

## افزایش نرخ کشف نفوذ به شبکه‌های کامپیوتری با استفاده از درخت‌های تصمیم

پریسا جعفرزاده<sup>1</sup>، شهرام جمالی<sup>2</sup>

<sup>1</sup> کارشناسی ارشد، دانشگاه آزاد اسلامی، واحد زنجان، گروه کامپیوتر، زنجان، ایران، [jafarzadeh\\_parisa@yahoo.com](mailto:jafarzadeh_parisa@yahoo.com)

<sup>2</sup> دانشیار دانشگاه محقق اردبیلی، گروه مهندسی کامپیوتر، اردبیل، ایران، [jamali@uma.ac.ir](mailto:jamali@uma.ac.ir)

### چکیده

امنیت، یکی از مسائل مهم در سیستم‌های رایانه‌ای مدرن امروزی است؛ از جمله چالش‌های مهم در این سیستم‌ها، تشخیص نفوذ است؛ نفوذ، به عنوان مجموعه‌ای از فعالیت‌ها تعریف می‌شود که هدف این فعالیت‌ها، به خطر انداختن یکپارچگی، قابلیت اطمینان سیستم و دسترسی غیرمجاز به منبع خاص می‌باشد. لذا سیستم‌های تشخیص نفوذ یکی از تکنیک‌هایی است که برای حفظ امنیت در شبکه‌های رایانه‌ای استفاده می‌شود. یکی از مشکلات این سیستم‌های امنیتی، گزارش نادرست هشدار نفوذ به سیستم است؛ این مقاله با استفاده از درخت تصمیم یک مکانیزم کشف نفوذ به شبکه طراحی می‌کند. نتایج شبیه‌سازی نشان می‌دهد که نه تنها در این روش مدت زمان فاز آموزش کاهش می‌یابد، بلکه نرخ هشدار نادرست و نرخ تشخیص نفوذ نیز بهبود نسبی پیدا می‌کند.

**کلید واژه‌ها:** داده‌کاوی، درخت تصمیم، سیستم تشخیص نفوذ.

را مورد شناسایی قرار می‌دهد و فعالیت‌های لازم را برای توقف حملات اتخاذ می‌کند. مسأله‌ی تشخیص نفوذ به طور وسیع در زمینه‌ی امنیت سیستم‌های رایانه‌ای مورد مطالعه قرار گرفته است و اخیراً در زمینه‌های یادگیری ماشین و داده‌کاوی توجه زیادی به آن شده است. سیستم‌های تشخیص نفوذ، حملات شناخته شده و حملات بالقوه را در ترافیک شبکه و یا داده‌های بازرسی<sup>2</sup> و یا فایل‌های گزارش (ثبت وقایع<sup>3</sup>) ثبت شده به وسیله‌ی میزبان‌ها جستجو می‌کند و اقدامات لازم را جهت حفاظت از سیستم انجام می‌دهد. لازم به ذکر است که اولین منبع اطلاعات که سیستم‌های تشخیص نفوذ از آنها استفاده می‌کنند داده‌های بازرسی سیستم‌عامل می‌باشند. داده‌های بازرسی توسط بخشی به نام بخش بازرسی که

### 1. مقدمه

رشد روز افزون استفاده از خدمات شبکه‌های رایانه‌ای از یک سو و حمله به این شبکه‌ها از سوی دیگر باعث شده است که تشخیص نفوذ به عنوان یک زمینه‌ی تحقیقاتی مهم در مسأله‌ی امنیت این شبکه‌ها مطرح شود.

سیستم تشخیص نفوذ<sup>1</sup>، ابزار امنیتی مؤثر و کارآمدی است که در شبکه‌های رایانه‌ای قرار می‌گیرد و با استفاده از یک سری قوانین از پیش تعریف شده، دسترسی کاربران را بررسی و محدود می‌نماید. این قوانین بر پایه‌ی دانش افراد متخصص حاصل شده است.

سیستم‌های تشخیص نفوذ، تمام نفوذهای کاربران

زیرمجموعه‌ی سیستم عامل است تهیه می‌شود. داده‌های بازرسی دربرگیرنده‌ی اطلاعاتی درباره‌ی فعالیت‌های سیستم می‌باشد. این اطلاعات برحسب زمان مرتب شده‌اند و دریک یا چند فایل به نام فایل بازرسی ذخیره می‌شوند. هر فایل بازرسی حاوی مجموعه‌ای از رکوردهای بازرسی است که هر یک بیانگر یک رویداد در سیستم هستند. این رکوردها توسط فعالیت‌های کاربر یا فرآیندها ایجاد می‌شوند.

اولین نیاز سیستم‌های تشخیص نفوذ، وجود منبع اطلاعات است. این منبع اطلاعات، به عنوان تولیدکننده‌ی رویداد در نظر گرفته می‌شود. درسیستم‌های تشخیص نفوذ، این منابع اطلاعاتی با توجه به مکانشان دسته‌بندی می‌شود. با توجه به این معیار، دو دسته‌بندی کلی برای سیستم‌های تشخیص نفوذ وجود دارد [1]:

مبتنی بر شبکه<sup>4</sup>: در این دسته، بسته‌های عبوری در سطح شبکه به عنوان منبع اطلاعات جمع‌آوری می‌شوند؛ این عمل، با قرار دادن کارت شبکه در حالت بی‌قاعده<sup>5</sup> صورت می‌گیرد. در این حالت چون سرآیند<sup>6</sup> بسته‌های شبکه مورد بررسی قرار می‌گیرد، محدودیت نگاه به سرآیند بسته‌ها باعث می‌شود این سیستم‌ها قابلیت تشخیص دامنه‌ی محدودتری از حملات را داشته باشند؛ لازم به ذکر است که منظور از حالت بی‌قاعده این است که در حالت طبیعی، یک سیستم بر روی شبکه، تنها ترافیکی که به طور مستقیم به آن ارسال می‌شود را می‌خواند و پاسخ می‌دهد، ولی در حالت بی‌قاعده، سیستم، تمام ترافیک شبکه را می‌خواند.

