

تشخیص باتنت برای شبکه‌های نظیر به نظیر

رضا عزمی^۱، مریم قلی‌نژاد^{۲*}، محسن صابری^۳

۱- استادیار، گروه مهندسی کامپیوتر، دانشکده فنی و مهندسی، دانشگاه الزهرا تهران

۲- دانشجو کارشناسی ارشد، گروه مهندسی کامپیوتر، دانشکده فنی و مهندسی، دانشگاه الزهرا تهران

۳- مربی، گروه مهندسی کامپیوتر، دانشکده فنی و مهندسی، دانشگاه بزرگمهر قاننات

(دریافت: ۹۴/۴/۱۶، پذیرش: ۹۴/۱۰/۲۲)

چکیده

باتنت‌ها از جمله جدیدترین نوع بدافزارها در مقیاس اینترنت می‌باشند که در سال‌های اخیر بیش‌ترین تهدیدات را متوجه سامانه‌های اینترنتی نموده‌اند. بات، رایانه‌ای آلوده شده به یک بدافزار است که بدون آگاهی و اراده کاربر و از راه دور توسط یک یا چند عامل انسانی کنترل می‌شود. به این عامل کنترل کننده، سرکرده یا مدیربات گویند و گاهی سیستم آلوده را قربانی نیز می‌نامند. باتنت نظیر به نظیر یکی از انواع باتنت است که از پروتکل‌های نظیر به نظیر برای کنترل بات‌های خود استفاده می‌کند و شناسایی این نوع از باتنت نسبت به انواع دیگر مشکل‌تر است. رویکرد پیشنهادی ما یک راه‌حل برای شناسایی چنین باتنت‌هایی است. این رویکرد با استفاده از تحلیل جریان شبکه و روش‌های خوشه‌بندی در داده‌کاوی، باتنت‌های نظیر به نظیر را شناسایی می‌کند. رویکرد ارائه شده، مبتنی بر جریان بوده و با مقایسه شباهت بین جریان‌ها و خوشه‌بندی جریان‌های مقایسه شده با استفاده از الگوریتم K-Means و در نهایت مقایسه خوشه‌های ترافیک جدید با خوشه‌های قبلی، می‌تواند وجود یا عدم وجود حمله را تشخیص دهد. به این ترتیب که ابتدا یک ترافیک از جنس حمله و یک ترافیک معمولی خوشه‌بندی می‌شود، سپس ترافیک جدید که نوع آن مشخص نیست خوشه‌بندی می‌شود. در نهایت خوشه‌بندی ترافیک جدید با دو ترافیک قبلی مقایسه شده و مشخص می‌شود که ترافیک جدید از چه نوعی است. رویکرد ارائه شده برای تشخیص باتنت از ترافیک ازدحام ناگهانی، عملکرد خوبی را از خود نشان می‌دهد و این خصوصیت، وجه تمایز الگوریتم پیشنهادی، نسبت به الگوریتم‌های مشابه آن است. در نهایت عملکرد رویکرد ارائه شده، با ترافیک‌های مختلف مورد بررسی قرار گرفته می‌شود.

واژه‌های کلیدی: باتنت، شناسایی باتنت، بات، باتنت نظیر به نظیر، تحلیل جریان شبکه

۱- مقدمه

بات‌ها می‌توانند نقش کارگزار^۶ فرمان- کنترل^۷ را بازی کنند و مهاجمین می‌توانند حملات خود را از هر رایانه در باتنت نظیر به نظیر، راه‌اندازی کنند؛ بنابراین شناسایی آن‌ها در مقایسه با دیگر باتنت‌ها مشکل‌تر است.

ویژگی منحصر به فردی که یک باتنت را از دیگر بدافزارها متمایز می‌سازد، زیرساخت ارتباطی آن یعنی کانال فرمان- کنترل است. مدیر بات از کانال فرمان- کنترل برای دستور دادن به بات‌ها به منظور انجام فعالیت‌های مخرب استفاده می‌کند [۱].

مدیران بات برای نگه داشتن باتنت خود از فن‌های غیرقانونی زیادی استفاده می‌کنند از جمله پنهان‌سازی کد^۸، رمزنگاری حافظه^۹، فتاوری پیاده‌سازی نظیر به نظیر و یا تقلید ازدحام ناگهانی^{۱۰}. در اینجا تمرکز ما روی حمله DDOS ای است که با استفاده از تقلید ازدحام ناگهانی صورت می‌گیرد. حمله DDOS

باتنت‌ها^۱ از جمله جدیدترین نوع بدافزارها^۲ در مقیاس اینترنت می‌باشند که در سال‌های اخیر بیش‌ترین تهدیدات را متوجه سامانه‌های اینترنتی نموده‌اند. بات، رایانه‌ای آلوده شده به یک بدافزار است که بدون آگاهی و اراده کاربر و از راه دور توسط یک یا چند عامل انسانی کنترل می‌شود. به این عامل کنترل کننده، سرکرده یا مدیربات^۳ گویند و گاهی سیستم آلوده را قربانی^۴ نیز می‌نامند.

باتنت‌ها اهداف مخرب زیادی دارند و اخیراً مهاجمین، از باتنت‌های نظیر به نظیر^۵ برای راه‌اندازی انواع حملات مانند حملات DDOS، پخش بدافزار، توزیع هرزنامه‌ها و غیره، استفاده می‌کنند. باتنت‌های نظیر به نظیر نسل جدیدی از باتنت‌ها هستند که از شبکه‌های نظیر به نظیر برای کنترل از راه دور بات‌های خود استفاده می‌کنند. در این نوع از باتنت‌ها تمامی

*ایمیل نویسنده پاسخگو: M.gholinezhad84@gmail.com

6- Server

7- Command and Control Server

8- Code Obfuscation

9- Memory Encryption

10- Flash Crowd Mimicking

1- Botnet

2- Malware

3- Bot Master

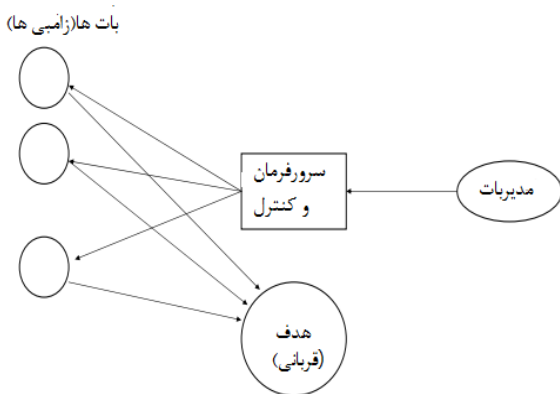
4- Victim

5- Peer to Peer

هستند که با ارسال دستورات خود بات‌ها را، از راه دور کنترل می‌کنند تا فعالیت‌های غیرقانونی یا مخرب انجام دهند.

- هسته ارتباطی بات‌ها، کارسازی به نام "فرمان-کنترل" است که بات‌ها فرمان‌های اجرایی خود را از آن دریافت می‌نمایند. کارگزار فرمان-کنترل، دستورات مدیران بات را برای بات‌ها در رایانه‌های آلوده شده توزیع می‌کند یا به عبارتی دستورات را از مدیر بات دریافت و برای بات‌ها ارسال می‌کند. البته همچنان‌که در بخش‌های بعدی به آن پرداخته خواهد شد، این معماری تنها معماری مرسوم در بات‌نت‌ها نیست ولی عمومی‌ترین معماری است.

عناصر بات‌نت در شکل (۱) نشان داده شده است:



شکل (۱). عناصر یک بات‌نت

عملیات بات‌نت به وسیله ارسال دستورات خاصی از سمت مدیر بات و از طریق کانال فرمان-کنترل آغاز می‌شود. عناصر یک بات‌نت ارتباط نزدیکی باهم دارند و برای این‌که حمله‌ای بخواهد اتفاق بیفتد، زیرساخت ارتباطی یعنی کانال فرمان-کنترل، نیاز خواهد بود. زیرساخت فرمان کنترل^۵ مهم‌ترین مؤلفه بات‌نت و وجه تمایز آن‌ها با دیگر بدافزارهاست که اجازه می‌دهد بات‌ها، دستورات مخرب را از جانب مدیر بات دریافت کنند و امکان هدایت و کنترل بات‌ها را درون یک بات‌نت برای مدیر بات فراهم می‌کند [۷].

۲- انواع معماری بات‌نت‌ها

بات‌نت‌ها با توجه به سازوکاری که برای فرمان-کنترل خود استفاده می‌کنند از دو معماری کلی استفاده می‌کنند. معماری متمرکز و نامتمرکز.

یک حمله بحرانی برای اینترنت است و اغلب، بات‌نت‌ها در پشت این حملات قرار دارند. مدیران بات ماهر، با تقلید الگوی ترافیک ازدحام ناگهانی، فناوری تشخیص را غیرفعال می‌کنند؛ که این مسئله باعث می‌شود برای کسانی که به دنبال دفع حملات بات‌نت هستند یک چالش عظیم ایجاد شود. نکته‌ای که وجود دارد این است که جریان‌های حمله جاری در مقایسه با جریان ازدحام ناگهانی، شباهت بیشتری به هم دارند؛ و طبق همین اصل، الگوریتمی ارائه شده که از ضریب هم‌بستگی جریان، به عنوان یک پارامتر شباهت برای جریان‌های مشکوک استفاده می‌شود.

ازدحام ناگهانی، یک اتفاق ناگهانی و ناخواسته اما در عین حال مشروع است، مثل هجوم ناگهانی برای خواندن یک خبر جدید. یک راه حل خوب برای رخنه‌گرها^۱ شبیه شدن به الگوی ازدحام ناگهانی است تا از دید آشکار کننده‌ها مخفی بماند که به این نوع، حمله ازدحام ناگهانی^۲ می‌گویند.

