

Using Genetic Algorithm to Improve Random Forest Algorithm in Order to Detect DDoS Attacks in Cloud Computing Platform

¹Ali Mahmodi Derakhsh, ²Parisa Daneshjoo, ³Changiz Delara

¹ Department of Computer Engineering, West Tehran Branch, Islamic Azad University, Tehran, Iran
Mahmodi.Ali@wtiau.ac.ir

² Department of Computer Engineering, West Tehran Branch, Islamic Azad University, Tehran, Iran
Daneshjoo.p@wtiau.ac.ir

³ Department of Computer Engineering, West Tehran Branch, Islamic Azad University, Tehran, Iran
Delara.c@wtiau.ac.ir

Abstract

Devices such as routers, switches or firewalls are the most vital connections in communication network among physical machines in a cloud computing environment. In the absence of security on the network, intruders are allowed to access the equipment and configure it in the way they want to. Hence, a method suggested to deal with denial-of-service (DoS) attacks in the cloud computing platform is one of the essential and most important security issues in this area. This study tends to provide a smart method based on random forest algorithm focusing on genetic algorithms for detecting DoS attacks. Through different network streams, network streams which trigger DoS and DDoS attacks are very important.

The main idea of this study is to use random forest algorithm to identify DoS attacks, which is the main reason for optimizing this algorithm using genetic algorithms. In this method, an optimal subset of the set of features is extracted using genetic algorithms, and this optimal subset is used for random forest learning. Results of the experiments carried out and comparison of the suggested method with other methods indicate proper accuracy and operation of the suggested method.

Keywords: Cloud computing, Network security, Denial-of-service attacks, Genetic algorithms, Random forest

بکارگیری الگوریتم ژنتیک برای بهبود الگوریتم جنگل تصادفی به منظور تشخیص حملات محروم سازی از سرویس توزیع شده در بستر رایانش ابری

علی محمودی درخش^۱، پریسا دانشجو^۲، چنگیز دل آرا^۳

^۱دانشجوی کارشناسی ارشد، گروه کامپیوتر، واحد تهران غرب، دانشگاه آزاد اسلامی، تهران
Mahmodi.Ali@wtiau.ac.ir

^۲استادیار و عضو هیئت علمی، گروه کامپیوتر، واحد تهران غرب، دانشگاه آزاد اسلامی، تهران
Daneshjoo.p@wtiau.ac.ir

^۳استادیار و عضو هیئت علمی، گروه کامپیوتر، واحد تهران غرب، دانشگاه آزاد اسلامی، تهران
Delara.c@wtiau.ac.ir

چکیده

تجهیزاتی همچون مسیریاب، سوئیچ یا دیوارهای آتش به عنوان حیاتی ترین اتصالات در شبکه ارتباطی میان ماشین های فیزیکی در یک محیط رایانش ابری هستند. در صورت عدم وجود امنیت در شبکه، به نفوذگران به شبکه اجازه داده می شود که با دستیابی به تجهیزات، امکان پیکربندی آن ها را به گونه ای که تمایل دارند عمل کنند. از این رو ارائه روشی برای مقابله با حملات محروم سازی از سرویس^۱ در بستر رایانش ابری یکی از ضروریات و مهم ترین مسائل امنیتی این حوزه می باشد. در این تحقیق سعی شده روشی هوشمند و مبتنی بر الگوریتم درخت تصادفی^۲ با تمرکز بر روی الگوریتم ژنتیک^۳ برای شناسایی حملات محروم سازی از سرویس ارائه شود. از میان جریان های شبکه ای متفاوت، جریان های شبکه ای که موجب بروز حملات DoS و DDos^۴ می شوند از اهمیت بسیار بالایی برخوردارند.

ایده اصلی در این تحقیق استفاده از روش های داده کاوی و الگوریتم درخت تصادفی برای شناسایی حملات محروم سازی از سرویس است که نکته اصلی در بهینه سازی این الگوریتم با استفاده از الگوریتم های ژنتیک می باشد. طی این روش زیرمجموعه بهینه از مجموعه ویژگی های با استفاده از الگوریتم های ژنتیک استخراج می شود و این زیرمجموعه بهینه برای یادگیری مدل درخت تصادفی مورد استفاده قرار می گیرد. نتایج آزمایش های انجام شده و مقایسه آن با روش معمول حاکی از دقت و عملکرد مناسب روش ارائه شده دارد.

کلمات کلیدی

رایانش ابری، امنیت شبکه، حملات محروم سازی از سرویس، الگوریتم های ژنتیک، جنگل تصادفی

دستیابی به تجهیزات، امکان پیکربندی آن ها را به گونه ای که تمایل دارند عمل کنند، از این طریق هرگونه نفوذ و سرقت اطلاعات و یا هر نوع صدمه دیگری به شبکه محیط رایانش ابری، توسط نفوذگر، امکان پذیر خواهد شد. برای جلوگیری از خطرهای محروم سازی از سرویس، تأمین امنیت تجهیزات بر روی شبکه الزامی است. توسط این حمله ها نفوذگران می توانند سرویس هایی را در شبکه از کار بیاندازند. از این رو ارائه روشی برای مقابله با حملات محروم سازی از سرویس در بستر شبکه محیط رایانش ابری یکی از ضروریات و مهم ترین مسائل امنیتی این حوزه می باشد. از این رو در این تحقیق سعی بر آن داریم تا روش هوشمند و مبتنی بر الگوریتم درخت تصمیم با