مبتنی بر میزبان<sup>7</sup>: داده‌های این نوع سیستم‌ها از فعالیت‌های متنوع کامپیوترهای میزبان که شامل رکوردهای بازرسی سیستم‌عاملی، فایل‌های ثبت رخداد سیستمی، اطلاعات برنامه‌های کاربردی و ... هستند، فراهم می‌شود. اطلاعات از منابع موجود در این میزبان‌ها جمع‌آوری می‌شود. همان‌گونه که اشاره شد، این منابع شامل فایل‌های ثبت وقایع و داده‌های

بازرسی هستند.

از آنجا که این سیستم‌ها از منابع خود دستگاه استفاده می‌کنند، باعث ایجاد سربار اضافی بر روی دستگاه می‌شوند، ولی از نظر تئوری، قابلیت تشخیص تقریباً تمام حملات را دارند. این سیستم‌های تشخیص نفوذ می‌توانند در مورد نفوذی بودن یا نبودن یک کامپیوتر میزبان به دقت داوری کنند، زیرا داده‌های آنها از داده‌های مورد بازرسی سیستم و فایل‌های ثبت رخدادها در سیستم‌عامل به دست می‌آیند و در مقایسه با سیستم‌های تشخیص نفوذ مبتنی بر شبکه می‌توانند حملات داخل شبکه و میزبان‌های نفوذی را شناسایی نمایند.

بعد از معرفی منابع اطلاعاتی، باید نوع تحلیل‌گر شبکه مشخص شود. در تحلیل‌گر، اطلاعات از منابع اطلاعاتی استخراج می‌شوند و با توجه به سیاست‌های امنیتی، انواع حملات مورد بررسی قرار می‌گیرند. بر اساس تحلیل‌گر، دو نوع دسته‌بندی داریم:

1) تشخیص رفتار غیرعادی: در این گونه سیستم‌ها تشخیص حمله به این صورت است که رفتار عادی سیستم را مدل می‌کنیم و رفتار و قواعد خاصی برای آنها در نظر گرفته می‌شود؛ رفتارهایی که از این الگوها پیروی کنند، عادی می‌باشد و رفتارهایی که مغایر با این الگوها هستند، به عنوان رفتار غیرعادی در نظر گرفته می‌شوند. این سیستم‌ها معمولاً دارای تعداد هشدار نادرست بالا هستند، ولی قابلیت تشخیص حمله‌های ناشناخته را نیز به سیستم می‌دهند. در این گونه دسته‌بندی، سیستم یک مرز نرمال از الگوهای متعارف استفاده از سیستم برای خود ترسیم و تعریف می‌کند. هر رفتار و رویدادی که به میزان زیادی از این الگوها دور باشد به عنوان اخلاک‌گر محتمل بر شبکه در نظر گرفته می‌شود. آنچه خلاف قاعده و نابهنجار فرض می‌شود می‌تواند متغیر باشد. اما معمولاً رویدادی که با تناوبی بیشتر یا کمتر از دو مرتبه انحراف از آمار نرمال به وقوع بپیوندد نابهنجار فرض می‌شود. چند مثال برای این شرایط عبارت است از:

## 2. داده‌کاوی

استخراج دانش از اطلاعات می‌باشد. در اینجا منظور از استخراج اطلاعات، دستیابی به اطلاعاتی است که قبلاً بدیهی نبوده و برای ما ناشناخته بوده‌اند. قبل از آغاز اجرای عملیات داده‌کاوی بر روی اطلاعات، لازم است که پیش‌پردازش‌هایی روی داده‌های خام و انبوه انجام دهیم؛ این عملیات پیش‌پردازش، شامل انتخاب خصیصه‌های مفید از داده‌ها، کاهش حجم و بُعد داده‌ها به منظور کاهش زمان محاسبات داده‌کاوی، نرمال‌سازی داده‌ها، گسسته‌سازی داده‌ها و ... می‌باشد [3].

بعد از انجام کاوش و استخراج الگوهای مفید از داده‌ها، ممکن است نیاز به انجام پردازش‌هایی روی این الگوها باشد. مجموعه‌ی این اعمال را استخراج دانش از پایگاه‌های داده نامیده می‌شود.

همان‌گونه که در شکل 1 مشاهده می‌شود، فرآیند کشف دانش متشکل از چند مرحله می‌باشد که از داده‌های خام، اطلاعات مفید را به دست می‌دهد. این مراحل عبارتند از:

پاکسازی داده‌ها<sup>15</sup>: این مرحله با عنوان پیرایش داده‌ها نیز شناخته می‌شود. در این فاز، داده‌های اضافی و نامربوط از مجموعه‌ی داده‌های موجود حذف می‌شوند.

تمامیت داده‌ها<sup>16</sup>: در این مرحله، تصمیم‌گیری بر روی داده‌ها متناسب با تحلیل، انجام می‌شود و این داده‌ها از مجموعه داده‌ها استخراج می‌شوند.

تبدیل صورت داده‌ها<sup>17</sup>: این مرحله با عنوان تثبیت داده نیز شناخته می‌شود که در این فاز، داده‌های انتخاب‌شده به صورتی متناسب جهت فرآیند داده‌کاوی تبدیل می‌شوند.

داده‌کاوی<sup>18</sup>: این مرحله بسیار مشکل می‌باشد و در آن از تکنیک‌های هوشمند برای استخراج الگوهای مفید بالقوه، استفاده می‌شود.

ارزیابی ارزش الگوها<sup>19</sup>: در این مرحله الگوهای صریح به منظور کشف دانش، شناسایی می‌شوند.