پژوهش‌های دیگری قبلاً صورت گرفته که تشخیص حمله DDOS را با استفاده از مشخص کردن و استخراج ویژگی‌های حمله DDOS انجام می‌دهند، چنین رویکردهایی به صورت فعال قادر به شناسایی حملات نیستند. راه حل‌های دیگری برای مقابله با چنین حملاتی استفاده می‌شود که در آن از جورچین‌های^۳ گرافیکی برای تفاوت قائل شدن بین انسان و بات‌ها استفاده می‌کنند. بدی این روش این است که استفاده از پاسخ‌های انسانی، می‌تواند برای کاربر آزاردهنده باشد. استفاده از CAPTCHA و رفتارهای انسانی در لایه کاربرد، خوب عمل می‌کند. با این حال، هیچ روش تشخیصی در لایه شبکه مشاهده نشده که بتواند سیستم دفاعی را توسعه دهد. یک تعداد پژوهش‌هایی انجام شده [۴-۲] که در آن‌ها با استفاده از هانی پات، بات‌ها به دام انداخته شده‌اند و روی آن‌ها تحلیل انجام داده شده تا فعالیت‌های بات‌نت را مورد مطالعه قرار دهند.

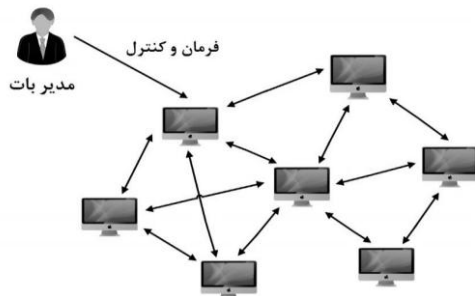
بات‌ها از سه عنصر بات، کارگزار فرمان-کنترل و مدیر بات تشکیل شده‌اند و تهدیدات آن‌ها نیز توسط این عناصر سازماندهی می‌شود. این عناصر به شرح زیر است [۶-۵]

- بات، یک برنامه نرم‌افزاری است که بر روی میزبان‌های آسیب پذیر نصب شده و قادر است اقدامات مخربی را انجام دهد. پس از این‌که برنامه بات بر روی رایانه نصب گردید، آن رایانه به یک بات یا زامبی^۴ تبدیل می‌شود.
- مدیر بات یا هدایت کننده بات، شخص یا گروهی از اشخاص

1- Hacker
2- Flash Crowd Attack
3- Puzzle
4- Zombie

۲-۲- معماری غیرمتمرکز

در معماری نامتمرکز بات‌ها به صورت یک شبکه نظیربه‌نظیر باهم در تعامل می‌باشند. سرکرده دستورات خود را به یک یا چند بات ارسال می‌کند و با استفاده از قراردادهای نظیربه‌نظیر این دستورات در تمام بات‌نت منتشر می‌شود. ویژگی اصلی بات‌نت‌های نظیربه‌نظیر این است که تمام هم‌تاها می‌توانند نقش یک کارگزار فرمان-کنترل را بازی کنند [۸] به عبارت دیگر در بات‌نت‌های نظیربه‌نظیر کارگزار فرمان-کنترل پنهان است [۹]. حسن این روش این است که بات‌ها به کارگزار فرمان-کنترل وابسته نیستند. توزیع ترافیک در این بات‌نت‌ها شناسایی آن‌ها را دشوار کرده است. این بات‌نت‌ها مقیاس‌پذیری بالایی دارند اما پیاده‌سازی آن‌ها پیچیده و دشوار است. بات‌ها به خاطر تأخیر در انتشار دستورات نمی‌توانند به سرعت آرایش بگیرند. حالت خاصی هم از بات‌نت‌های نظیربه‌نظیر موجود است که در آن هر بات فقط با یک بات دیگر در ارتباط است. در نتیجه گراف اتصال حاصل از این شبکه به صورت زنجیر است. این پیاده‌سازی گرچه از مدل کامل نظیربه‌نظیر ساده‌تر است ولی تأخیر انتشار در آن زیاد است و تضمینی نیز برای رسیدن پیام به همه بات‌ها وجود ندارد زیرا این روش از این ضعف رنج می‌برد که در صورت قطع شدن زنجیر بخشی از بات‌نت از دسترس خارج می‌گردد در نتیجه چنین ساختاری از استحکام مطلوبی برخوردار نیست.

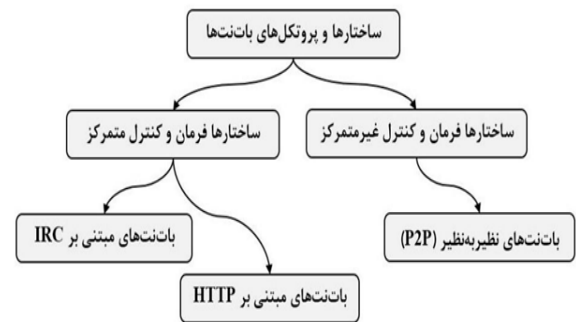


شکل (۴). ساختار فرمان-کنترل غیرمتمرکز

جدول (۱) ساختار فرمان-کنترل متمرکز و غیرمتمرکز را باتوجه به عوامل بیان‌شده، مورد مقایسه قرار داده است.

جدول (۱). مقایسه ساختار فرمان-کنترل متمرکز و غیرمتمرکز

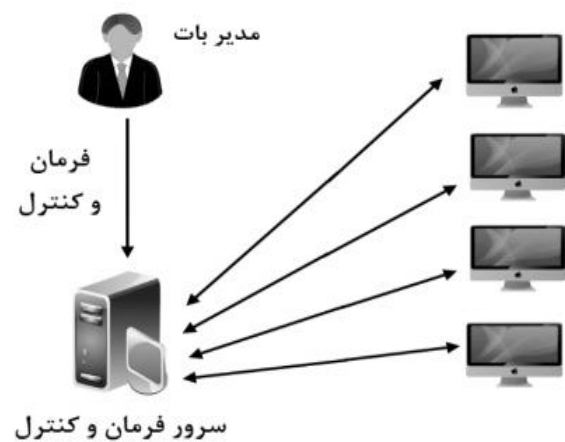
غیرمتمرکز	متمرکز	ساختار فرمان و کنترل عامل مقایسه
زیاد	کم	پیچیدگی طراحی
زیاد	کم	زمان تأخیر پیام
کم	زیاد	قدرت آشکارسازی
زیاد	کم	انعطاف‌پذیری



شکل (۲). ساختارها و پروتکل‌های بات‌نت‌ها

۲-۱- معماری متمرکز

در معماری متمرکز^۱ که عمومی‌ترین معماری نیز است، یک یا چند کارگزار فرمان-کنترل وجود دارد و تمام بات‌ها دستورات خود را از این کارگزار دریافت می‌کنند. حسن معماری متمرکز پیاده‌سازی آسان و امکان آرایش‌دهی سریع بات‌ها است به این معنی که با ارسال دستورات در مدت‌زمان کوتاهی می‌توان بات‌ها را برای هدف خاصی آرایش داد. همچنین بازخورد مستقیم، مدیر بات را قادر می‌سازد به راحتی وضعیت بات‌نت را بررسی کند و به اطلاعات مهمی درباره بات‌نت از قبیل تعداد بات‌های فعال، دست‌یاب [۶]. عیب اساسی این روش این است که اگر به‌دلیلی کارگزار فرمان-کنترل از دسترس خارج شود، عملاً بات‌نت از کار می‌افتد. دوم این‌که این روش به راحتی قابل شناسایی است. زیرا دستورات مشابه‌دادن به یک تعداد گره در شبکه و یا اتصال چندین گره از یک شبکه به یک سیستم خاص، شک دستگاه‌های تشخیص نفوذ^۲ (IDS) را برخواهد انگیخت.



شکل (۳). ساختار فرمان-کنترل متمرکز

۳- چرخه حیات بات‌نت

به‌منظور این‌که یک میزبان آسیب‌پذیر به یک بات و بخشی از یک بات‌نت تبدیل شود، باید از یک چرخه از مراحل عبور کند. این مراحل تحت عنوان چرخه حیات بات‌نت^۱ بیان می‌شوند. چرخه حیات بات‌نت‌ها در شکل (۵) نشان داده شده است:



شکل (۵). چرخه حیات یک بات‌نت

چرخه حیات بات‌نت شامل مراحل زیر است [۱۰]:

۱. شکل‌گیری: مرحله شکل‌گیری مرحله‌ای است که مهاجم از آسیب‌پذیری مشخصی در رایانه هدف، سوءاستفاده کرده و از دسترسی به دست آمده برای نصب برنامه مخرب، استفاده می‌کند. پس از نصب برنامه مخرب، رایانه قربانی برنامه را اجرا کرده و به یک بات تبدیل می‌شود.
۲. فرمان - کنترل (C&C): در مرحله فرمان - کنترل بات‌ها تلاش می‌کنند ارتباطی را با کارگزار فرمان - کنترل خود برقرار سازند و از طریق ارتباط ایجاد شده به بات‌نت بپیوندند.
۳. حمله: پس از این‌که بات‌ها ارتباط خود را با کارگزار فرمان - کنترل برقرار کردند، مرحله حمله آغاز می‌شود. بات‌ها از طریق کانال فرمان - کنترل، دستورات مدیربات را دریافت می‌کنند و فعالیت‌های مخرب را براساس دستورات دریافت شده، انجام می‌دهند.
۴. پس از حمله: در این مرحله ممکن است مدیربات، برنامه بات را برای بهبود قابلیت‌ها یا برخی تغییرات دیگر به‌روزرسانی کند. همچنین ممکن است روش‌های جدیدی را برای جلوگیری از تشخیص بات‌ها و مخفی نگه‌داشتن آن‌ها به‌کار ببرد.