۱- مقدمه

برای تأمین امنیت بر روی یک شبکه، یکی از بحرانی ترین و خطرناک ترین مراحل، تأمین امنیت دسترسی و کنترل تجهیزات شبکه است. تجهیزات همچون مسیریاب، سوئیچ یا دیوارهای آتش که به عنوان حیاتی ترین اتصالات در شبکه ارتباطی میان ماشین های فیزیکی در یک محیط رایانش ابری هستند، موضوع امنیت تجهیزات به دو علت اهمیت ویژه ای می یابد [1] عدم وجود امنیت تجهیزات در شبکه، به نفوذگران به شبکه اجازه می دهد که با

در [18] روشی هوشمند و مبتنی بر الگوریتم دسته‌بندی‌کننده بیز ساده^۶ با تمرکز بر روی الگوریتم ژنتیک برای شناسایی حملات محروم‌سازی از سرویس ارائه شده که شباهت بسیار زیادی از لحاظ نوآوری و نحوه پیاده سازی به روش ارائه شده در این مقاله دارد و فقط از لحاظ الگوریتم یادگیری ماشین به کار رفته در روش ارائه شده، متفاوت است.

۲- روش ارائه شده

همان‌طور که در بخش‌های قبل نیز به آن اشاره شد، روش ارائه شده در این تحقیق مبتنی بر یادگیری ماشین و بهینه‌سازی الگوریتم‌های ژنتیک است. در این تحقیق از طبقه‌بند جنگل تصادفی استفاده می‌شود. هدف این است که طبقه‌بند فوق توسط الگوریتم‌های ژنتیک تا حد ممکن بهینه شود و در انتها، دقت شناسایی حملات DDOS توسط این الگوریتم‌ها با یکدیگر مقایسه شده و بهترین الگوریتم انتخاب شود.

۲-۱- روند کلی الگوریتم‌های ژنتیکی

در شکل ۱ یک الگوریتم ژنتیکی استاندارد و در شکل ۲ نمودار گردش الگوریتم‌های ژنتیکی نشان داده شده است. قبل از این که یک الگوریتم ژنتیکی بتواند اجرا شود، ابتدا باید کدگذاری یا نمایش مناسبی برای مسئله مورد نظر پیدا شود. همچنین یک تابع برازندگی نیز باید ابداع شود تا به هر راه‌حل کدگذاری شده ارزشی را نسبت دهد.

در طی اجراء والدین برای تولیدمثل انتخاب می‌شوند و با استفاده از عملگرهای آمیزش و جهش باهم ترکیب می‌شوند تا فرزندان جدیدی تولید کنند. این فرآیند چندین بار تکرار می‌شود تا نسل بعدی جمعیت تولید شود. سپس این جمعیت بررسی می‌شود و در صورتی که ضوابط همگرایی برآورده شوند، فرآیند فوق خاتمه می‌یابد.

```

BEGIN /* genetic algorithm */
Generate initial population
Compute fitness of each individual

WHILE NOT finished DO
BEGIN /* produce new generation */

FOR population_size / 2 DO
BEGIN /* reproductive cycle */
Select two individuals from old generation for mating
/* biased in favor of the fitter ones */
Recombine the two individuals to give two offspring
Compute fitness of the two offspring
Insert offspring in new generation
END
END

```

شکل ۱: الگوریتم استاندارد ژنتیک [9]

۲-۱-۱- روش انتخاب تورنمنت

ایده اصلی در روش‌های انتخاب این است که افراد بهتر بر افراد بدتر ترجیح داده شوند، که بهتر و بدتر بودن افراد توسط تابع برازندگی f تعریف می‌شود. روش‌های انتخاب متعددی برای استفاده در الگوریتم‌های ژنتیکی پیشنهاد شده‌اند. یکی از ویژگی‌های خوب روش‌های انتخاب این است که این روش‌ها مستقل از نمایش افراد جمعیت هستند و در آن‌ها تنها مقادیر برازندگی افراد در نظر گرفته می‌شود.

تمرکز بر روی الگوریتم‌های ژنتیک برای شناسایی حملات محروم‌سازی از سرویس ارائه دهیم. ایده اصلی در این تحقیق استفاده از الگوریتم درخت تصمیم برای شناسایی حملات محروم‌سازی از سرویس است که نکته اصلی در بهینه‌سازی این الگوریتم‌ها با استفاده از الگوریتم‌های ژنتیک می‌باشد که روشی کاملاً جدید و نو می‌باشد.

۱-۱- کارهای پیشین

در [2] به بررسی روش‌های مختلف یادگیری ماشین و داده‌کاوی که در شناسایی حملات محروم‌سازی از سرویس از آن‌ها استفاده شده است، می‌پردازد. روش‌هایی مانند ماشین بردار پشتیبان، شبکه‌های بیزین و درخت تصمیم از جمله مهم‌ترین الگوریتم‌هایی است که در این مقاله مروری مورد ارزیابی قرار گرفته است. این مقاله برای آشنایی با روش‌های مقابله با حملات محروم‌سازی از سرویس که از فن‌های داده‌کاوی استفاده می‌کنند، منبعی جامع، کامل و جدید هست.

در [3] به شناسایی حملات محروم‌سازی از سرویس با استفاده از یک طبقه‌بند پرداخته شده است. طبقه‌بند مورد استفاده در این روش ماشین بردار پشتیبان است که در نوع خود یکی از کاراترین طبقه‌بندهای ممکن هست. نکته اصلی در این مقاله استفاده از این طبقه‌بند برای شناسایی حملات محروم‌سازی از سرویس در محیط شبکه مجهز به SDN^۵ می‌باشد.