کاربری که به جای یک یا دو بار ورود و خروج از سیستم در طول روز 20 بار این کار را انجام داده است. کامپیوتری که در ساعت 2 بعد از نیمه شب مورد استفاده قرار گرفته است، در صورتی که قرار نبوده بعد از ساعت اداری روشن باشد. در یک سطح دیگر این تکنیک می‌تواند الگوهایی در مورد کاربران از جمله برنامه‌هایی که به اجرا در می‌آورند را مورد استفاده قرار دهد. مثلاً اگر کاربری از بخش گرافیک یک سازمان ناگهان شروع به دستیابی به برنامه‌های حسابداری یا کامپایل کردن کد بنماید، سیستم می‌تواند یک هشدار به مدیر مسؤو امنیت<sup>8</sup> شبکه ارسال کند [2].

2) تشخیص سوء استفاده<sup>9</sup>: در این حالت، الگوهای نفوذ از پیش ساخته‌شده، به عنوان قانون نگهداری می‌شوند. در واقع، در این حالت، هدف، تشخیص حملات شناخته‌شده به همراه تغییرات کوچک در آنها می‌باشد. از مزایای این روش، دقت در تشخیص حملات شناخته‌شده و کاهش هشدارهای نادرست می‌باشد. این سیستم، توانایی تشخیص حملات ناشناخته را ندارد [2].

سیستم‌های تشخیص نفوذ، از روش‌های مختلفی همانند درخت‌های تصمیم‌گیری<sup>10</sup>، الگوریتم ژنتیک<sup>11</sup>، منطق فازی<sup>12</sup>، خوشه‌بندی<sup>13</sup>، شبکه‌ی بیزین<sup>14</sup> و ... برای کشف نفوذ استفاده می‌کند.

مابقی این مقاله به این ترتیب سازماندهی شده است: در بخش 2، به طور مختصر به مفاهیم داده‌کاوی و چرایی استفاده از آن می‌پردازیم. در بخش 3، کاربرد داده‌کاوی را در ساخت سیستم‌های تشخیص نفوذ کارآمد بیان می‌کنیم. در ادامه و در بخش 4، به راهکار یادگیری ماشین که ارتباط نزدیکی با تکنیک‌های داده‌کاوی دارد، اشاره می‌کنیم. همچنین، به الگوریتم C4.5 می‌پردازیم که یکی از روش‌های یادگیری ماشین است. در بخش 5، به شبیه‌سازی و ارزیابی نتایج حاصل می‌پردازیم. در نهایت، در بخش 6، نتیجه‌گیری خود را از این مقاله ارائه می‌کنیم.

اکتشاف دانش<sup>20</sup>: در این مرحله نهایی، دانش مورد نظر کشف شده، به صورت بصری به کاربر ارائه می شود. این مرحله مهم از تکنیک های موصّرسازی برای کمک به کاربران در درک و تفسیر نتایج داده کاوی، استفاده می کند.

فرآیند اکتشاف علوم، یک فرآیند تکرارشونده است. زمانی که علوم کشف و به کاربران ارائه شد، سنجش هایی جهت ارزشیابی انجام می شود و امکان پاک شدن مجدد علوم به دست آمده در طی عملیات داده کاوی با داده های جدید و یا ترکیبی از منابع داده ای جدید وجود دارد و ارائه ی نتایج بهتری را به ارمغان می آورد.

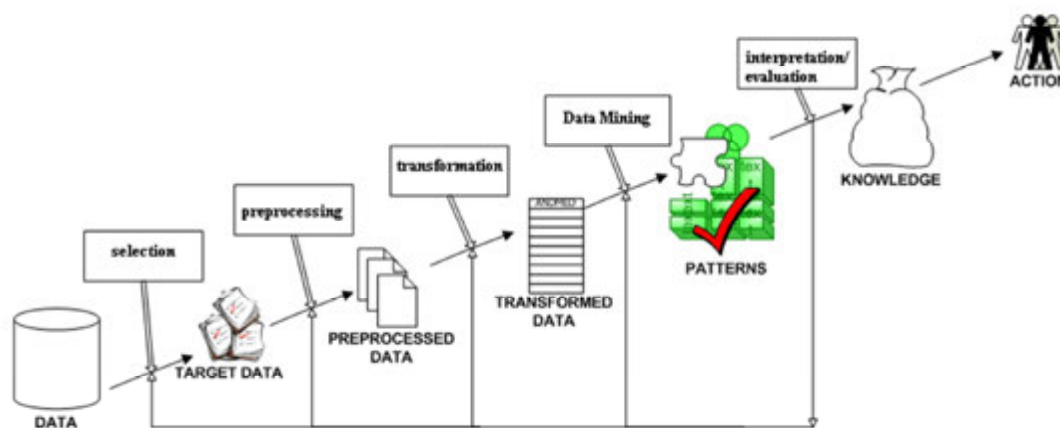
### 3. کاربرد داده کاوی در تشخیص نفوذ

نفوذ، در حقیقت به مجموعه اقدامات غیرقانونی اطلاق می شود که صحت، محرمانگی و یا دسترسی به یک منبع را به خطر می اندازد. عوامل نفوذی را می توان به دو گروه خارجی و داخلی تقسیم کرد نفوذی های خارجی کسانی هستند که اجازه ی استفاده از منابع سیستم را ندارند، اما سعی می کنند آنها را مورد دستیابی قرار دهند و نفوذی های داخلی کسانی هستند که برای دستیابی به سیستم اختیارات محدودی دارند، اما سعی می کنند به منابعی که اجازه ی استفاده از آنها را ندارند، دسترسی پیدا کنند. به منظور مقابله با نفوذکنندگان به سیستم ها و شبکه های کامپیوتری توسط هر دو دسته کاربران داخلی و حمله کنندگان خارجی، روش های متعددی پیشنهاد شده است که

تکنیک های تشخیص نفوذ نامیده می شوند. هدف از تشخیص نفوذ این است که استفاده ی غیرمجاز، سوءاستفاده و آسیب رساندن به سیستم ها و شبکه های کامپیوتری شناسایی و با آنها مقابله شود. یکی از پرکاربردترین تکنیک های تشخیص نفوذ، یادگیری ماشین می باشد که در بخش 4 به طور مختصر معرفی می شود.