۴- کاربرد داده‌کاوی در شناسایی بات‌نت‌ها

امروزه شرکت‌ها و سازمان‌ها حجم انبوهی از داده‌ها را در پایگاه‌های داده خود در اختیار دارند. در این راستا نیاز است تا روش‌های خودکاری وجود داشته باشد تا بتواند دانش مفید را از

میان انبوه داده‌ها، کشف و استخراج کنند. داده‌کاوی ابزار مؤثری است که برای این منظور مورد استفاده قرار می‌گیرد. در زمینه‌ی شناسایی بات‌نت‌ها، روش‌های داده‌کاوی می‌توانند با کشف الگوهای رفتار مخرب بات‌نت‌ها، آن‌ها را از ترافیک عادی شبکه متمایز کرده و شناسایی کنند [۱۱]. روشی که در اینجا برای شناسایی بات‌نت از آن استفاده شده است روش خوشه‌بندی داده‌هاست.

در خوشه‌بندی، داده‌ها براساس قاعده "حداکثر مشابهت درون کلاسی و حداقل مشابهت برون کلاسی" خوشه‌بندی یا گروه‌بندی می‌شوند. درنهایت خوشه‌هایی از داده‌ها شکل می‌گیرند که داده‌های درون یک خوشه بالاترین مشابهت را با یکدیگر و کم‌ترین مشابهت را با داده‌های دیگر خوشه‌ها دارند. هر خوشه‌ای که شکل گرفته، می‌تواند به‌عنوان یک کلاس از داده‌ها در نظر گرفته شود و قواعدی از آن استنتاج شوند [۱۲].

خوشه‌بندی، فرایند گروه‌بندی مجموعه‌ای از داده‌ها به کلاس دارای داده‌های مشابه است. در خوشه‌بندی، هیچ اطلاعی از کلاس‌های موجود درون داده‌ها وجود ندارد و به‌عبارتی خود خوشه‌ها نیز از داده‌ها استخراج می‌شود [۱۳].

۵- دسته‌بندی رویکردهای شناسایی بات‌نت

براساس روش استفاده‌شده برای تشخیص

از این دیدگاه، دو رویکرد کلی می‌توان برای تشخیص بات‌نت در نظر گرفت:

۱- استفاده از هانی‌پات: هانی‌پات شبکه‌ای از کوزه‌های عسل است که به‌عنوان طعمه در داخل شبکه قرار می‌گیرد و اگر به یک بات آلوده شد، تمام جزئیات تعاملی و رفتاری آن را ثبت می‌نماید و می‌تواند این آلودگی را گزارش کند، که منجر به شناسایی بات‌نت خواهد شد.

۲- استفاده از شنود غیرفعال: این رویکرد با شنود و تحلیل ترافیک شبکه وجود شبکه‌های بات را تشخیص می‌دهد. این رویکرد نیز به چهار دسته تقسیم می‌شود:

- مبتنی بر امضا: در این روش سعی می‌شود که با تحلیل بسته‌های ورودی/خروجی و یافتن الگوهای^۸ خاصی در آن‌ها به‌نام امضا و مقایسه این الگوها با بانک الگوی از قبل آماده‌شده، وجود بات‌نت بررسی شود. این روش در دستگاه‌های تشخیص نفوذ مانند snort به‌کار برده می‌شود.
- مبتنی بر ناهنجاری: در این روش از ناهنجاری‌هایی که فعالیت بات‌نت بر روی ترافیک طبیعی شبکه ایجاد می‌کند، سعی در شناسایی آن دارند. مقاله بینکلی^{۱۰} و

6- Honeynet
7- Signature
8- Pattern
9- Anomaly
10- Binkley

1- Botnet Life Cycle
2- Formation
3- Command and Control
4- Attack
5- Post-Attack

۶-۱- شناسایی مبتنی بر تجزیه و تحلیل محتوا

در روش‌های شناسایی مبتنی بر تجزیه و تحلیل محتوا^۳، محتوای بسته‌های مبادله‌شده در شبکه، تجزیه و تحلیل می‌شوند تا براساس تجزیه و تحلیل انجام‌شده و ویژگی‌های به‌دست‌آمده، ترافیک بات‌نت شناسایی شود. این روش‌ها اگرچه دقت تشخیص بالایی را نشان می‌دهند، اما از وجود برخی محدودیت‌ها رنج می‌برند [۶ و ۱۷]:

- در روش‌های شناسایی مبتنی بر تجزیه و تحلیل محتوا، نیاز است حجم زیادی از داده‌های بسته‌های ترافیک تجزیه و تحلیل شوند، بنابراین روش‌هایی زمان‌بر هستند و حجم پردازشی^۴ زیادی را دربردارند.
- روش‌های شناسایی مبتنی بر تجزیه و تحلیل محتوا نمی‌توانند بات‌نت‌هایی را که از بسته‌های رمزنگاری‌شده در کانال‌های ارتباطی خود استفاده می‌کنند، شناسایی کنند. زیرا بسته‌های ترافیک رمزنگاری‌شده و اطلاعات محتوای بسته در دسترس نیست.
- در روش‌های شناسایی مبتنی بر تجزیه و تحلیل محتوا، ممکن است به دلیل مسائل قانونی و حریم خصوصی، امکان خواندن محتوای بسته‌های ترافیک و دستیابی به اطلاعات محتوا وجود نداشته باشد.

۶-۲- شناسایی مبتنی بر تجزیه و تحلیل رفتار

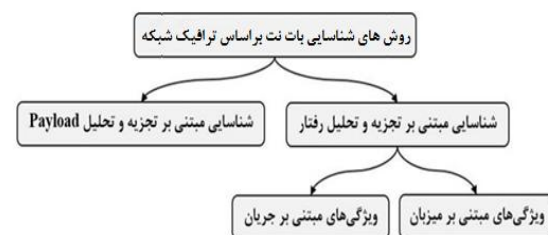
میزبان‌های آلوده به بات در یک بات‌نت به‌خاطر از پیش کدشدن، رفتار ترافیکی مشابهی را از خود نشان می‌دهند. در روش‌های شناسایی مبتنی بر تجزیه و تحلیل رفتار ترافیک از این ایده استفاده می‌شود که بات‌ها درون یک بات‌نت معمولاً رفتار ترافیکی یکنواختی دارند و الگوهای ترافیکی خاصی را به‌منظور برقراری ارتباط در شبکه از خود نشان می‌دهند [۱۷]. این روش‌ها با استخراج ویژگی‌های مختلفی که در ترافیک شبکه رخ می‌دهد، مانند اندازه بسته، مدت‌زمان جریان و غیره، رفتارها و الگوهای ترافیکی بات‌نت‌ها را مشخص کرده و ترافیک آن‌ها را از ترافیک عادی شبکه متمایز می‌سازند [۱۱]. مزیت روش‌های شناسایی مبتنی بر تجزیه و تحلیل رفتار ترافیک در این است که به محتوای بسته‌های ترافیک وابسته نیستند و می‌توانند بات‌نت‌های با ترافیک رمزگذاری‌شده را شناسایی کنند. همچنین، اطلاعات مربوط به ترافیک شبکه معمولاً می‌تواند به‌راحتی از دستگاه‌های مختلف شبکه بازیابی شود، بدون این‌که تأثیر قابل‌توجهی بر عملکرد شبکه یا دسترسی خدمات داشته باشد [۱۸] ویژگی‌ها می‌توانند براساس خصوصیات میزبان‌ها و یا

همکاران [۱۴] که توضیح داده خواهد شد نمونه‌ای از این روش‌ها است.

- مبتنی بر DNS: این روش در حقیقت یک روش مبتنی بر ناهنجاری است ولی از آنجاکه در بسیاری موارد، پرس‌وجوهای DNS عاملی برای تشخیص بات‌نت است به‌عنوان روشی مستقل نیز از آن یاد می‌شود. بات‌ها برای یافتن کارگزار فرمان-کنترل خود، درخواست‌هایی به DNS می‌دهند. از آنجاکه اغلب کارگزار فرمان-کنترل از فن‌آوری DNS پویا برای مخفی نگه‌داشتن خود استفاده می‌کند، درخواست‌های مکرر بات‌ها یک نوع ناهنجاری در ترافیک DNS ایجاد می‌کند که تشخیص این ناهنجاری به شناسایی بات‌نت منتهی خواهد شد.
 - مبتنی بر داده‌کاوی: این روش‌ها گستره‌ای را در برمی‌گیرند و معمولاً بر تحلیل‌های آماری و داده‌کاوی و هوش مصنوعی استوار می‌باشند. به‌عنوان نمونه هسنا^۱ و همکاران [۱۵] در حملات هرزنامه‌ای ناشی از بات‌نت‌ها ویژگی‌هایی را مدنظر قرار داده‌اند و توانسته‌اند با تحلیل مؤلفه‌های اصلی^۲ (PCA)، ویژگی‌های با واریانس بیش‌تر را یافته و حملات را تشخیص دهند. الگوریتم پیشنهادی که در فصل‌های بعدی ارائه خواهد شد، نیز مبتنی بر داده‌کاوی است.
- روش‌های نامبرده هرکدام نقاط قوت و ضعف خود را دارند. جدول (۲)، مقایسه‌ای کلی بر روش‌های مطرح‌شده نشان می‌دهد.

۶-۳- دسته‌بندی رویکردهای شناسایی بات‌نت براساس ترافیک شبکه

روش‌های شناسایی بات‌نت‌ها، روش‌های شناسایی مبتنی بر ترافیک شبکه هستند که به نظارت بر ترافیک شبکه تمرکز دارند. این روش‌ها به دودسته تجزیه و تحلیل مبتنی بر محتوا و تجزیه و تحلیل مبتنی بر رفتار تقسیم‌بندی می‌شوند [۱۶] [۱۷] که در ادامه بیش‌تر توضیح داده شده است.