در [4] به ارائه چارچوب و چارچوبی برای شناسایی حملات محروم‌سازی از سرویس با استفاده از الگوریتم‌های ژنتیک پرداخته است. با توجه به الگوریتم‌های ژنتیک برای حل مسائل بهینه‌سازی مورد استفاده قرار می‌گیرند، در این مقاله شناسایی حملات محروم‌سازی از سرویس به یک مسئله بهینه‌سازی نگاشت شده است و از الگوریتم‌های ژنتیک برای حل آن استفاده می‌شود. آزمایش‌ها نشان می‌دهد که روش ارائه شده از دقت مناسبی برخوردار است.

در [5] از روش انتخاب ویژگی برای یافتن ویژگی‌های بهینه برای تشخیص نفوذ و شناسایی حملات DDOS استفاده شده است. ویژگی اصلی این روش در استفاده از فیلترها برای انتخاب مجموعه ویژگی بهینه است. در نهایت نیز از الگوریتم ماشین بردار پشتیبان برای تشخیص نفوذ استفاده شده است. نتایج حاکی از دقت بسیار مناسب این روش می‌باشد.

در [6] از فیلترها برای کاهش ویژگی و شناسایی مجموعه ویژگی بهینه استفاده شده است. تفاوت این روش با پاراگراف قبل در این است که در این روش اولاً از فیلترهای چندگانه استفاده شده است و دوماً روش ارائه شده یک روش بلادرنگ است و قابلیت اجرا بر روی داده‌های جریانی را دارد. این روش تمرکز خود را بر روی تشخیص نفوذ در محیط رایانش ابری قرار داده است.

در [7] نیز تمرکز اصلی بر روی انتخاب و شناسایی ویژگی است. برای این منظور از روش آنتروپی فازی استفاده شده است. علاوه بر این برای بهینه‌سازی و انتخاب ویژگی‌های بهینه نیز از روش کلونی مورچگان استفاده شده است. ویژگی مهم دیگری که این روش دارد بلادرنگ بودن آن است که تمایز آن را با سایر روش‌ها ایجاد می‌کند. در مجموع این روش دارای نوآوری‌های مناسب و نتایج بسیار قابل قبولی می‌باشد.

۲-۲- الگوریتم جنگل تصادفی

این روش یک فن یادگیری جمعی^۹ برای طبقه‌بندی است که مجموعه‌ای از درختان تصمیم‌گیری^{۱۰} را در زمان آموزش می‌سازد و خروجی نهایی، مد و یا میانگین خروجی این درختان می‌باشد.

جنگل تصادفی چالش overfitting در درختان تصمیم معمول را مرتفع می‌سازد. به عبارت دیگر مجموعه داده اولیه به چندین قسمت تقسیم می‌شود و هر قسمت بر روی یکی از درختان آموزش داده می‌شود و خروجی نهایی گرفته می‌شود. اکنون با توجه به سیاست‌های مختلف میانگین یا مد خروجی-های درخت‌ها به عنوان خروجی نهایی در نظر گرفته می‌شود [8].

هدفی که در این تحقیق دنبال می‌شود، استفاده از الگوریتم‌های ژنتیک برای بهبود طبقه‌بند است. هدف اصلی در استفاده از الگوریتم‌های ژنتیک، انتخاب ویژگی است. با توجه به این که مجموعه داده حملات DDoS، عموماً دارای تعداد بسیار زیادی ویژگی هستند، استفاده از روشی مناسب برای انتخاب ویژگی نیز تأثیر بسزایی بر کارایی طبقه‌بند جنگل تصادفی دارد. روال کلی استفاده از الگوریتم‌های ژنتیک برای انتخاب ویژگی از مجموعه ویژگی‌ها به صورت زیر است:

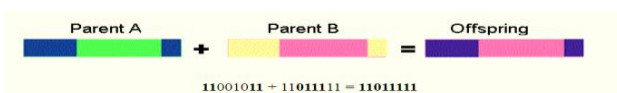
جمعیت اولیه به صورت تصادفی تولید می‌شوند. هر نمونه از جمعیت شامل n ژن است که برابر تعداد ویژگی‌ها در مجموعه داده است. به عبارت دیگر هر ژن مشخص می‌کند که آیا ویژگی متناسب با آن در ساخت مدل استفاده شده است یا خیر، اگر استفاده شده باشد مقدار آن ۱ و در غیر این صورت مقدار آن صفر خواهد بود. در نتیجه هر نمونه در جمعیت نشان‌دهنده انتخاب برای ویژگی‌های موجود است. برای هر نمونه در جمعیت جاری، مدل مربوطه ایجاد می‌شود.

بعد از این که مدل جنگل تصادفی مربوطه ایجاد شد، این مدل با مجموعه داده ارزیابی^{۱۱}، ارزیابی می‌شود و میزان classification error شامل آن استخراج می‌شود. جنگل تصادفی ای که مقدار کمتری classification error rate داشته باشد نمونه مناسب‌تر و مطلوب‌تری است.

