	22	0	0	1	11633
	Test	2421	7	0	0	0	2413

جدول ۵: ارزیابی دقت روش پیشنهادی با معیارهای مختلف

Records Type		Accuracy	Recall	Precision	F-measure
Normal	Train	99.92	99.89	99.94	99.92
	Test	99.82	99.77	99.87	99.82
DOS	Train	99.93	99.89	99.96	99.93
	Test	99.85	99.73	99.97	99.85
U2R	Train	92.27	99.92	86.66	92.82
	Test	98.68	99.83	97.58	98.69
R2L	Train	99.26	99.93	98.61	99.27
	Test	99.64	99.86	99.42	99.64
Probe	Train	99.87	99.93	99.80	99.87
	Test	99.77	99.82	99.71	99.77

جدول ۶: مقایسه روش پیشنهادی با سایر روش‌های تشخیص نفوذ به تفکیک نوع حمله

Classifier	Normal	DOS	U2R	R2L	Probe	Accuracy
SIPSO [17]	-	99.80	97.50	82.50	99.70	-
Adaptive IDS [6]	66.51	88.64	66.51	20.88	99.15	75.15
GA [10]	-	-	-	-	-	68.51
CSM [11]	-	-	-	-	-	99.79
AdaBoost+ RandomTree +IG [5]	99.85	100.0	88.45	98.00	99.70	97.20
Ant Colony [12]	97.41	99.78	93.51	99.17	74.65	98.9
MARS [18]	99.71	99.97	76.00	98.75	99.85	92.75
Proposed Method	99.87	99.97	97.52	99.42	99.71	99.81

۵- نتیجه‌گیری

دقت الگوریتم‌های داده‌کاوی بستگی به انتخاب ویژگی‌های مناسب و همچنین تعداد رکوردهای مورد نظر برای یادگیری است. الگوریتم ژنتیک ارائه شده به خوبی ویژگی‌های مناسب را انتخاب می‌کند. با تنظیم دقیق و تطبیقی پارامترهای میزان تشابه، مجموعه ویژگی‌های هنجار و ناهنجار با سرعت بیشتری شناسایی شده و کارایی روش پیشنهادی بالا می‌رود. روش‌های داده‌کاوی معرفی شده نیز به خوبی دسته‌بندی داده‌ها را انجام می‌دهند. با این حال استفاده از الگوریتم‌های داده‌کاوی در سیستم تشخیص نفوذ دارای چالش‌هایی می‌باشند. داده‌ها و ترافیک تولید شده در شبکه، حجم بسیار بالایی دارند و تعداد ویژگی‌های آنها زیاد و اغلب ناهمگن هستند و نیاز به الگوریتم‌های داده‌کاوی با کارایی بالایی دارند. داده‌هایی که در شبکه سیر می‌کنند به صورت ذاتی حالت جریانی دارند و نیاز است تا بتوان آنها را به صورت برخط شناسایی کرد. حملاتی که امروزه بر علیه شبکه صورت می‌پذیرد، عموماً از چندین مبدا انجام می‌گیرد و هدف آنها نیز ممکن است چندین مقصد باشد. بنابراین نیاز است تا بتوان داده‌های شبکه را در چندین نقطه تحلیل و نفوذ حملات توزیعی را تشخیص داد.

۶- کارهای آینده

هدف این تحقیق، یافتن مجموعه ویژگی‌های بهینه برای سیستم‌های تشخیص نفوذ بوده است. با توجه به اینکه سیستم‌های تشخیص نفوذ، با داده‌های حجیم برای تحلیل مواجه هستند، کاهش مجموعه داده می‌تواند راه حل مناسبی برای افزایش دقت تشخیص آنها باشد. نیازمندی دیگری که در سیستم‌های تشخیص نفوذ مطرح می‌باشد، دانستن مجموعه ویژگی‌های بهینه برای هر نوع حمله است. چرا که در اینصورت، سیستم تشخیص نفوذ قادر خواهد بود برای تشخیص هر نوع حمله، تنها از مجموعه ویژگی‌های متناسب با آن حمله استفاده کند. همچنین در ادامه راه این تحقیق پیشنهاد می‌شود راه‌حلی به منظور یافتن تعداد ویژگی‌های مطلوب ابتدایی (DNF)، به صورت خودکار ارائه شود. به عنوان مثال می‌توان الگوریتم ژنتیک را با طول رشته متغیر بهبود بخشید. در این حالت هر کروموزوم در جمعیت دارای تعداد ویژگی‌های متفاوتی خواهد بود.