شکل (۶). دسته‌بندی روش‌های شناسایی بات‌نت براساس ترافیک شبکه

3- Payload
4- Overhead

1 - Husna
2- Principal Component Analysis

ویژگی‌های مبتنی بر جریان: ویژگی‌هایی هستند که از جریان‌های شبکه استخراج می‌شوند و برای نسبت دادن جریان‌ها به کلاس خاصی از ترافیک شبکه مثل ترافیک بات یا ترافیک غیربات، مورد استفاده قرار می‌گیرند. جریان به مجموعه‌ای از بسته‌ها اطلاق می‌شود که آدرس IP مبدأ، درگاه مبدأ، آدرس IP مقصد، درگاه مقصد و پروتکل یکسان دارند [۱۶]. استخراج ویژگی مبتنی بر جریان، زمان کمتری را نسبت به ویژگی‌های مبتنی بر میزبان صرف می‌کند، زیرا فضای تجزیه و تحلیل، فقط اطلاعات مربوط به جریان است.

جریان‌های ترافیک شبکه باشند و به دو دسته ویژگی‌های مبتنی بر میزبان و ویژگی‌های مبتنی بر جریان گروه‌بندی می‌شوند [۴ و ۱۱] که در ادامه در مورد آن‌ها توضیح داده شده است. ویژگی‌های مبتنی بر میزبان: ویژگی‌هایی هستند که از الگوهای ارتباطی میزبان‌ها استخراج می‌شوند. این ویژگی‌ها در ارتباطات بین یک میزبان و دیگر میزبان‌های شبکه رخ می‌دهند و برای شناسایی میزبان‌های با الگوهای ارتباطی مشترک، مفید هستند. استخراج ویژگی‌های مبتنی بر میزبان، نیازمند تجزیه و تحلیل هر بسته متعلق به یک میزبان خاص است و زمان زیادی را نسبت به استخراج ویژگی‌های مبتنی بر جریان، صرف می‌کند.

جدول (۲). مقایسه روش‌های تشخیص بات‌نت (۱)

روش تشخیص	مزایا	معایب
مبتنی بر امضا	الگوریتم ساده و دقیق	عدم شناسایی بات‌های جدید افزایش نرخ منفی کاذب برای بات‌های جدید افزایش نرخ مثبت کاذب هنگام تنوع بات‌ها سنگین بودن الگوریتم به دلیل داشتن پایگاه داده بزرگ از امضاها عبور بات‌نت از سیستم، با استفاده از فن‌های مبهم کدگذاری
مبتنی بر ناهنجاری	داشتن عمومیت قادر بودن به شناسایی بات‌های جدید مقاوم در برابر ترافیک رمز شده سبک‌تر نسبت به روش مبتنی بر امضا	سربار محاسباتی بیش‌تر نسبت به الگوریتم‌های مبتنی بر امضا اعلان خطای بالا
مبتنی بر DNS	داشتن تمام مزایای مربوط به الگوریتم‌های مبتنی بر ناهنجاری پیاده‌یازی راحت‌تر نسبت به الگوریتم‌های مبتنی بر ناهنجاری	داشتن تمام معایب مربوط به الگوریتم‌های مبتنی بر ناهنجاری
مبتنی بر داده‌کاوی	شناسایی بات‌های جدید مقاوم در برابر رمزنگاری	پیچیدگی الگوریتم سرعت عمل کم‌تر

استفاده‌شده مبتنی بر میزبان بوده و استخراج آن‌ها زمان زیادی را صرف می‌کند. رویکرد آن‌ها برای شناسایی نیاز دارد تا فعالیت مخرب توسط بات انجام شود و برای شناسایی در مرحله فرمان - کنترل مناسب نیست. به عبارت دیگر، امکان شناسایی بات در مرحله فرمان - کنترل وجود ندارد. نتایج ارزیابی آن‌ها نشان‌دهنده نرخ تشخیص ۹۵/۷ درصد و نرخ مثبت کاذب ۳/۵ درصد برای شناسایی بات‌های نظیربه‌نظیر است.

استوانویک^۴ و همکاران [۲۱] یک سیستم تشخیص مبتنی بر جریان را برای شناسایی ترافیک بات‌نت‌های نظیربه‌نظیر ارائه داده‌اند. سیستم آن‌ها از دو جزء تشکیل شده است: پیش‌پردازش و طبقه‌بند^۵. در مرحله پیش‌پردازش، با تجزیه و تحلیل ترافیک شبکه در سطح جریان، ویژگی‌های استخراج شده‌اند و در مرحله طبقه‌بند، جریان‌های ترافیک با استفاده از الگوریتم‌های یادگیری ماشین^۶ نظارتی، به صورت مخرب و غیرمخرب طبقه‌بندی می‌شوند. در سامانه‌ی آن‌ها، تعداد ۳۹ ویژگی از جریان‌های ترافیک استخراج شده و عملکرد ۸ طبقه‌بند، مورد ارزیابی

۷- معرفی چند پژوهش برای تشخیص بات‌نت

نظیربه‌نظیر

تارنج^۱ و همکاران [۱۹] یک روش شش مرحله‌ای را برای شناسایی جریان‌های ترافیک بات‌نت‌های نظیربه‌نظیر در مرحله فرمان - کنترل از چرخه حیات ارائه داده‌اند که مبتنی بر محتوا است. برای تقسیم‌بندی جریان‌های ترافیک استفاده شده است. در ادامه تعدادی از رویکردهای شناسایی بات‌نت‌های نظیربه‌نظیر که مبتنی بر تجزیه و تحلیل رفتار ترافیک هستند، بیان شده است. چن^۲ و همکاران [۲۰] رویکردی را برای شناسایی بات‌های نظیربه‌نظیر مبتنی بر میزبان، ارائه داده‌اند. رویکرد آن‌ها تشخیص رفتارهای مخرب و ارتباطات نظیربه‌نظیر را هدف قرار داده است. آن‌ها برای تشخیص رفتارهای مخرب، فراخوانی توابع API تولیدشده توسط برنامه بات بر روی میزبان را تجزیه و تحلیل کرده‌اند. برای شناسایی ارتباطات نظیربه‌نظیر، از ویژگی‌های رفتاری میزبان مانند آدرس IP، شماره پورت و نوع پروتکل استفاده شده است. در رویکرد آن‌ها، ویژگی‌های

3- False Positive

4 - Stevanovic

5- Classifier

6- Machine Learning

1- Tarng

2- Chen

محاسباتی، امکان شناسایی بات‌نت‌ها به‌صورت سریع و بلادرنگ وجود داشته باشد. مسئله این پژوهش شناسایی بات‌نت‌های نظیربه‌نظیر است. باتوجه به این‌که شناسایی براساس تجزیه‌وتحلیل رفتار ترافیک صورت می‌گیرد، شناسایی مستقل از محتوای بسته‌های ترافیک بوده و امکان تشخیص بات‌نت‌های نظیربه‌نظیر با ترافیک رمزنگاری شده وجود دارد. در این پژوهش از ویژگی‌های مبتنی بر جریان برای شناسایی بات‌نت استفاده شده است. و بر اساس میزان شباهت^۲ بین جریان‌ها، طبق قواعد خاصی کار شناسایی انجام می‌شود. نکته‌ای که باعث شده این پژوهش نسبت به پژوهش‌های مشابه انجام‌شده، برتری داشته باشد این است که با استفاده از این الگوریتم می‌توان ترافیک بات‌نت را از ترافیک ازدحام ناگهانی تشخیص داد. ترافیک ازدحام ناگهانی ترافیک مشروع اما ناخواسته‌ای است که به‌دلیل سیل عظیم همزمان کاربران به یک سایت به‌وجود می‌آید. این مسئله در بخش‌های بعدی به‌طور مفصل شرح داده خواهد شد. در ادامه چند پژوهش دیگر در زمینه شناسایی بات‌نت معرفی می‌شود.

جدول (۳) خلاصه رویکردهای بیان‌شده برای شناسایی

بات‌نت‌ها را نشان می‌دهد.

جدول (۳). خلاصه رویکردهای شناسایی بات‌نت‌ها

نرخ تشخیص	نوع معماری شبکه بات‌نت	مرحله شناسایی	روش شناسایی از دیدگاه دوم	روش شناسایی از دیدگاه اول	رویکرد
-	نظیربه‌نظیر	فرمان- کنترل	مبتنی بر محتوا	مبتنی بر داده‌کاوی	تارنج و همکاران [۱۹]
نرخ تشخیص ۹۵/۷ درصد و نرخ مثبت کاذب ^۳ ۲/۵ درصد	نظیربه‌نظیر	حمله	مبتنی بر میزبان	مبتنی بر ناهنجاری	چن و همکاران [۲۰]
نرخ تشخیص ترافیک رابین تا ۹۵/۵ درصد	نظیربه‌نظیر	فرمان- کنترل	مبتنی بر جریان	مبتنی بر داده‌کاوی	استوانویک و همکاران [۲۱]
نرخ تشخیص بالای ۹۰ درصد و نرخ خطای کم‌تر از ۷ درصد	نظیربه‌نظیر	فرمان- کنترل	مبتنی بر میزبان و جریان	مبتنی بر داده‌کاوی	سعد و همکاران [۱۸]
نرخ تشخیص ۸۵ درصد	نظیربه‌نظیر	فرمان- کنترل	مبتنی بر جریان	مبتنی بر داده‌کاوی	روش پیشنهادی

قرارگرفته است. سه طبقه‌بند درخت تصمیم C4.5، Random Forest و Random Tree عملکرد بهتری را نسبت به دیگر طبقه‌بندها نشان داده‌اند. عملکرد طبقه‌بندها، نرخ تشخیص ترافیک را بین ۹۵/۵ تا ۹۶/۵ درصد نشان داده است.

سعد^۱ و همکاران [۱۸] رویکردی را برای شناسایی بات‌نت‌های نظیربه‌نظیر در مرحله فرمان- کنترل با استفاده از رفتار ترافیک شبکه ارائه داده‌اند. در رویکرد آن‌ها از ویژگی‌های مبتنی بر میزبان و مبتنی بر جریان برای شناسایی ترافیک بات‌نت‌های نظیربه‌نظیر استفاده شده است. آن‌ها تعداد ۱۷ ویژگی را از جریان‌های شبکه و الگوهای ارتباطی میزبان‌ها استخراج کرده‌اند. استخراج ویژگی‌های مبتنی بر میزبان زمان بیشتری را نسبت به ویژگی‌های مبتنی بر جریان صرف می‌کند. در رویکرد آن‌ها، عملکرد پنج روش یادگیری ماشین برای طبقه‌بندی ترافیک شبکه، بررسی و ارزیابی شده است. نتایج رویکرد آن‌ها به‌منظور تشخیص ترافیک فرمان- کنترل بات‌نت‌های نظیربه‌نظیر، نرخ تشخیص بالای ۹۰ درصد و نرخ خطای کم‌تر از ۷ درصد را برای طبقه‌بندهای ANN، SVM و Nearest Neighbors نشان می‌دهد.