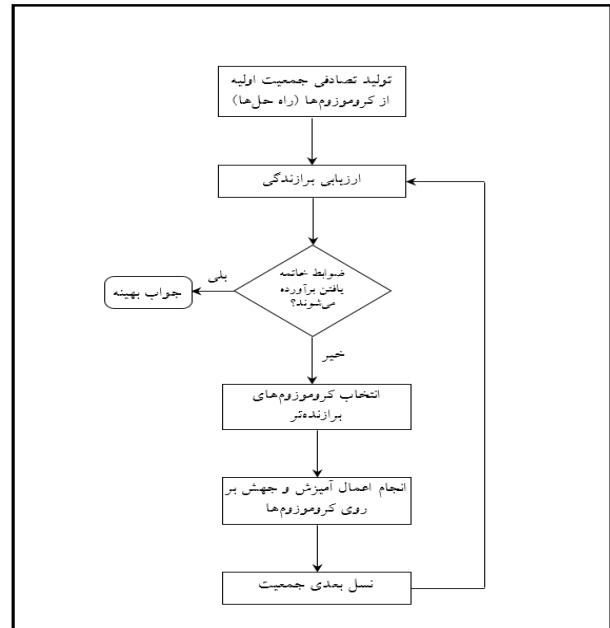
زمانی که مقدار تابع ارزیابی یا همان classification error rate برای تمام نمونه‌های جمعیت محاسبه شد، الگوریتم ژنتیک نسل بعدی را به صورت زیر می‌سازد:

(۱) انتخاب نمونه‌ها برای ساخت نسل بعدی با استفاده از Rank selection method که در ادامه توضیح داده می‌شود.

(۲) از روش two point crossover برای ساخت فرزندان استفاده می‌شود (شکل ۳)، این روش به صورت زیر محاسبه می‌شود:
در این روش در کروموزوم والدین دو نقطه انتخاب می‌شود و هر آنچه میان این دو نقطه قرار دارد جابجا می‌شود تا فرزندان جدید تولید شوند.



شکل ۳: روش Two-point crossover [11]



شکل ۲: نمودار گردش الگوریتم‌های ژنتیک [9]

انتخاب تورنمنت^۷ مشابه با انتخاب رتبه‌ای برحسب فشار انتخاب است، اما از نظر محاسباتی کارتر و برای پیاده‌سازی‌های موازی مناسب‌تر است. در این روش، دو فرد از جمعیت به صورت تصادفی انتخاب می‌شوند. سپس، یک غلذ تصادفی r بین صفر و یک انتخاب می‌شود. k یک پارامتر است و اگر $r < k$ باشد، برای مثال اگر $0.75/k$ باشد، فرد برآورده‌تر و در غیر این صورت فردی که برآوردگی کمتری دارد، به عنوان والد انتخاب می‌شود. این دو سپس به جمعیت اولیه بازگردانده می‌شوند و دوباره در فرآیند انتخاب شرکت داده می‌شوند.

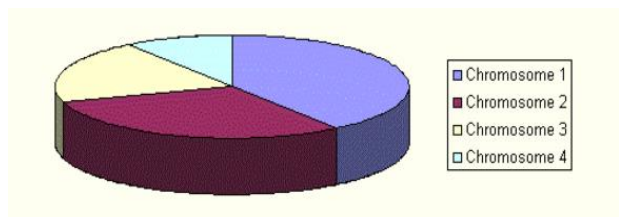
۲-۱- شرط پایان الگوریتم

در الگوریتم‌های تکاملی غالباً اجرای برنامه برای تعداد نسل‌های از پیش تعیین شده‌ای صورت می‌گیرد؛ اما شرط دیگری نیز برای پایان الگوریتم‌های ژنتیک توسط Grefenstette [10] ارائه شده است که آن میزان پراکندگی بیت‌ها درون جمعیت^۸ است. این محک نشان‌دهنده میزان همگرا شدن کد اعضای جمعیت است. اگر کد یک عنصر دارای طول یک بیت باشد و به صورت $\bar{a}_i (a_{i1}, \dots, a_{i\mu})$ نشان داده شود و μ تعداد اعضای جمعیت است که بین $\{1, \dots, \mu\}$ می‌باشد، مقدار پراکندگی بیت جمعیت P ، $b(P)$ به صورت زیر تعریف می‌شود:

$$b(p) = \frac{1}{\mu \cdot 1} \sum_{j=1}^1 \text{Max} \left\{ \sum_{i=1}^{\mu} (1 - a_{ij}), \sum_{i=1}^{\mu} a_{ij} \right\} \in [0.5, 1.0] \quad (1)$$

هراندازه میزان b بزرگ‌تر باشد میزان پراکندگی بیت‌ها درون جمعیت کمتر خواهد بود. در حالت ویژه اگر $b(P)=1$ باشد به این معنی است که کد همه اعضای جمعیت یکسان است. شرط پایان به صورت $b(p) > b_{max}$ تعریف می‌شود که معمولاً $b_{ma} \approx 0.95$ است.

مشکل اصلی این روش در همگرایی کند آن است، زیرا میان بهترین کروموزوم و دیگر کروموزومها هیچ تفاوتی قائل نشده است.



شکل ۵: چرخ رولت کروموزومها برای روش Rank Selection

همانطور که ملاحظه می شود هرکدام از روشها فوق به نوعی دارای کمبود و یا نقصی در عملکرد خود هستند. در این پژوهش از روش ارائه شده در [12] استفاده می شود که طی آزمایشها و تحلیلهایی که بر روی این روش انجام شده است دارای کارایی مناسبی است.