### 4. یادگیری ماشین

یادگیری ماشین، شاخه ای از هوش مصنوعی است که در حقیقت برخی از تکنولوژی هایی را که در آنها کامپیوتر، شرایط یادگیری اتوماتیک را دارد، گسترش می دهد. بحث اصلی در یادگیری ماشین این است که چگونه از کامپیوترها و روش های آماری استفاده کنیم تا بتوانیم اطلاعات مفید را از میان حجم انبوه داده ها استخراج کنیم. بنابراین یادگیری ماشین و داده کاوی هر دو با روش های آماری و تئوری دانش کامپیوتر در ارتباط هستند. تکنولوژی های رایج یادگیری ماشین شامل موارد زیر هستند [5]: تئوری تصمیم بیز<sup>21</sup>، روش های چند متغیری<sup>22</sup>، خوشه بندی، درخت های دسته بندی<sup>23</sup>، تفکیک خطی<sup>24</sup>، مشاهدات چندلایه<sup>25</sup>، مدل های محلی<sup>26</sup>، مدل های مارکوف پنهان<sup>27</sup>، یادگیری تقویت<sup>28</sup>. تکنیک های یادگیری ماشین در شرایطی مناسب است که هیچ گونه دانش اولیه در مورد الگوهای داده ها وجود ندارد؛



شکل (1): مراحل داده‌کاوی و استخراج دانش [4]

این سؤال، شاخه وجود دارد که هر شاخه با مقدار آن جواب مشخص می‌شود. برگ‌های این درخت با یک کلاس و یا یک دسته از جواب‌ها مشخص می‌شوند. مهمترین و اساسی‌ترین الگوریتم‌هایی که در درخت‌های تصمیم استفاده می‌شوند ID3 و C4.5 که هر دو متعلق به ساختمان داده‌ی درخت هستند و الگوریتم آنها به صورت زیر می‌باشد [6]:

تمام نمونه‌های داده‌های آموزشی در مکان ریشه‌ی درخت قرار می‌گیرند.

اگر گرهی دارای هیچ داده‌ای نباشد و یا داده‌های یک گره از نوع یکسان باشند، در این صورت گره، یا به عنوان یک گره‌ی برگ خالی می‌باشد و یا به عنوان برگ‌ی از نمونه‌های یکسان در نظر گرفته می‌شود. اگر گره دارای نمونه‌های بیشتر از یک نوع باشد، باید بتوانیم طبق یک تابع ارزیابی معین به تمام خصیصه‌های داده‌ها دست یابیم و بتوانیم خصیصه‌های مناسب را انتخاب کنیم. بر مبنای تعداد این خصیصه‌ها، نمونه‌های یک گره به  $N$  قسمت تقسیم می‌شوند و هر قسمت به عنوان یک گره‌ی جدید است که به گره‌ی ریشه متصل شده است. این فرآیند، تقسیم گره<sup>30</sup> نامیده می‌شود.

بعد از تقسیم گره، تصمیم‌گیری در مورد برگ بودن گره‌های جدید انجام می‌گیرد. اگر برگ نباشند، این گره‌ها به عنوان ریشه برای زیردرخت‌های جدید

به همین دلیل گاهی به این روش‌ها، پائین به بالا می‌گویند. مزیت مهم این تکنیک‌ها این است که معمولاً به انسان‌های خبره برای تعیین ملزومات مورد نظر به منظور تشخیص نفوذ نیازی نیست به همین دلیل بسیار سریع عمل می‌کنند و مقرون به صرفه هستند.

در بخش بعد به بررسی درخت تصمیم می‌پردازیم و نوآوری خود را که در مورد انتخاب ویژگی‌های تأثیرگذار مجموعه داده‌ی به کار رفته در فاز آموزش است تشریح می‌کنیم.

## 5. درخت دسته‌بندی

دسته‌بندی، یک مدل پیشگویی در یادگیری ماشین است و درخت تصمیم نیز نامیده می‌شود که یک گراف با الگوی درختی و شبیه به ساختمان داده‌ی نمودار گردشی<sup>29</sup> است.

درخت تصمیم درختی است که در آن نمونه‌ها را به نحوی دسته‌بندی می‌کند که از ریشه به سمت پائین رشد می‌کنند و در نهایت به گره‌های برگ می‌رسد: هر گره‌ی داخلی یا غیربرگ بایک ویژگی مشخص می‌شود. این ویژگی سؤالی را در رابطه با مثال ورودی مطرح می‌کند.

در هر گره‌ی داخلی به تعداد جواب‌های ممکن با

3) قابلیت ترکیب با روش‌های دیگر: نتیجه درخت تصمیم را می‌توان با تکنیک‌های تصمیم‌سازی دیگر ترکیب کرده و نتایج بهتری به دست آورد.

با توجه به مزایای فوق و بررسی کارهای انجام‌شده [2] در رابطه با تشخیص نفوذ متوجه شدیم که درخت تصمیم دارای نرخ هشدار نادرست کم و نرخ تشخیص نفوذ بالایی است؛ ولی ایراد این روش این است که مدت زمان آموزش این الگوریتم در فاز آموزش نسبتاً زیاد است بنابراین بایستی حجم مجموعه داده‌ای که در فاز آموزش برای ساخت مدل استفاده می‌شود کاهش یابد. به این منظور سعی کردیم که از بین تمام ویژگی‌های موجود در مجموعه داده‌ی به کار رفته در فاز آموزش درخت تصمیم، تنها مهمترین و تأثیرگذارترین ویژگی‌ها را انتخاب کنیم و در ساخت مدل به کار ببریم.