مراجع

- [1] Liao, Hung-Jen, Chun-Hung Richard Lin, Ying-Chih Lin, and Kuang-Yuan Tung. (2013), *Intrusion detection system: A comprehensive review. Journal of Network and Computer Applications* 36, no. 1 : 16-24.
- [2] Pan, Shengyi, Thomas Morris, and Uttam Adhikari.(2015). *Developing a hybrid intrusion detection system using data mining for power systems. IEEE Transactions on Smart Grid* 6.6 : 3104-3113.
- [3] Zuech, Richard, Taghi M. Khoshgoftaar, and Randall Wald.(2015). *Intrusion detection and big heterogeneous data: a survey. Journal of Big Data* 2.1 : 1.
- [4] Denning, Dorothy E., and Peter G. Neumann. (1985), *Requirements and model for IDES—a real-time intrusion detection expert system. Document A005, SRI International* 333.
- [5] خدایار، محمد؛ علیرضا عصاره و منصور امینی لاری، ۱۳۹۳، بکارگیری الگوریتم های ترکیبی یادگیری ماشین در بهبود سیستمهای تشخیص نفوذ، همایش ملی مهندسی رایانه و مدیریت فناوری اطلاعات، تهران، شرکت علم و صنعت طلوع فرزین.
- [6] Goyal, Anup, and Chetan Kumar, (2008). *GA-NIDS: a genetic algorithm based network intrusion detection system. Northwestern university.*
- [7] Muda, Z., W. Yassin, M. N. Sulaiman, and N. I. Udzir, (2011). *Intrusion detection based on K-Means clustering and Naïve Bayes classification. In Information Technology in Asia (CITA 11), 2011 7th International Conference on, pp. 1-6.*
- [8] Sriparna Saha, Ashok Singh Sairam, Asif Ekbal,(2012). *Genetic Algorithm Combined with Support Vector Machine for Building an Intrusion Detection System, International Conference on Advances in Computing, Communications and Informatics (ICACCI-2012)*
- [9] Zhu, Shuxin, and Bin Hu, (2013). *Hybrid feature selection based on improved GA for the intrusion detection system. Indonesian Journal of Electrical Engineering and Computer Science* 11, no. 4 : 1725-1730.
- [10] Benaicha, Salah Eddine, Lalia Saoudi, Salah Eddine Bouhouita Guermeche, and Ouarda Lounis, (2014). *Intrusion detection system using genetic algorithm. In Science and Information Conference (SAI), 2014, pp. 564-568.*
- [11] Chae, Hee-su, Byung-oh Jo, Sang-Hyun Choi, and Twaekyung Park, (2015). *Feature Selection for Intrusion Detection using NSL-KDD. Recent Advances in Computer Science, ISBN : 978-960.*
- [12] Aghdam, Mehdi Hosseinzadeh, and Peyman Kabiri, (2016). *Feature selection for intrusion detection system using ant colony optimization. International Journal of Network Security* 18.3 : 420-432.
- [13] Park, Chan Hee, and Seoung Bum Kim. (2015), *Sequential random k-nearest neighbor feature selection for high-dimensional data." Expert Systems with Applications* 42.5 : 2336-2342.
- [14]
- [15] Hidayati, R., Kanamori, K., Feng, L., & Ohwada, H. (2016). *Combining Feature Selection with Decision Tree Criteria and Neural Network for Corporate Value Classification. In Pacific Rim Knowledge Acquisition Workshop (pp. 31-42). Springer International Publishing.*
- [16] Nsl-kdd dataset for network based intrusion detection systems. Available on: <http://nsl.cs.unb.ca/KDD/NSL-KDD.html>, March 2009.
- [17] Tavallae M, Stakhanova N and Ghorbani AA., (2010). *Towards credible evaluation of anomaly based intrusion detection methods , IEEE Transaction on System, Man and Cybernetics, Part-c, Applications and Reviews; 40(5):516-524.*
- [18] Warsi, Sana, Yogesh Rai, and Santosh Kushwaha, (2015). *Selective Iteration based Particle Swarm Optimization (SIPSO) for Intrusion Detection System. International Journal of Computer Applications* 124.17.
- [19] Mubarak, Shaik Liyakhat, (2016). *Intrusion Detection System using SVM, SOM & NN.*

improving intrusion detection systems by feature reducing based on genetics algorithm and data mining techniques

Mehdi Keshavarzi^{1*}, hossein Momenzadeh^{2*}

^{1*} Mehdikeshavarzi2000@gmail

^{2*} Momenzadeh.hossein@gmail

ABSTRACT:

The network-based computer systems play critical role in our modern society; so there is highly chance these systems might be target of intrusion and attacks. In order to implement full-scale security in a computer network, firewalls and other intrusion prevention mechanisms aren't always enough and needs other systems called intrusion detection systems. An Intrusion detection system can be set of tools, algorithms and evidence that help to identify, locate and report illegal or not approved activities by the network. Intrusion detection systems can be established by software or hardware systems and each have their own advantages and disadvantages. Because of various characteristics of intrusion detection data, in this research we select effective characteristics using improved genetic algorithm. Then by means of standard data mining techniques, we present a model for data classification. For performance evaluation of this suggested method, we used **NSL-KDD database that has more realistic records than other intrusion detection data.**

KEYWORDS: system intrusion detection, data mining, feature selection, genetic algorithms, database NSL-KDD