مراحل این الگوریتم به صورت زیر است:

- ۱) جمعیت کروموزومها را با توجه به مقدار تابع fitness آنها sort می کنیم.
- ۲) جمعیت را به دو قسمت تقسیم می کنیم. قسمت با مقادیر fitness بالاتر (HF) و قسمت با مقادیر fitness پایین تر (LF) که محل شکستن جمعیت نیز توسط پارامتر m که به صورت درصدی بیان می شود تعیین می گردد.
- ۳) همیشه والد اول (P1) را از HF و والد دوم (P2) را از LF انتخاب می شود.
- ۴) فرزندان P1 و P2 محاسبه می شوند.
- ۵) فرزندان به انتهای جمعیت اضافه می شوند.
- ۶) جمعیت مجدد بر اساس مقدار تابع fitness مرتب می شود و کروموزومهای انتهایی حذف می شوند تا طول جمعیت ثابت بماند.

۲-۳- تحلیل روش ارائه شده

از الگوریتمهای ژنتیک برای شناسایی و انتخاب زیرمجموعه مؤثر از ویژگیها برای طبقه بندی استفاده می شود. به عبارت دیگر سعی شده است تا ویژگیهایی که بیشترین تأثیر را در جداسازی دادهها دارند از میان ویژگیهای اولیه انتخاب و سپس طبقه بندی بر اساس این ویژگیها انجام شود. برای نشان دادن تأثیر ویژگیها از الگوریتم Support Vector Decomposition [13] استفاده شده است تا برای هر ویژگی در مجموعه داده، مقدار ویژه آن محاسبه شود. هرچه مقدار ویژه بیشتر باشد، نشان از تأثیر بیشتر ویژگی مربوطه در جداسازی و یا تعریف دادهها دارد. در این تحقیق از مجموعه داده NSL-KDD 2013 [14] استفاده شده است که شامل ۴۱ ویژگی است و توسعه یافته مجموعه داده KDDCup99 [15] می باشد.

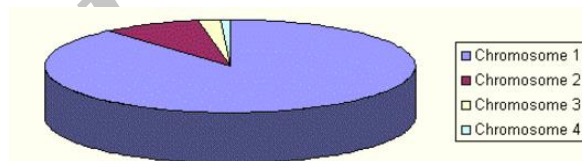
شایان ذکر است که در فرآیند پیش پردازش، ویژگیهایی که دارای مقادیر عددی نیستند، خود به چندین ویژگی عددی شکسته می شوند تا بتوان از این ویژگیها نیز در الگوریتم جنگل تصادفی استفاده نمود، در نتیجه پس از پیش پردازش تعداد ویژگیها بسیار بیشتر از ۴۱ ویژگی اولیه می شود. شکل ۶ مقادیر ویژه برای ویژگیهای این مجموعه داده را به صورت نزولی نشان می

۳) در هر فرزند تولید شده از روش bit level mutation برای چشم استفاده می کنیم. در این روش یک بیت به صورت تصادفی انتخاب می شود و مقدار آن جابجا می شود یعنی اگر صفر باشد، یک می شود و اگر یک باشد، صفر می شود. احتمال انتخاب هر بیت برای چشم در این روش برابر $\frac{1}{l}$ است که l طول کروموزوم است. دو نمونه برتر نگه داشته می شوند و تمام جمعیت جاری با offspringها جایگزین می شوند.

مراحل فوق به صورت مکرر و تا زمان رسیدن به حداکثر تعداد تکرارها که ۱۰۰۰ دفعه در نظر گرفته شده، اجرا می شوند.

نکته اصلی که در الگوریتمهای ژنتیک وجود دارد و می تواند تأثیر بسزایی در عملکرد آن داشته باشد و مانع از همگرایی سریع^{۱۲} شود، روش انتخاب کروموزومهای برتر است که در ادامه چندین روش انتخاب معمول که در الگوریتمهای ژنتیک از آنها استفاده می شود، توضیح داده می شود.

فرض کنید برای تمام کروموزومها، چرخ رولت^{۱۳} شکل ۴، را به نحوی می سازیم، که هر کروموزومی که مقدار تابع fitness بیشتری دارد مساحت بیشتری را به خود اختصاص دهد.



شکل ۴: چرخ رولت کروموزومها [11]

در روش معمول چرخ رولت، کروموزوم با بیشترین مقدار تابع fitness شانس بیشتری برای انتخاب شدن دارد که نحوه انتخاب آن به صورت زیر است:

- ۱) مجموع تمام مقادیر تابع fitness برای هر کروموزوم محاسبه می شود و در متغیر S قرار داده می شود.
- ۲) یک مقدار تصادفی r که بین $(0, S)$ قرار دارد انتخاب می شود.

۳) بر روی جمعیت کروموزومها و مقادیر تابع fitness از صفر تا S گذر می شود تا زمانی که مقدار جاری از r بزرگتر باشد، در این حالت عمل متوقف می شود و کروموزوم مربوطه برگشت داده می شود.

بدیهی است که در این حالت، کروموزومها با مقادیر fitness بیشتر شانس بالاتری برای انتخاب دارند؛ اما اگر یک کروموزوم دارای مقدار fitness بسیار زیادی، مثلاً ۹۰ درصد مجموع کل مقادیر تابع fitness باشد، در این حالت شانس بسیار کمی برای انتخاب دیگر کروموزومها به وجود می آید. در روش Rank Selection برای حل این مشکل، به هر کروموزوم یک مقدار rank اختصاص می یابد که ۱ برای بدترین کروموزوم، ۲ برای دومین بدترین و در نهایت N برای بهترین کروموزوم است. در این حالت، تمام کروموزومها شانس یکسانی برای انتخاب شدن دارند (شکل ۵)، اما

۳- نتایج

۳-۱- مجموعه داده

مجموعه داده‌ای که در این تحقیق برای پیاده‌سازی روش ارائه شده از آن استفاده می‌شود، مجموعه داده NSL-KDD [۱۴] است. این مجموعه داده شامل سه فایل اصلی زیر است:

- ۱) مجموعه داده آموزش که شامل ۱۲۵۹۷۳ داده است که هر داده نیز شامل ۴۱ ویژگی می‌باشد.
- ۲) مجموعه داده آزمون اول که شامل ۲۲۵۴۴ داده است.
- ۳) مجموعه داده آزمون دوم که شامل ۱۱۸۵۰ داده است.