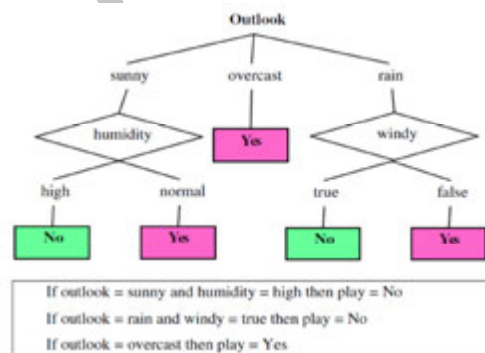
همان‌گونه که قبلاً نیز عنوان شد برای کاهش مدت زمان ساخت مدل در این الگوریتم و همچنین افزایش سرعت و جامعیت سیستم مورد نظر، به جای استفاده از تمام 41 ویژگی موجود در مجموعه داده‌ی KDDCup99، از تعدادی ویژگی موجود در این مجموعه داده استفاده کردیم که این ویژگی‌ها با استفاده از الگوریتم‌های انتخاب ویژگی موجود در نرم‌افزار داده‌کاوی Weka انتخاب شدند. این تعداد ویژگی قابل تغییر می‌باشد و بدیهی است که انتخاب ویژگی‌های مناسب، تأثیر به‌سزایی بر عملکرد خوب سیستم تشخیص نفوذ دارد. نحوه‌ی عملکرد هر کدام از این الگوریتم‌ها به طور مختصر شرح داده می‌شود و برای توضیحات بیشتر به سایت [7] Weka مراجعه شود:

- CfsSubsetEval: در این روش، ارزش پارامترها به طور مکاشفه‌ای با استفاده از همبستگی ارزیابی می‌شود. در این روش بهترین زیرمجموعه‌ی ویژگی‌ها، ویژگی‌هایی هستند که بیشترین همبستگی را برای پیش‌بینی ویژگی هدف (کلاس) دارند. اما از سوی دیگر، درجه‌ی

قرار می‌گیرند و برای تشکیل زیردرخت‌های جدید استفاده می‌شوند.

مراحل فوق به روش بازگشتی بسط داده می‌شوند، تا زمانی که تمام گره‌های جدید به عنوان گره‌های برگ باشند. تابع ارزیابی نیز معمولاً شامل تابع آنتروپی<sup>31</sup>، بهره‌ی اطلاعاتی<sup>32</sup>، تابع جینی<sup>33</sup>، پیش‌هرس درخت<sup>34</sup> و پس‌هرس درخت<sup>35</sup> می‌باشد. بعد از محاسبه‌ی مقادیر این توابع، ویژگی با بیشترین مقدار انتخاب می‌شود.

شکل 2، ساختار یک درخت تصمیم‌گیری و قوانین استخراج شده از آن را نشان می‌دهد. خوانندگان این مقاله برای دریافت اطلاعات بیشتر می‌توانند به مرجع [6] مراجعه نمایند.



شکل 2: ساختار یک درخت تصمیم‌گیری و قوانین استخراج شده از آن [6]

از مزایای درخت‌های دسته‌بندی می‌توان به موارد زیر اشاره کرد:

- 1) کارکردن با داده‌های بزرگ و پیچیده: درخت تصمیم در عین سادگی می‌تواند با داده‌های پیچیده به راحتی کار کند و از روی آنها تصمیم بسازد.
- 2) استفاده‌ی مجدد آسان: در صورتی که درخت تصمیم برای یک مسأله ساخته شد، نمونه‌های مختلف از آن مسأله را می‌توان با آن درخت تصمیم محاسبه کرد.

سپس معیار بهره‌ی اطلاعات را محاسبه می‌کند. از توانایی‌های دیگر این الگوریتم، برخورد با ارزش‌های مفقود شده به عنوان یک ارزش جداگانه می‌باشد. در این روش، ارزش مفقود شده با توجه به نسبت فراوانی ارزش‌های دیگر محاسبه می‌گردد.

بهره‌ی اطلاعات، قدرت دسته‌بندی یک ویژگی را، در مجموعه‌ی داده‌ای بیان می‌کند. هدف از دسته‌بندی آن است که مجموعه‌ی ورودی به دسته‌هایی تقسیم گردد که هر دسته تنها معرف یک ارزش باشد. به عبارت دیگر خلوص هر دسته بیشینه باشد. از آنجا که ویژگی‌های دارای بهره‌ی اطلاعات بیشتر، مجموعه داده‌ای را به تعداد قسمت‌های یک متر با حداکثر خلوص تقسیم می‌کنند، بنابراین این ویژگی‌ها برای دسته‌بندی، ویژگی‌های مناسبی به شمار می‌روند.

#### 6. شبیه‌سازی و ارزیابی نتایج

برای شبیه‌سازی، از نرم‌افزار داده‌کاوی Weka استفاده کردیم. این نرم‌افزار، یک واسط همگون برای بسیاری از الگوریتم‌های یادگیری متفاوت، فراهم کرده است که از طریق آن روش‌های پیش‌پردازش، پس‌پردازش و ارزیابی نتایج طرح‌های یادگیری روی همه‌ی مجموعه داده‌های موجود، قابل اعمال است. این سیستم به زبان جاوا نوشته شده است و تقریباً روی هر پلتفرمی اجرا می‌شود و پلت‌فرم جاوایی Weka باعث می‌شود تا اجرای آن روی بیشتر سیستم‌عامل‌ها امکان‌پذیر باشد. بستر آزمایشی که شبیه‌سازی را روی آن انجام دادیم یک رایانه با سیستم‌عامل لینوکس و نسخه‌ی فدورا 11 بود. این رایانه دارای حافظه‌ی 2 گیگاهرتزی و پردازنده‌ی دو هسته‌ای شرکت اینتل می‌باشد. به منظور شبیه‌سازی نیاز به دو مجموعه داده داریم. یکی از این مجموعه داده‌ها را برای آموزش دادن الگوریتم‌ها استفاده می‌کنیم و مجموعه داده‌ی دیگر را برای آزمایش کردن الگوریتم‌ها استفاده می‌کنیم.