پژوهش‌های مختلفی به‌منظور شناسایی بات‌نت‌های نظیربه‌نظیر موردبررسی و ارزیابی قرار گرفتند. نتایج ارزیابی نشان می‌دهد که اگرچه تاکنون رویکردهای مختلفی برای شناسایی بات‌نت‌های نظیربه‌نظیر ارائه شده است، اما این رویکردها با برخی مشکلات و محدودیت‌ها مواجه‌اند که عبارتند از:

برخی رویکردها (مانند رویکرد چن و همکاران)، بات‌نت‌ها را باتوجه به فعالیت مخرب آن‌ها شناسایی می‌کنند. چنین رویکردهایی فقط قادرند بات‌نت را در مرحله حمله از چرخه حیات شناسایی کنند و امکان شناسایی در مرحله فرمان- کنترل وجود ندارند. رویکردی مؤثر است که بتواند بات‌ها را پیش از این که اقدام به حمله کنند، شناسایی کند.

برخی رویکردها (مانند رویکرد تارنج و همکاران)، شناسایی را براساس محتوای بسته‌های ترافیک انجام می‌دهند. این رویکردها نمی‌توانند بات‌نت‌هایی را که از بسته‌های رمزنگاری شده در کانال‌های ارتباطی خود استفاده می‌کنند، شناسایی کنند. برخی رویکردها که شناسایی را براساس ویژگی‌های مبتنی بر جریان انجام می‌دهند، تعداد زیادی ویژگی را از جریان‌های شبکه استخراج کرده‌اند (مانند رویکرد استوانویک و همکاران که تعداد ۳۹ ویژگی را از جریان‌های ترافیک استخراج کرده‌اند). نکته مهم در شناسایی بات‌نت‌ها، انتخاب بهترین و کم‌ترین تعداد ویژگی است تا با کم‌ترین حجم

در داده‌های متوالی است. به‌عنوان مثال اگر اندازه جریان شبکه داده شده X_i باشد، N باشد، جریان شبکه به‌صورت زیر نمایش داده می‌شود:

$$X_i = \{x_i[1], x_i[2], \dots, x_i[N]\} \quad (1-8)$$

که در آن، $x_i[k]$ ، تعداد بسته‌هایی است که در k امین بازه زمانی برای جریان شبکه، شمارش شده است. شدت جریان^۶؛ برای یک جریان شبکه X_i که اندازه جریان شبکه، N باشد ($N \geq 1$)، امیدریاضی جریان را به‌عنوان شدت جریان تعریف می‌کنیم:

$$E[X_i] = \frac{1}{N} \sum_{n=1}^N x_i[n] \quad (2-8)$$

شدت جریان، میانگین نرخ بسته‌های یک جریان شبکه را نشان می‌دهد. اگر X_i یک جریان حمله DDOS باشد، به $E[X_i]$ شدت حمله می‌گوییم.

اثر انگشت جریان^۷؛ برای جریان شبکه داده شده X_i با اندازه N ، اثر انگشت جریان X_i نمایش یک‌دستی از X_i است:

$$X_i' = \{x_i'[1], x_i'[2], \dots, x_i'[N]\} = \left\{ \frac{x_i[1]}{N.E[X_i]}, \frac{x_i[2]}{N.E[X_i]}, \dots, \frac{x_i[N]}{N.E[X_i]} \right\} \quad (3-8)$$

طبق این تعریف می‌دانیم که، $\sum_{k=1}^N x_i'[k] = 1$ اگر X_i و X_j ($i \neq j$) دو جریان شبکه با اندازه یکسان N باشند همبستگی بین دو جریان^۸ به‌صورت زیر تعریف می‌شود:

$$r_{x_i x_j} = \frac{1}{N} \sum_{n=1}^N x_i[n] x_j[n] \quad (4-8)$$

همبستگی جریان برای نشان دادن میزان شباهت بین جریان‌های مختلف استفاده می‌شود. در بعضی موارد ممکن است همبستگی، مقدار صفر را نشان دهد درحالی‌که دو جریان کاملاً هم بسته باشند، دلیل این امر تفاوت فاز بین دو جریان است، بنابراین این تعریف به‌صورت زیر توسعه داده شد:

$$r_{x_i x_j}[k] = \frac{1}{N} \sum_{n=1}^N x_i[n] x_j[n+k] \quad (5-8)$$

که در آن، k ($k=0,1,2,3,\dots,N-1$) مکان جریان X_j بعد از اعمال تغییر فاز است.

ضریب همبستگی جریان^۹؛ اگر X_i و X_j ($i \neq j$) دو جریان شبکه با اندازه یکسان N باشند، ضریب همبستگی جریان را به‌صورت زیر تعریف می‌کنیم:

$$\rho_{x_i x_j}[k] = \frac{r_{x_i x_j}[k]}{\frac{1}{N} \left[\sum_{n=1}^N x_i^2[n] \sum_{n=1}^N x_j^2[n] \right]^{1/2}} \quad (6-8)$$

۸- طرح پیشنهادی (روش تشخیص مبتنی بر شباهت)

ایده اولیه رویکردی که در این پژوهش ارائه شده، برگرفته از کاری است که شوئی یو^۱ و همکاران [۲۲] انجام داده‌اند. ابتدا کار شوئی یو و همکاران مطرح می‌شود سپس الگوریتم پیشنهادی که توسعه شده الگوریتم شوئی یو و همکاران است مطرح خواهد شد.

حقیقی در زمینه بات‌نت وجود دارد که در طول این پژوهش مدنظر قرار گرفته شده‌اند:

- ابزارهای حمله معمولاً در یک بات‌نت مشابه است. مدیر بات یک فرمان را برای تمام بات‌های بات‌نتش برای شروع یک حمله در نظر می‌گیرد.
- جریان‌های حمله‌ای که در سمت قربانی مشاهده می‌کنیم، تراکم تعداد زیادی از جریان‌های حمله اولیه است و جریان حمله متراکم شده^۲ یک انحراف استاندارد^۳ مشابه را به اشتراک می‌گذارد. انحراف استاندارد جریان حمله اغلب کوچک‌تر از جریان ازدحام ناگهانی اصلی است. دلیل این مسئله این است که تعداد بات‌های زنده بات‌نت به‌مراتب کم‌تر از تعداد کاربران همزمان قانونی ازدحام ناگهانی است (بات‌نت حدود چندصد یا سطح‌های پایینی از چندهزار است و کاربران همزمان ازدحام ناگهانی در سطح صدها هزار نفر هستند).

پس نتیجه این‌که، انحراف استاندارد در جریان حمله کم‌تر از جریان ازدحام ناگهانی است پس شباهت بین جریان‌ها در حمله DDOS بیش‌تر از ازدحام ناگهانی است. بنابراین روش تشخیصی ارائه می‌شود که در آن از ضریب همبستگی جریان برای تشخیص استفاده شده است. در این پژوهش، شناسایی بات‌نت نظیر به‌نظیر در مرحله فرمان - کنترل از چرخه حیات، براساس تجزیه و تحلیل ترافیک انجام می‌شود که در ادامه با جزئیات بیش‌تر شرح داده می‌شود.

۸-۱- تعاریف

جریان شبکه: برای یک مسیر یاب^۴ داده شده در یک شبکه محلی، بسته‌های^۵ شبکه را کلاستر بندی می‌کنیم، به‌گونه‌ای که آن‌هایی آدرس‌های مقصد و مبدا یکسانی دارند به‌عنوان یک جریان شبکه تعریف می‌شوند.

جریان شبکه را با $X_i[n]$ نشان می‌دهیم که در آن i ($i \geq 1$) آندکس جریان‌های شبکه است و n نشان‌دهنده n امین المان

6- Flow Strength
7- Flow Fingerprint
8- Flow Correlation
9- Flow Correlation Coefficient

1- Shui Yu
2- Aggregated Attack Flows
3- Standard Deviation
4- Router
5- Packet

DDOS و ازدحام ناگهانی از ضریب هم‌بستگی جریان متفاوتی استفاده می‌کنند. لازم است که یک حد بالا برای ضریب هم‌بستگی جریان به دست بیاوریم که اگر ضریب هم‌بستگی از این مقدار بیش‌تر باشد این‌طور در نظر می‌گیریم که حمله رخ داده است.

در رابطه‌های (۷-۸) و (۹-۸)، زمانی که از حد آستانه استفاده می‌کنیم، یک تعداد از جریان‌ها را به‌طور کامل کنار می‌گذاریم. در واقع این اطلاعات کلاً دور ریخته می‌شود که این باعث می‌شود دقت الگوریتم کم‌تر شود. برای مقابله با این مشکل، الگوریتم فوق را کمی بهبود داده‌ایم که در ادامه آن را توضیح می‌دهیم.