باتوجه به این که روش ارائه شده در این تحقیق برای شناسایی حملات محروم سازی از سرویس در بستر رایانش ابری می‌باشد، برای ارزیابی، الگوریتم را به طور مستقل و فارغ از نیازمندی‌های محیط رایانش ابری با روش‌های معمول مقایسه کردیم تا میزان دقت، کارایی و عملکرد روش پیشنهادی در مقایسه با روش‌ها پایه مشخص گردد. به عبارت دیگر از لحاظ کردن شرط و یا فرضی که به کلیات مسئله حملات محروم سازی از سرویس آسیبی بزند، اجتناب شده است.

۳-۲- پیش پردازش مجموعه داده

مجموعه داده فوق برای آماده شدن برای اجرای الگوریتم درخت تصمیم و الگوریتم‌های ژنتیک، نیاز به فرآیند پیش‌پردازش دارد که در قالب فازهای زیر انجام می‌شود:

- ۱) حذف مقادیر تهی
برای حذف مقادیر تهی از جایگزین کردن این مقادیر با میانگین مقادیر همان ویژگی در مجموعه داده استفاده شده است.

۲) نرمال سازی ویژگی‌ها

با توجه به این که دامنه مقادیر ویژگی‌ها متفاوت است و اختلاف زیادی با یکدیگر دارد، می‌بایست یکسان سازی دامنه این مقادیر برای بهبود عملکرد الگوریتم‌های یادگیری ماشین صورت پذیرد که این امر نیز بر روی مجموعه داده اعمال شده است. برای نرمال سازی ویژگی‌ها از قالب نرمال سازی Z-score [۱۷] استفاده شده است.

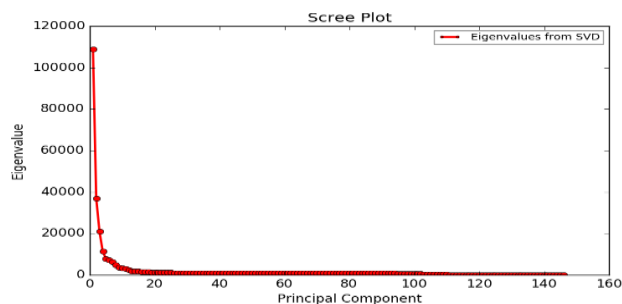
۳-۳- پیاده سازی روش ارائه شده

برای پیاده سازی روش ارائه شده از زبان برنامه نویسی پایتون نسخه ۳.۵ استفاده شده است. مشخصات سخت‌افزاری که الگوریتم بر روی آن اجرا شده است در جدول ۱ آمده است. برای پیاده سازی روش از کتابخانه‌های مربوط به الگوریتم‌های یادگیری ماشین، الگوریتم‌های ژنتیک و محاسبات ریاضی ماتریسی که به زبان پایتون هستند، استفاده شده است. لیست این کتابخانه‌ها به همراه توضیحات عملکرد آن‌ها در جدول ۲ آورده شده است.

۳-۴- پارامترهای ارزیابی

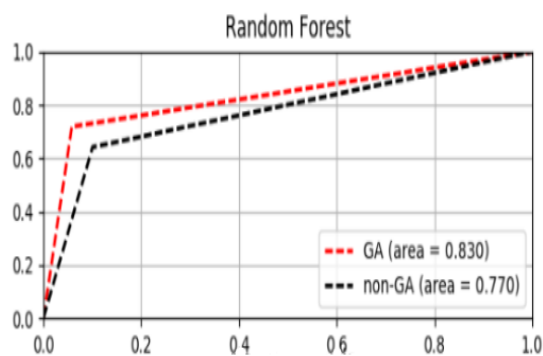
برای ارزیابی روش ارائه شده از مجموعه پارامترهای ماتریس درهم‌ریختگی^{۱۷} استفاده شده است.

دهد. همان‌طور که در شکل نیز مشخص است، حدود ۲۰ ویژگی از حدود ۱۵۰ ویژگی، پس از انجام پیش‌پردازش‌های لازم در این مجموعه داده دارای مقدار ویژه بزرگ‌تر از صفر هستند و این نشان از عدم تأثیر حجم زیادی از ویژگی‌ها، حدود ۱۳۰ ویژگی در تعریف و یا توزیع داده‌ها است، بنابراین در این تحقیق از فن الگوریتم ژنتیک برای شناسایی ویژگی‌های با تأثیر بالا در عملکرد طبقه‌بند استفاده شد.