همبستگی بین خود ویژگی‌های موجود در زیر مجموعه کم می‌باشد که از فرمول رابطه‌ی 1 برای به دست آوردن درجه‌ی شایستگی این زیرمجموعه‌ها استفاده می‌شود [8]:

$$\text{Merit} = \frac{\text{درجه‌ی همبستگی ویژگی مورد نظر با ویژگی هدف}}{\text{درجه‌ی همبستگی ویژگی مورد نظر با ویژگی‌های موجود در زیرمجموعه}} \quad (1)$$

- **ChiSquaredAttributeEval**: این الگوریتم برای رتبه‌بندی ویژگی‌ها از معیار Chi-Squared استفاده می‌کند. ویژگی‌هایی که رتبه‌ی بالاتری دارند، از قدرت جداکنندگی بیشتری نیز برخوردار هستند.
- **ClassifierSubsetEval**: این الگوریتم ارزش ویژگی‌ها را با استفاده از کلاس‌بندی که کاربر مشخص می‌کند، ارزیابی می‌کند. به عنوان مثال، کاربر کلاس‌بند OneR را انتخاب می‌کند.
- **InfoGainAttributeEval**: این روش، یکی از پرکاربردترین روش‌های استاندارد در انتخاب ویژگی‌های مرتبط با ویژگی هدف می‌باشد که میزان شایستگی ارتباطی که ویژگی با ویژگی هدف را با اندازه‌گیری معیار بهره‌ی اطلاعات<sup>36</sup> مشخص می‌کند. هرچه میزان بهره‌ی اطلاعات یک ویژگی با ویژگی هدف نسبت به بقیه بیشتر باشد، ارتباط ویژگی مورد نظر نیز با ویژگی هدف نسبت به بقیه بیشتر می‌باشد. به بیان دیگر، ویژگی هدف را به وسیله‌ی ویژگی انتخاب شده توسط این مدل با میزان خطای کمتری می‌توان پیش‌بینی کرد. این الگوریتم قابلیت کار با داده‌های عددی و اسمی را دارد. در برخورد با داده‌های عددی پیوسته، ابتدا داده‌های مورد نظر را به صورت باینری در نظر می‌گیرد و

آزمایشی که برای شبیه‌سازی در نظر گرفته‌ایم، کار کند.

در ادامه با استفاده از تکنیک‌های انتخاب ویژگی، از بین 41 ویژگی، 7 ویژگی را انتخاب کردیم و مجموعه داده‌های مورد استفاده در فاز آموزش و آزمایش را بر اساس این ویژگی‌ها آماده نمودیم. این ویژگی‌ها عبارتست از:

Duration, flag, src\_bytes, dst\_bytes,  
num failed\_login, is\_hot\_login و  
is\_guest\_login

جدول 7: تعداد رکوردهای کلاس حمله و نرمال در مجموعه

	داده‌ها	
	مجموعه داده‌ی آموزشی	مجموعه داده‌ی آزمایشی
تعداد رکوردهای کلاس حمله	25995	25006
تعداد رکوردهای کلاس عادی	81504	6097
تعداد کل نمونه‌ها	107499	31103

همان‌گونه که در جدول 1 مشاهده می‌کنیم، مشخصات مجموعه داده‌های مورد استفاده برای شبیه‌سازی نشان داده شده است.

## 6.2. معیارهای ارزیابی

کارایی سیستم‌های تشخیص نفوذ و حتی در حالت کلی‌تر کارایی دسته‌بندی‌کننده‌ها، غالباً از طریق دو معیار نرخ تشخیص و نرخ خطای مثبت نادرست ارزیابی می‌شود. برای ارائه تعریف این دو معیار ابتدا باید به تعریف دو نوع خطای ممکن در سیستم‌های تشخیص نفوذ بپردازیم. این دو نوع عبارتند از:

- خطای مثبت نادرست: اگر یک اتصال نرمال به عنوان حمله شناسایی شود، یک خطای مثبت نادرست رخ داده است.
  - خطای منفی نادرست: اگر یک اتصال با برچسب حمله به عنوان یک اتصال نرمال شناسایی شود، یک خطای منفی نادرست رخ داده است.
- با توجه به تعاریف فوق داریم:

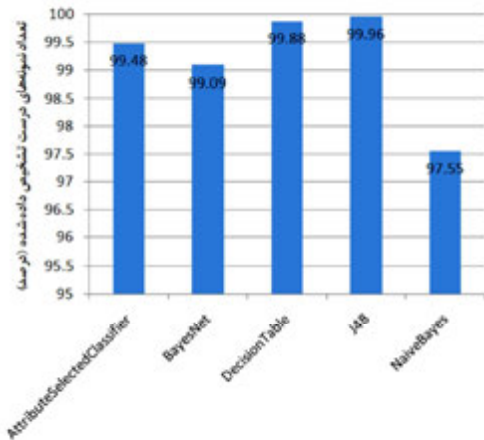
## 6.1. آماده‌سازی مجموعه داده‌های آموزشی و آزمایشی

مجموعه داده‌های فراوانی به منظور ارزیابی سیستم‌های تشخیص نفوذ وجود دارند. مجموعه داده‌های [9KDD Cup 1999] یکی از پرکاربردترین مجموعه داده‌ها برای ارزیابی سیستم‌ها و مدل‌های ارائه به منظور تشخیص نفوذ می‌باشد. مجموعه داده‌های آموزشی KDDCup1999، شامل 4,898,431 رکورد یا بردار اتصالی است که هر کدام از آنها دارای 41 ویژگی هستند.