بعد از این که هم‌بستگی بین هر دو جفت جریان محاسبه شد، حجم زیادی از ضرایب هم‌بستگی بین جفت جریان‌ها را داریم که با استفاده از الگوریتم K-Means خوشه‌بندی می‌شوند. و بهترین خوشه‌بندی با تعداد مراکز مناسب انتخاب می‌شود. روش کار الگوریتم K-Means و چگونگی انتخاب بهترین خوشه‌بندی در فصل بعد بیشتر توضیح داده خواهد شد. این کار برای دو ترافیک بات‌نت و معمولی به‌صورت جدا انجام می‌شود. ترافیک معمولی می‌تواند ترافیک ازدحام ناگهانی و یا یک ترافیک معمولی از شبکه باشد. چراکه الگوریتم ارائه‌شده برای هر دو حالت می‌تواند جواب‌گو باشد. مراحل کار به این صورت است که ابتدا جریان‌ها جدا می‌شوند و هر جریان به‌شکل رابطه (۱-۸) نمایش داده می‌شود، یعنی تعداد بسته‌های جریان برای بازه‌های زمانی مختلف شمارش می‌شود. سپس هم‌بستگی بین هر دو جریان از جریان‌های ترافیک داده‌شده، طبق رابطه (۵-۸) محاسبه می‌شود و با داشتن مقدار هم‌بستگی، ضریب هم‌بستگی بین دو جریان طبق رابطه (۶-۸) محاسبه می‌شود. و در نهایت بعد از اعمال الگوریتم K-Means روی این ضریب‌ها، بهترین خوشه‌بندی انتخاب می‌شود. تمام مراحل فوق برای هر دو ترافیک بات‌نت و ترافیک معمولی انجام می‌شود. مراحل اجرای الگوریتم در شکل (۷) نشان داده شده است:

تا اینجا بهترین خوشه‌بندی برای ترافیک‌ها مشخص شد. حال برای این که بتوان برای یک ترافیک جدید، حمله‌بودن یا نبودن آن را تشخیص داد، به دو طریق می‌توان عمل کرد: اول این که برای ترافیک جدید به‌ازای هر دو جریان موجود در آن، ضریب هم‌بستگی دو جریان محاسبه شود و به خوشه‌بندی‌های ترافیک بات‌نت موجود اضافه شود. چنان‌چه باعث افزایش پراکندگی خوشه‌های موجود شد، بدین معنی است که ترافیک جدید از جنس ترافیک بات‌نت نیست. راه حل دوم که نسبت به راه‌حل قبلی عملکرد بهتری دارد، این است که ترافیک جدید

ضریب هم‌بستگی جریان، برای مشخص کردن شباهت بین دو جریان استفاده می‌شود. زمانی که دو جریان مشابه داشته باشیم که تفاوت فازی داشته باشند، مقدار این ضریب کاهش پیدا می‌کند. در اینجا طبق رابطه (۵-۸) این مشکل با تغییر فاز حل می‌شود.

۲-۸- شناسایی مبتنی بر شباهت

در اینجا روشی برای تشخیص بات‌نت ارائه می‌شود که اساس آن مبتنی بر شباهت جریان‌هاست، و از آن برای مقابله با حمله ازدحام ناگهانی استفاده می‌شود.

هنگامی که یک هشدار حمله DDOS در شبکه اعلام می‌شود، رویترها با استفاده از شمارش تعداد بسته‌ها در یک بازه زمانی داده شده (مثلاً ۱۰۰ میلی‌ثانیه)، شروع به نمونه‌گیری از جریان‌های مشکوک می‌کنند. هنگامی که اندازه جریان موردنظر N به دست آمد، ضریب هم‌بستگی جریان بین جریان‌های مشکوک محاسبه می‌شود.

فرض کنیم که M جریان شبکه وجود دارد، X_1 و X_2 و X_3 ، X_4 ، \dots و X_M ، بنابراین می‌توان ضریب هم‌بستگی را بین هر دو جریان، X_i و X_j ($1 \leq i \leq M$) و ($1 \leq j \leq M$) به دست آورد.

به‌عنوان شاخص شباهت بین دو جریان X_i و X_j در نظر گرفته می‌شود، این شاخص می‌تواند دو مقدار داشته باشد: ۱- برای حمله DDOS و صفر برای حالات دیگر. ∂ نیز به‌عنوان حد آستانه برای تصمیم‌گیری در نظر گرفته می‌شود به‌طوری که:

$$I_{X_i X_j} = \begin{cases} 1, & \rho_{X_i X_j} [k] \geq \partial \\ 0, & \text{otherwise} \end{cases} \quad (7-8)$$

که در آن، $1 \leq i, j \leq M$ ، $i \neq j$. ضریب هم‌بستگی بین هر دو جریان متفاوت طبق رابطه (۷-۸) محاسبه می‌شود و مجموع آن‌ها طبق رابطه (۸-۸) به‌عنوان احتمال وجود حمله در نظر گرفته می‌شود.

$$\Pr(I_A = 1) = \frac{\sum_{1 \leq i, j \leq M, i \neq j} I_{X_i X_j}}{\binom{M}{2}} \quad (8-8)$$

که در آن، I_A شاخص حمله DDOS است. در اینجا هم یک حد آستانه در نظر گرفته می‌شود که اگر مقدار احتمال از این حد آستانه بیش‌تر شود نشان‌دهنده وجود حمله در شبکه است.

$$I_A = \begin{cases} 1, & \Pr(I_A = 1) \geq \partial \\ 0, & \Pr(I_A = 1) < \partial \end{cases} \quad (9-8)$$

به‌عنوان مثال، اگر حد آستانه ۰/۶ باشد، در صورتی حمله DDOS داریم که حداقل ۶۰ درصد مقایسات مثبت باشد. حمله

$$T = \sum_{i=1}^m D_i \quad (12-8)$$

که در آن، V_i مرکز خوشه i ام از ترافیک جدید و V_j مرکز خوشه j ام از ترافیک باتنت، m تعداد خوشه‌های ترافیک جدید و n تعداد خوشه‌های ترافیک باتنت می‌باشد.

مقدار T را نسبت به ترافیک جدید و ترافیک معمولی نیز محاسبه می‌کنیم. بدیهی است هرچه T کمتر باشد ترافیک‌ها به هم شبیه‌ترند. حال با مقایسه این مقدار برای دو ترافیک معمولی و باتنت می‌توان نتیجه گرفت که ترافیک جدید از چه نوعی است. اگر مقدار این شاخص برای ترافیک جدید نسبت به ترافیک باتنت کمتر از ترافیک معمولی باشد نتیجه می‌گیریم که ترافیک جدید از نوع باتنت است. برای روشن‌تر شدن موضوع، در فصل بعد به صورت عملی این را نشان خواهیم داد.

اما مساله‌ای هست که تا اینجا در نظر گرفته نشده است و آن، مساله‌ی جمعیت هر خوشه است. هرچه جمعیت دو خوشه از دو ترافیک متفاوت به هم نزدیک‌تر باشد دلیلی بر شباهت بیشتر این دو خوشه است. از آنجایی که جمعیت کلی، بین دو ترافیک متفاوت است، برای مقایسه، جمعیت هر خوشه تقسیم بر کل جمعیت را به عنوان جمعیت نسبی در نظر می‌گیریم و جمعیت نسبی را به عنوان عامل مقایسه در نظر می‌گیریم:

$$P_i = \frac{\|c_i\|}{\binom{M}{2}} \quad (13-8)$$

که در آن، M تعداد جریان‌های شبکه است و $\|c_i\|$ تعداد عناصر خوشه i ام.

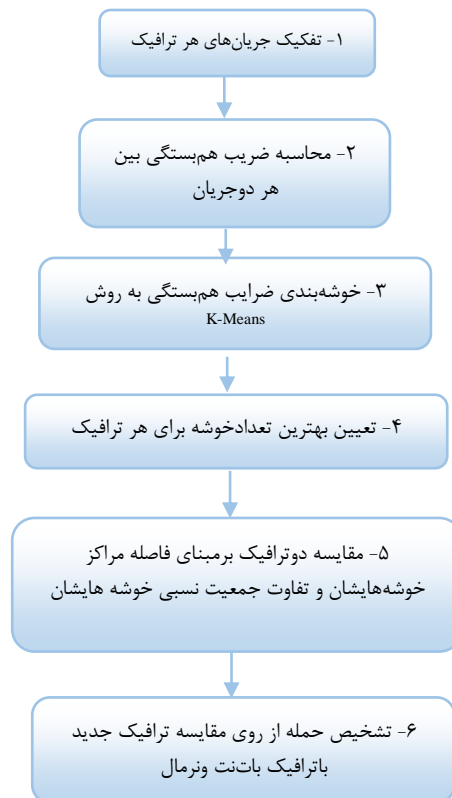
پس اگر بخواهیم دو خوشه را شبیه به هم در نظر بگیریم اولاً مراکز دو خوشه باید به هم شبیه باشند و ثانیاً، جمعیت نسبی دو خوشه باید به هم نزدیک باشد. بنابراین بعد از این که طبق رابطه (۸-۱۰) خوشه‌هایی که مراکز آن‌ها به هم نزدیک است، مشخص شد، جمعیت این دو خوشه با هم مقایسه می‌شود و با اضافه کردن این مقدار به عنوان یک ضریب در رابطه (۸-۱۲) تأثیر جمعیت در نظر گرفته می‌شود. یعنی احتمال شبیه‌بودن دو خوشه، به وسیله تفاوت جمعیت نسبی‌شان مشخص می‌شود. پس رابطه (۸-۱۲) به صورت زیر تعمیم داده می‌شود:

$$T_A = \sum_{i=1}^m A_i D_i \quad (14-8)$$

$$A_i = |P_i - P_j| \quad (15-8)$$

که در آن، $1 \leq i \leq m$ و $1 \leq j \leq n$ و Z اندیس خوشه‌ای است که به ازای آن d_{ij} کمترین مقدار را دارد. بدیهی است که هرچه مقدار فوق کمتر باشد، دو ترافیک به هم شبیه‌ترند و برای تعیین این که ترافیک جدید از چه نوعی است باید مشخص کنیم که به کدام

همانند ترافیک باتنت و معمولی، خوشه‌بندی شود و خوشه‌های جدید با خوشه‌های ترافیک معمولی و باتنت مقایسه شود. و براساس شباهتی که با هر دو این ترافیک‌ها دارد می‌توان تشخیص داد که از کدام نوع است. برای تعیین شباهت بین دو ترافیک می‌توان چند معیار را بررسی کرد. ممکن است این سؤال پیش آید که برای مقایسه، تعداد خوشه‌ها را باید ثابت در نظر گرفت؟ جواب خیر است، چرا که بهترین تعداد خوشه‌ها طبق آزمایش، به دست آمده و ممکن است برای هر ترافیک تعداد خوشه متفاوتی داشته باشیم.