شکل ۶: مقادیر ویژه برای ویژگی‌های مجموعه داده NSL-KDD

در ادامه تأثیر الگوریتم‌های ژنتیک بر نتیجه طبقه‌بندی در نظر گرفته شده بررسی می‌شود. در شکل ۷ برای طبقه‌بند جنگل تصادفی، تفاوت کارایی طبقه‌بند با و بدون اعمال الگوریتم‌های ژنتیک در حالت اولیه مورد بررسی قرار گرفته است. برای ارزیابی عملکرد طبقه‌بند از پارامتر ROC [۱۶] استفاده شود که نموداری برحسب نرخ خطای مثبت^{۱۵} و نرخ مثبت واقعی^{۱۶} است. هرچه مساحت سطح زیر نمودار ROC بیشتر باشد نشان‌دهنده این است که به عبارتی مربوطه با دقت و احتمال بیشتری جواب صحیح تولید می‌کند و دیگر، هرچه سطح زیر نمودار بیشتر باشد، طبقه‌بند مربوطه کارایی بهتری دارد. همان‌طور که در شکل نیز مشخص است، در طبقه‌بند درخت تصمیم در نظر گرفته شده، در صورت اعمال الگوریتم‌های ژنتیک شاهد عملکرد بهتری هستیم. با توجه به توضیحات ارائه شده، کارایی این روش در بهبود دقت طبقه‌بند درخت تصمیم واضح و روشن می‌شود.



شکل ۷: مقایسه کارایی طبقه‌بند جنگل تصادفی با و بدون اعمال الگوریتم‌های ژنتیک در حالت اولیه

۴-۱- پیشنهادات

با توجه به رشد روزافزون فناوری و حملات سایبری، لزوم استفاده از روش‌های مبتنی بر تحلیل داده و آمار اطلاعات می‌تواند نقش مهمی را در کاهش آسیب‌های ناشی از حملات سایبری ایفا کند. از این رو پیشنهادهای زیر به منظور ادامه پژوهش جاری ارائه می‌گردد:

(۱) مقیاس‌پذیر بودن

با توجه به بزرگ شدن شبکه‌ها و افزایش اطلاعات، نیاز ضروری است تا روش‌های ارائه شده قابلیت اعمال و پشتیبانی حجم زیادی از داده‌ها را داشته باشند.

(۲) تمرکز بر روی پیش‌پردازش داده

از آنجایی که کیفیت داده‌ها، چگونگی ذخیره‌سازی و قالب نمایش آن‌ها در منابع و حوزه‌های مختلف متفاوت است، ادغام و ترکیب آن‌ها برای اجرای الگوریتم‌های یادگیری ماشین دارای چالش‌هایی خواهد بود که می‌تواند زمینه تحقیقات آینده باشد. با توجه به این که داده‌های شبکه معمولاً دارای قالب‌ها ناهمگون هستند، این امر می‌تواند کمک شایانی به افزایش دقت شناسایی حملات سایبری نماید.

جدول ۳: مقایسه نسخه‌های مختلف با مجموعه داده آزمون اول

KDDTest ⁺				نام روش
Accuracy	F-Measure	Recall	Precision	
۰.۷۳۸۶	۰.۷۴	۰.۷۴	۰.۷۸	Random Forest
۰.۷۸۶۹	۰.۷۹	۰.۷۹	۰.۸۳	Random Forest+ GA

جدول ۴: مقایسه نسخه‌های مختلف با مجموعه داده آزمون دوم

KDDTest ⁻²¹				نام روش
Accuracy	F-Measure	Recall	Precision	
۰.۵۳۵۱	۰.۵۹	۰.۵۴	۰.۷۶	Random Forest
۰.۵۴۴۲	۰.۵۹	۰.۵۴	۰.۸۰	Random Forest + GA

جدول ۵: میزان زمان اجرای الگوریتم‌ها به دقیقه

Run Time (min)	نام روش
۱	Random Forest
۴۵۶	Random Forest + GA

جدول ۱: مشخصات سخت افزار پیاده سازی

سیستم‌عامل	ویندوز ۷، نسخه ۶۴ بیتی
زبان برنامه‌نویسی	پایتون نسخه ۳.۵.۱
ویرایشگر متن	JetBrains Pycharms 2016
پردازنده	Core i7-4510U
حافظه اصلی	8 GB

جدول ۲: کتابخانه‌های استفاده شده

نام کتابخانه	کاربرد
Scikit-learn	الگوریتم‌های یادگیری ماشین
deap	الگوریتم‌های ژنتیک
numpy	محاسبات ریاضی به صورت ماتریسی

۳-۵- ارزیابی روش ارائه شده

در این تحقیق از الگوریتم جنگل تصادفی برای شناسایی حملات محروم سازی از سرویس توزیع شده استفاده شده است. نوآوری اصلی تحقیق در استفاده از الگوریتم‌های ژنتیک برای انتخاب زیرمجموعه بهینه از ویژگی‌های مجموعه داده است تا عملکرد الگوریتم جنگل تصادفی بهبود یابد. در جدول ۳ و ۴ نتایج ارزیابی الگوریتم فوق با و بدون اعمال الگوریتم‌های ژنتیک اعمال شده بر روی مجموعه داده‌های KDDTest⁺ و KDDTest⁻²¹ آورده شده است.

همان‌طور که در جدول‌های شماره ۳ و ۴ نیز مشخص است، الگوریتم جنگل تصادفی که با الگوریتم‌های ژنتیک بهینه شده‌اند، نسبت به حالت معمول الگوریتم جنگل تصادفی از کارایی و عملکرد بهتری برخوردار هستند. از نظر زمان اجرای الگوریتم در شرایط مختلف نیز مقایسه‌ای صورت گرفته است که در جدول ۵ این نتایج آمده است. البته دقت کنید این زمان، مدت زمانی است که الگوریتم برای آموزش نیاز دارند که زمان آفلاین محسوب می‌شود.