این ویژگی‌ها به چهار دسته تقسیم می‌شوند: ویژگی‌های پایه<sup>37</sup>، ویژگی‌های محتوایی<sup>38</sup>، ویژگی‌های ترافیک زمانی<sup>39</sup>، و ویژگی‌های ترافیک میزبان<sup>40</sup>. اتصالات به دو نوع مخرب و عادی تقسیم می‌شوند. مجموعه داده‌های آزمایشی نیز شامل 311,027 رکورد می‌باشد و مجموعه داده‌های آموزشی در کل دارای 23 نوع حمله و مجموعه داده‌های تست نیز شامل 37 نوع حمله هستند که 14 نوع حمله بیشتر از داده‌های آموزشی هستند. نکته مهم این است که داده‌های آزمایشی دارای احتمال توزیع یکسان با داده‌های آموزشی نیستند و شامل حملات خاصی هستند که در داده‌های آموزشی موجود نیستند، و این امر کار را بسیار واقعی جلوه می‌دهد.

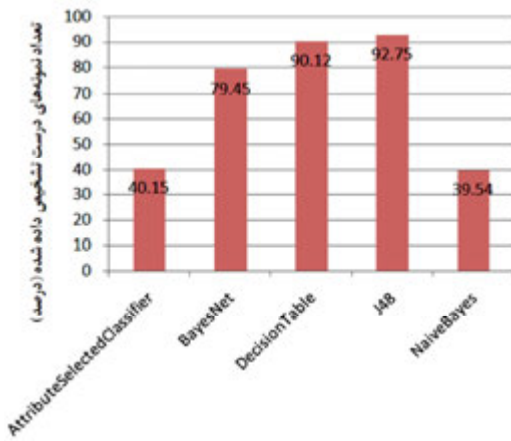
به دلیل تعداد بالای نمونه‌های موجود در این مجموعه داده، نرم‌افزار Weka قادر به نمونه‌برداری از این تعداد رکوردها نمی‌باشد؛ به همین جهت، از نرم‌افزار SQL Server 2010 برای نمونه‌گیری تصادفی استفاده می‌کنیم. چون دارای قدرت پردازشی بالایی می‌باشد. لازم به ذکر است که چون مجموعه داده‌های مورد استفاده در فاز آموزش و آزمایش الگوریتم‌ها تا حدودی حجیم هستند، مقدار حافظه‌ی تخصیص‌یافته به پلتفرم جاوا را، که نرم‌افزار Weka روی آن اجرا می‌شود، از 128 مگا بایت به 2 گیگا بایت افزایش دادیم. با اِعمال این تغییر، نرم‌افزار Weka به راحتی می‌تواند با مجموعه داده‌های آموزشی و





الگوریتم‌های دسته‌بندی

شکل (3): درصد دقت نمونه‌هایی که در فاز آموزش الگوریتم‌ها به درستی تشخیص داده شده‌اند.



الگوریتم‌های دسته‌بندی

شکل 4: درصد تعداد نمونه‌هایی که در فاز آزمایش الگوریتم‌ها به درستی تشخیص داده شده‌اند.

شکل 5 را مشاهده نمایید. نتایج آزمایشات همچنین نشان می‌دهد که در زمان آموزش دادن الگوریتم، الگوریتم C4.5 دارای نرخ هشدارهای اشتباه به مراتب کمتری نسبت به سایر الگوریتم‌ها می‌باشد و حتی می‌توان از این نرخ هشدار اشتباه نیز چشم‌پوشی کرد.

همچنین شکل (6) نشان می‌دهد که الگوریتم درخت تصمیم‌گیری C4.5 باز هم دارای نرخ هشدارهای اشتباه کمی می‌باشد. نرخ هشدارهای اشتباه الگوریتم‌هایی نظیر AttributeSelectedClassifier و NaiveBayes

$$\text{False positive rate} = \frac{\text{Number of false positives}}{\text{Total number of attack connections}} \quad (2)$$

$$\text{False negative rate} = \frac{\text{Number of false negatives}}{\text{Total number of attack connections}} \quad (3)$$

$$\text{Detection rate} = 1 - \frac{\text{Number of false positives}}{\text{Total number of attack connections}} \quad (4)$$

بنابراین نرخ خطای نادرست مثبت (منفی) عبارت است از درصدی از اتصالات نرمال (حمله) که به غلط، برچسب حمله (نرمال) به آنها نسبت داده می‌شود. در آزمایشات انجام شده، از دو معیار نرخ تشخیص و نرخ خطای مثبت نادرست برای ارزیابی راهکار ارائه شده استفاده شده است.

### 6.3. ارزیابی تشخیص نفوذ مبتنی بر درخت تصمیم

همان‌گونه که در شکل (3) مشاهده می‌شود، الگوریتم C4.5 (J48) نسبت به سایر الگوریتم‌های دسته‌بندی، دارای نرخ تشخیص حملات کمی بیشتر در زمان آموزش می‌باشد. این الگوریتم کارایی خود را در زمان آزمایش نشان می‌دهد؛ و همان‌طور که در شکل 4 واضح است، الگوریتم C4.5 بیشترین دقت را در تشخیص حملات از خود نشان می‌دهد.

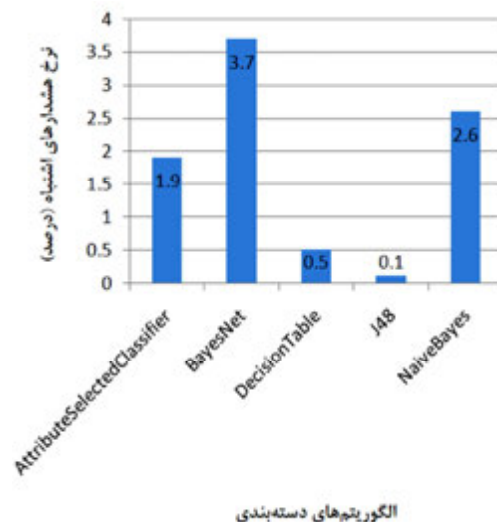
متوقف نمودن حملات اتخاذ می‌کند. مسأله‌ی تشخیص نفوذ به طور وسیعی در زمینه‌ی امنیت شبکه‌های رایانه‌ای مورد مطالعه قرار است و اخیراً در زمینه‌های یادگیری ماشین و داده‌کاوی توجه زیادی به آن شده است.