شکل (۷). مراحل اجرای الگوریتم تشخیص باتنت

برای مقایسه ترافیک ناشناخته با ترافیک باتنت فاصله هر یک از مراکز خوشه‌های ترافیک جدید با تمام مراکز خوشه‌های ترافیک باتنت محاسبه می‌شود. مرکزی که کمترین مقدار را دارد، شبیه‌ترین خوشه به این مرکز است. در نهایت مجموع فاصله مراکز ترافیک جدید به شبیه‌ترین خوشه‌هایشان، محاسبه می‌شود:

$$\forall 1 \leq i \leq m, 1 \leq j \leq n \quad (10-8)$$

$$D_i = \min(d_{ij})$$

$$d_{ij} = d(v_i, v_j) = |V_i - V_j| \quad (11-8)$$

$$1 \leq i \leq m, 1 \leq j \leq n$$

ترافیک شبیه‌تر است

در ادامه به بیان چگونگی پیاده‌سازی الگوریتم، جمع‌آوری بسته‌های ترافیک، استخراج جریان‌ها و محاسبه همبستگی می‌پردازیم. سپس با استفاده از الگوریتم K-Means داده‌ها را خوشه‌بندی خواهیم کرد و براساس پارامترهای مشخص شده، عملکرد الگوریتم را با یک نمونه ترافیک نشان می‌دهیم.

۸-۳- استخراج جریان

برای اجرای الگوریتم نیاز به ترافیک بات‌نت و ترافیک معمولی داریم. برای ترافیک بات‌نت، از مجموعه داده‌ای^۱ که $ISCX^2$ ، در اختیار عموم قرار داده^۲، استفاده شده است. قسمتی از این مجموعه داده که حاوی ۲۰۰,۰۰۰ بسته می‌باشد و در بازه زمانی ۸۰۰ ثانیه‌ای تولید شده، برای این کار انتخاب شده است. طبق آن‌چه که برای تعریف جریان قبلاً گفتیم، بسته‌هایی که مبدأ و مقصد یکسان دارند، جدا شده و تحت عنوان یک جریان شناخته می‌شوند. بدین ترتیب ۶۶۰ جریان در این بازه زمانی از مجموعه داده مذکور، استخراج گردید.

ترافیک معمولی شبکه هم ترکیبی از ترافیک برنامه‌های مختلف را شامل می‌شود که عبارتند از: برنامه‌های کاربردی (BitTorrent)، برنامه‌های چت (Skype) و ترافیک وب (Web Browsing). برای این کار، از ترافیک دانشگاه بزرگمهر قاین که حدود ۱۲۰۰۰ دانشجو دارد، در اولین ساعات انتخاب واحد، استفاده شده است. طبق لاگی که در اختیار داریم، در یک بازه زمانی یک ساعته یعنی ۳۶۰۰ ثانیه‌ای، حدود ۲۰۰,۰۰۰ بسته در آن ردوبدل شده است که از این تعداد بسته، ۳۸۰ جریان استخراج شده است. استخراج جریان و اجرای الگوریتم ارائه شده در محیط پایتون انجام شده است و برای ذخیره اطلاعات از پایگاه داده SQL Server استفاده شده است.

برای هر جریان، بازه‌های زمانی که بسته‌ها در آن بازه شمارش می‌شوند، ۵ ثانیه‌ای در نظر گرفته شده، بنابراین مقدار N ، برای ترافیک بات‌نت ۱۶۰ و برای ترافیک معمولی ۷۲۰ خواهد بود. (طبق رابطه (۸-۱)). پس از این‌که جریان‌ها مشخص شد، وضعیت هر بسته، در مورد اینکه بسته مربوط به کدام جریان است، مشخص می‌شود. سپس برای هر دو جریان از ترافیک موجود، مقدار همبستگی بین آن‌ها طبق رابطه (۸-۵) محاسبه می‌شود. مقدار k در این محاسبات صفر در نظر گرفته شده است، چرا که جریان‌های موجود در این مجموعه داده هم‌فاز هستند و چنان‌چه روی آن تغییر فاز اعمال شود، مقادیر نادرستی را برای همبستگی نتیجه می‌دهند. مقدار ضریب

همبستگی بین دو جریان براساس مقدار همبستگی، طبق رابطه (۸-۶) محاسبه می‌شود. حال برای هر دو جریان داده شده، مقدار ضریب همبستگی جریان را داریم که عددی بین ۰ تا ۱ است. الگوریتم K-Means را برای کلاستر بندی جریان‌ها، براساس ضریب همبستگی‌شان استفاده می‌کنیم.

روش K-Means یکی از روش‌های خوشه‌بندی داده‌ها در داده‌کاوی است. این روش علی‌رغم سادگی آن، یک روش پایه برای بسیاری از روش‌های خوشه‌بندی دیگر (مانند خوشه‌بندی فازی) محسوب می‌شود. برای این الگوریتم شکل‌های مختلفی بیان شده است. ولی همه آن‌ها دارای روالی تکراری هستند که برای تعدادی ثابت از خوشه‌ها سعی در تخمین به‌دست آوردن نقاطی به‌عنوان مراکز خوشه‌ها دارند. این نقاط در واقع همان میانگین نقاط متعلق به هر خوشه هستند. هر نمونه داده، به یک خوشه نسبت داده می‌شود، به‌گونه‌ای که آن داده، کم‌ترین فاصله تا مرکز آن خوشه را دارا باشد.

۸-۴- تعیین بهترین تعداد خوشه‌ها

هدف از اعتبارسنجی خوشه‌ها یافتن خوشه‌هایی است که بهترین تناسب را با داده‌های موردنظر داشته باشند. دو معیار پایه اندازه‌گیری پیشنهاد شده برای ارزیابی و انتخاب خوشه‌های بهینه عبارتند از: [۲۳]

- تراکم^۳: داده‌های متعلق به یک خوشه بایستی تا حد ممکن به یک‌دیگر نزدیک باشند. معیار رایج برای تعیین میزان تراکم داده‌ها، واریانس داده‌ها است.
- جدایی^۴: خوشه‌ها خود بایستی به اندازه کافی از یک‌دیگر جدا باشند. سه راه برای سنجش میزان جدایی خوشه‌ها مورد استفاده قرار می‌گیرد که عبارتند از:
 - فاصله^۵ بین نزدیک‌ترین داده‌ها از دو خوشه.
 - فاصله^۶ بین دورترین داده‌ها از دو خوشه.
 - فاصله^۷ بین مراکز خوشه‌ها.

جدول (۴). مجموعه علائم به‌کاررفته در این بخش

تعداد خوشه‌ها	n_c
فاصله بین دو داده	$d(x,y)$
نقطه مرکز خوشه i ام	v_i
امین خوشه I	c_i
تعداد داده‌های درون امین خوشه i	$ c_i $

هم اعمال شده و نتایج نهایی برای ترافیک معمولی هم در انتها نشان داده شده است.

۹-۱- اجرای الگوریتم K-Means

برای اجراهای مختلف الگوریتم با تعداد مراکز متفاوت، نتایج نهایی را در جداول زیر نشان می‌دهیم. در هر یک از موارد زیر، برای تست پایدار بودن مراکز، اجرای الگوریتم را تا ده مرتبه بعد از پایداری، ادامه داده‌ایم. S_i^* طبق رابطه (۸-۱۷) محاسبه شده است.

جدول (۵). اجرای الگوریتم K-Means با دو مرکز (تعداد تکرار

الگوریتم برای رسیدن به پایداری: ۵ بار)

فاصله میانگین (S_i^*)	جمعیت $\ C_i\ $	مراکز نهایی	مراکز اولیه
۰/۰۸۳	۱۲۹۱۷۲	۰/۲۰	۰/۳۳
۰/۱۷۹	۸۸۲۹۸	۰/۷۴	۰/۴۶

جدول (۶). اجرای الگوریتم با سه مرکز (تعداد تکرار الگوریتم برای

رسیدن به پایداری: ۹ بار)

فاصله میانگین (S_i^*)	جمعیت $\ C_i\ $	مراکز نهایی	مراکز اولیه
۰/۰۶۱	۱۰۰۳۸۷	۰/۱۴	۰/۳۳
۰/۰۹۵	۸۰۴۱۶	۰/۵۲	۰/۷۸
۰/۰۵۸	۳۶۶۶۷	۰/۹۶	۰/۹۲

جدول (۷). اجرای الگوریتم با چهار مرکز (تعداد تکرار الگوریتم برای

رسیدن به پایداری: ۱۳ بار)

فاصله میانگین (S_i^*)	جمعیت $\ C_i\ $	مراکز نهایی	مراکز اولیه
۰/۰۳۲	۵۹۹۸۳	۰/۰۶	۰/۱۱
۰/۰۵۹	۶۶۲۵۱	۰/۳۲	۰/۲۳
۰/۰۷۲	۵۶۳۹۵	۰/۵۹	۰/۵۴
۰/۰۴۴	۳۴۸۴۱	۰/۹۷	۰/۸۷

الگوریتم را تا ۸ مرکز اجرا کردیم که از بیان جزئیات آن صرفنظر می‌شود و تنها نتایج نهایی محاسبات را برای این اجراها در زیر نشان می‌دهیم.