۴- بحث و نتیجه گیری

در این تحقیق، به ارائه روشی مبتنی بر الگوریتم‌های ژنتیک برای شناسایی حملات محروم سازی از سرویس توزیع شده در بستر رایانش ابری پرداخته شد.

رویکرد اصلی در این تحقیق بهینه ساختن عملکرد طبقه‌بند جنگل تصادفی با استفاده از الگوریتم‌های ژنتیک است. به عبارت دیگر از الگوریتم‌های ژنتیک برای یافتن زیرمجموعه‌ای از ویژگی‌ها استفاده شد که این زیرمجموعه بیشترین دقت در شناسایی حملات محروم سازی از سرویس توزیع شده در الگوریتم درخت تصادفی را سبب می‌شود، در نهایت نیز روش ارائه شده بر روی مجموعه داده معتبری اجرا و نتایج حاصل شده حاکی از عملکرد و دقت بیشتر روش ارائه شده دارد.

- [14] Revathi, S. and A. Malathi, *A detailed analysis on NSL-KDD dataset using various machine learning techniques for intrusion detection*. 2013.
- [15] Cup, K., Dataset. available at the following website <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>, 1999. 72.
- [16] Hanley, J.A. and B.J. McNeil, *The meaning and use of the area under a receiver operating characteristic (ROC) curve*. *Radiology*, 1982. 143(1): p. 29-36.
- [17] Cheadle, C., et al., *Analysis of microarray data using Z score transformation*. *The Journal of molecular*
- [18] Mahmoodi Derakhsh, A., Daneshjoo, P., and Delara, C., *Using Genetic Algorithm to Improve Bernoulli Naïve Bayes Algorithm in Order to Detect DDos Attacks in Cloud Computing Platform*. *International Journal of Science and Engineering Investigations*, Volume7, Issue72, January 2018.

زیر نویس

- 1 Denial-of-service attack (DOS)
- 2 Random Forest
- 3 Genetic Algorithm
- 4 Distributed Denial-of-service attack (DDoS)
- 5 Software Defined Networking (SDN)
- 6 Bernoulli Naïve Bayes
- 7 Tournament Selection
- 8 Bit Diversity
- 9 Ensemble Learning
- 10 Decision Trees
- 11 Validation Dataset
- 12 Rapid Convergence
- 13 Roulette wheel
- 14 Receiver Operating Characteristic
- 15 False Positive Rate (Fall-out)
- 16 True Positive Rate (Sensitivity)
- 17 Confusion Matrix
- 18 Deep Learning

۳) استفاده از فن یادگیری عمیق^{۱۸} یکی از حوزه‌های تحقیقاتی که اخیراً مورد اقبال پژوهشگران قرار گرفته است، مباحث مربوط به یادگیری عمیق است. در مساله تحقیق جاری نیز می‌توان از قدرت شبکه‌های عصبی یادگیری عمیق برای استخراج و یادگیری ویژگی استفاده نمود.

سپاسگزاری

از استادان گرامیم آقای دکتر چنگیز دل آرا و خانم دکتر پریسا دانشجو بسیار سپاسگذارم چرا که بدون راهنمایی‌های ایشان به سرانجام رساندن این پژوهش کاری بسیار مشکل بود.

مراجع

- [1] Wang, B., et al., *DDoS attack protection in the era of cloud computing and software-defined networking*. *Computer Networks*, 2015. 81: p. 308-319.
- [2] Tama, B.A. and K.-H. Rhee, *Data Mining Techniques in DoS/DDoS Attack Detection: A Literature Review*. *International Information Institute (Tokyo)*. Information, 2015. 18(8): p. 3739.
- [3] Li, X., et al. *DDoS Detection in SDN Switches using Support Vector Machine Classifier*. in 2015 Joint International Mechanical, Electronic and Information Technology Conference (JIMET-15). 2015. Atlantis Press.
- [4] Malhi, A.K. and S. Batra, *Genetic-based framework for prevention of masquerade and DDos attacks in vehicular ad-hoc networks*. *Security and Communication Networks*, 2016. 9(15): p. 2612-2626.
- [5] Ambusaidi, M.A., et al., *Building an intrusion detection system using a filter-based feature selection algorithm*. *IEEE transactions on computers*, 2016. 65(10): p. 2986-2998.
- [6] Osanaiye, O., et al., *Ensemble-based multi-filter feature selection method for DDos detection in cloud computing*. *EURASIP Journal on Wireless Communications and Networking*, 2016. 2016(1): p. 130.
- [7] Varma, P.R.K., V.V. Kumari, and S.S. Kumar, *Feature Selection Using Relative Fuzzy Entropy and Ant Colony Optimization Applied to Real-time Intrusion Detection System*. *Procedia Computer Science*, 2016. 85: p. 503-510.
- [8] Liaw, A. and M. Wiener, *Classification and regression by randomForest*. *R news*, 2002. 2(3): p. 18-22.
- [9] Holland, J.H. and J.S. Reitman, *Cognitive systems based on adaptive algorithms*. *ACM SIGART Bulletin*, 1977(63): p. 49-49.
- [10] Grefenstette, J.J., *Optimization of control parameters for genetic algorithms*. *IEEE Transactions on systems, man, and cybernetics*, 1986. 16(1): p. 122-128.
- [11] Mitchell, M., *An introduction to genetic algorithms*. 1998: MIT press.
- [12] Ali, E. and E. Elamin, *A proposed genetic algorithm selection method*. 2006.
- [13] Pereira, F. and G. Gordon. *The support vector decomposition machine*. in *Proceedings of the 23rd international conference on Machine learning*. 2006. ACM.