در این مقاله به طور مختصر در مورد داده‌کاوی و یادگیری ماشین صحبت کردیم و به نحوه‌ی اعمال این تکنیک‌ها به سیستم‌های تشخیص نفوذ اشاره کردیم. همچنین، از طریق شبیه‌سازی یکی از الگوریتم‌های معروف یادگیری ماشین، C4.5، نشان دادیم که این الگوریتم در قیاس با سایر الگوریتم‌های معروف، دارای نرخ کشف نفوذ به مراتب بالاتری است. برای کار آینده، بحث بر روی تکنیک‌های کاهش بُعد داده‌ها و انتخاب مؤثرترین ویژگی‌های مجموعه داده‌ها، می‌تواند در دستور کار پژوهشگران قرار بگیرد. این پژوهش می‌تواند افق روشنی را در باب کاهش زمان تشخیص نفوذ به سیستم‌ها و نیز افزایش جامعیت مجموعه داده‌ها بگشاید.

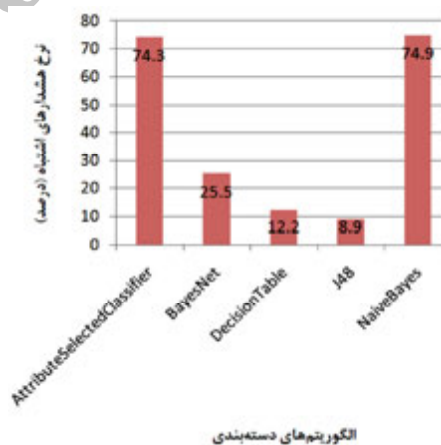
### 8. مراجع

- [1] Lin, Ying and Zhang, Yan and Ou, Yang-jia. "The design and implementation of host-based intrusion detection system". In Proceedings of the 3rd International Symposium on Intelligent Information Technology and Security Informatics, pp. 595–598, 2010, April.
- [2] Sabahi, F., & Movaghar, a.. "Intrusion detection: A survey". In Proceedings of the 3<sup>rd</sup> International Conference on Systems and Networks Communications pp. 23–26, New York, NY, USA, 2008, October.
- [3] Chang-tien Lu, Arnold P Boedihardjo, PrajwalManalwar and Falls Church. "Exploiting Efficient Data Mining Techniques to Enhance Intrusion Detection Systems". IRI 2005 IEEE International Conference on Information Reuse and Integration Conf 2005.

می‌تواند ناشی از یادگیری بیش از حد<sup>41</sup> در فاز آموزششان باشد.



شکل (5): نرخ هشدارهای اشتباه در فاز آموزش الگوریتم‌ها



شکل (6): نرخ هشدارهای اشتباه در فاز آزمایش الگوریتم‌ها

### 7. نتیجه‌گیری

همان‌گونه که در این مقاله عنوان شد، سیستم تشخیص نفوذ، با استفاده از یک‌سری قوانین از پیش تعریف‌شده دسترسی کاربران را بررسی و محدود می‌نماید. این قوانین، بر پایه‌ی دانش افراد مشخص شده است. این سیستم، تمام نفوذهای مهاجمان را مورد شناسایی قرار می‌دهد و فعالیت‌های لازم را برای

- <sup>34</sup>Pre-pruning
- <sup>35</sup>Post-pruning
- <sup>36</sup>Information gain
- <sup>37</sup>Basic features
- <sup>38</sup>Content features
- <sup>39</sup>Time-based traffic features
- <sup>40</sup>Host-based traffic features
- <sup>41</sup>Overfitting

- [4] [http://www.kmining.com/info\\_definitions.html](http://www.kmining.com/info_definitions.html)
- [5] Ming Xue, Changjun Zhu. "A Study and Application on Machine Learning of Artificial Intelligence". International Joint Conference on Artificial Intelligence, pp. 272-274, 2009.
- [6] S.Y. Wua and , E. Yen b, "Data mining-based intrusion detectors", Expert Systems with Applications: An International Journal Volume 36 Issue 3, pp. 5605-5612, 2009, April.
- [7] <http://www.cs.waikato.ac.nz/ml/weka/>
- [8] M. A. Hall, "Correlation-based Feature Subset Selection for Machine Learning. Hamilton", New Zealand, 1998.
- [9] <http://kdd.ics.uci.edu/databases/kddcup99/>

زیرنویس

- <sup>1</sup>Intrusion Detection Systems(IDS)
- <sup>2</sup>Audit data
- <sup>3</sup>Log
- <sup>4</sup>Network based
- <sup>5</sup>Promiscuous
- <sup>6</sup>Header
- <sup>7</sup>Host based
- <sup>8</sup>Administrator
- <sup>9</sup>Misuse detection
- <sup>10</sup>Decision tree
- <sup>11</sup>Genetic algorithm
- <sup>12</sup>Fuzzy logic
- <sup>13</sup>Clustering
- <sup>14</sup>Bayesian network
- <sup>15</sup>Data cleaning
- <sup>16</sup>Data integration
- <sup>17</sup>Data transformation
- <sup>18</sup>Data mining
- <sup>19</sup>Pattern evaluation
- <sup>20</sup>Knowledge discovery
- <sup>21</sup>Bayesian decision theory
- <sup>22</sup>Multi-variate methods
- <sup>23</sup>Classification trees
- <sup>24</sup>Linear discrimination
- <sup>25</sup>Multilayer perception
- <sup>26</sup>Local models
- <sup>27</sup>Hidden Markov models
- <sup>28</sup>Reinforcement learning
- <sup>29</sup>Flow chart
- <sup>30</sup>Node splitting
- <sup>31</sup>Entropy function
- <sup>32</sup>Information gain
- <sup>33</sup>Gini function