۹-۲- تعیین بهترین خوشه‌بندی

تا اینجا برای تعداد خوشه‌های مختلف الگوریتم اجرا شد. حال باید بهترین تعداد خوشه‌ها انتخاب شود. برای تعیین بهترین تعداد خوشه، مقدار DB طبق رابطه (۸-۱۹) برای خوشه‌بندی‌های فوق محاسبه گردیده و نتایج در جدول (۸)

۸-۵- شاخص دیویس بولدین^۱

این معیار از شباهت بین دو خوشه (R_{ij}) استفاده می‌کند که براساس پراکندگی یک خوشه (S_i^*) و عدم شباهت بین دو خوشه (d_{ij}) تعریف می‌شود. شباهت بین دو خوشه را می‌توان به صورت‌های مختلفی تعریف کرد ولی بایستی شرایط زیر را دارا باشد.

$$\bullet R_{ij} \geq 0$$

$$\bullet R_{ij} = R_{ji}$$

• اگر S_i^* و S_j^* هر دو برابر صفر باشند آن‌گاه R_{ij} نیز برابر صفر باشد.

$$\bullet \text{اگر } S_j^* > S_k^* \text{ و } d_{ij} = d_{ik} \text{ آن‌گاه } R_{ij} > R_{ik}$$

$$\bullet \text{اگر } S_j^* = S_k^* \text{ و } d_{ij} < d_{ik} \text{ آن‌گاه } R_{ij} > R_{ik}$$

معمولاً شباهت بین دو خوشه به صورت زیر تعریف می‌شود:

$$R_{ij} = \frac{S_i^* + S_j^*}{d_{ij}} \quad (۸-۱۶)$$

که در آن، d_{ij} طبق رابطه (۸-۱۱) و S_i^* (فاصله میانگین) با روابط زیر محاسبه می‌شوند.

$$S_i^* = \frac{1}{\|C_i\|} \sum_{x \in C_i} d(x, v_i) \quad (۸-۱۷)$$

که در آن، v_i نقطه مرکز خوشه A_i ، C_i ، خوشه A_i و $\|C_i\|$ تعداد داده‌های خوشه A_i است.

باتوجه به مطالب بیان شده و تعریف شباهت بین دو خوشه، شاخص دیویس بولدین به صورت زیر تعریف می‌شود.

$$DB = \frac{1}{n_c} \sum_{i=1}^{n_c} R_i \quad (۸-۱۸)$$

که در آن، n_c تعداد خوشه‌هاست و R_i هم به صورت زیر محاسبه می‌شود:

$$R_i = \max_{j \rightarrow 1..n_c, i \neq j} (R_{ij}), i = 1..n_c \quad (۸-۱۹)$$

این شاخص در واقع میانگین شباهت بین هر خوشه با شبیه‌ترین خوشه به آن را محاسبه می‌کند. می‌توان دریافت که هرچه مقدار این شاخص کم‌تر باشد، خوشه‌های بهتری تولید شده است.

۹- پیاده‌سازی

الگوریتم فوق را با تعداد مراکز مختلف برای هر دو ترافیک معمولی و بات‌نت اجرا کردیم. اجرای الگوریتم بر روی ترافیک بات‌نت، همراه با جزئیات در ادامه شرح داده شده‌است. تمام مراحل که روی این ترافیک انجام شده، بر روی ترافیک معمولی

مراحل فوق را بر روی ترافیک معمولی در دسترس، نیز اجرا کردیم و پس از گذشت تمام مراحل، نتیجه آزمایشات این بود که بهترین خوشه‌بندی برای ترافیک معمولی، خوشه‌بندی با ۴ مرکز می‌باشد. جدول (۱۰) خوشه‌بندی ترافیک معمولی با چهار خوشه را نشان می‌دهد.

جدول (۱۰). خوشه‌بندی ترافیک معمولی

مرکز خوشه (V_i)	جمعیت خوشه ($\ C_i\ $)	فاصله میانگین (S_i^*)	جمعیت نسبی (P_i)
۰/۰۵	۲۶۳۱۰	۰/۰۷۴	۰/۳۶
۰/۱۳	۲۱۴۵۲	۰/۰۸۳	۰/۲۹
۰/۳۶	۱۸۷۴۷	۰/۰۷۵	۰/۲۹
۰/۵۷	۵۵۰۱	۰/۱	۰/۰۷

۹-۳- تشخیص حمله

تا اینجا بهترین خوشه‌بندی را برای یک ترافیک بات‌نت و یک ترافیک معمولی پیدا کردیم. حال می‌خواهیم با استفاده از داده‌هایی که داریم، برای یک ترافیک جدید مخرب بودن یا نبودن آن را تعیین کنیم. ابتدا ترافیک جدید را همانند روال قبلی بررسی کرده و پس از محاسبه ضریب همبستگی جریان‌های ترافیک، الگوریتم K-Means را روی آن اجرا می‌کنیم. همان‌طور که در فصل قبل گفته شد، سپس فاصله بین خوشه‌های ترافیک ناشناخته را با خوشه‌های ترافیک بات‌نت و معمولی طبق رابطه (۸-۱۲) محاسبه می‌کنیم. نتایج محاسبات به صورت زیر می‌باشد:

جدول (۱۱). خوشه‌بندی ترافیک جدید

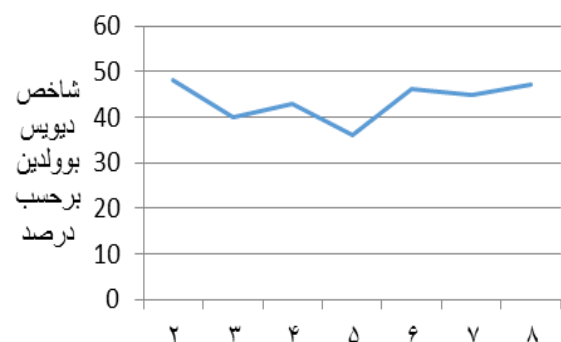
مرکز خوشه (V_i)	جمعیت خوشه ($\ C_i\ $)	فاصله میانگین (S_i^*)	جمعیت نسبی (P_i)
۰/۰۳	۲۶۳۷۶	۰/۰۴۸	۰/۲۱
۰/۱۸	۲۰۶۳۹	۰/۰۴۳	۰/۱۶
۰/۴۵	۳۲۲۱۲	۰/۰۳۷	۰/۲۶
۰/۶۳	۱۹۶۵۹	۰/۰۴۱	۰/۱۶
۰/۸۹	۲۵۸۶۴	۰/۰۳۲	۰/۲۱

ترافیک ناشناخته جدید حاوی ۸۰۰۰۰ بسته در بازه زمانی ۲۰۰۰ ثانیه‌ای است که ۵۰۰ جریان از آن استخراج شده‌است. برای ترافیک جدید بهترین تعداد مراکز همان‌طور که در جدول نشان داده شده ۵ مرکز می‌باشد.

نشان داده شده است: همان‌طور که گفته شد هرچه این شاخص کم‌تر باشد، خوشه‌بندی‌های بهتری تولید شده است. طبق نمودار شکل (۸) کم‌ترین مقدار شاخص در خوشه‌بندی با ۵ مرکز به دست آمده است، یعنی بهترین تعداد خوشه برای ترافیک بات‌نت ۵ خوشه می‌باشد. P_i طبق رابطه (۸-۱۳) و S_i^* طبق رابطه (۸-۱۷) محاسبه شده‌است. نتایج در جدول (۹) نشان داده شده است.

جدول (۸). مقدار شاخص دیویس بولدین برای تعداد مراکز مختلف ترافیک بات‌نت

تعداد مراکز	شاخص دیویس بولدین (DB)
۲	۰/۴۸
۳	۰/۴۰
۴	۰/۴۳
۵	۰/۳۶
۶	۰/۴۶
۷	۰/۴۵
۸	۰/۴۷



تعداد خوشه‌ها

شکل (۸). نمودار مقدار شاخص دیویس بولدین بر حسب تعداد مراکز

جدول (۹). بهترین خوشه‌بندی برای ترافیک بات‌نت

مرکز خوشه (V_i)	جمعیت خوشه ($\ C_i\ $)	فاصله میانگین (S_i^*)	جمعیت نسبی (P_i)
۰/۰۵	۵۴۶۶۱	۰/۰۲۸	۰/۲۵
۰/۲۸	۵۷۹۵۱	۰/۰۵	۰/۲۶
۰/۹۴	۴۷۴۲۳	۰/۰۸۴	۰/۲۲
۰/۶۹	۲۵۶۰۶	۰/۰۶۴	۰/۱۲
۰/۹۸	۳۱۸۲۹	۰/۰۲۰	۰/۱۵

جدول (۱۲). فاصله مراکز ترافیک جدید نسبت به ترافیک بات‌نت

D ₁	۰/۰۲
D ₂	۰/۱۰
D ₃	۰/۰۵
D ₄	۰/۰۶
D ₅	۰/۰۹

پس مقدار T_A را برای ترافیک جدید نسبت به ترافیک بات‌نت، براساس داده‌های فوق محاسبه می‌کنیم:

$$T_A = (0.02 * 0.04) + (0.1 * 0.1) + (0.14 * 0.04) \quad (1-9)$$

$$+ (0.06 * 0.04) + (0.09 * 0.06) = 242 * 10^{-4}$$

با انجام همین اعمال روی ترافیک معمولی، مقدار T_A برای 644 * 10⁻⁴ نسبت به ترافیک معمولی به دست آمد. طبق آنچه قبلاً هم گفته شد، هرچه این مقدار کم‌تر باشد به معنی شبیه‌تر بودن دو ترافیک است. پس نتیجه می‌گیریم که ترافیک جدید باید از نوع بات‌نت باشد. و در واقع هم همین‌طور است. ترافیک تست‌شده، ترافیک یک بات‌نت است که در دانشگاه CTU گرفته شده است.^۱

همان‌طور که قبلاً گفته شد یکی از چالش‌هایی که در زمینه پژوهش پیرامون بات‌نت با آن مواجه هستیم داشتن مواد اولیه است که در اینجا مواد اولیه ما ترافیک بات‌نت واقعی، ترافیک ازدحام ناگهانی واقعی و ترافیک معمولی شبکه است. به دست آوردن ترافیک معمولی شبکه کار مشکلی نیست اما برای ترافیک بات‌نت از ترافیک‌هایی که توسط دانشگاه CTU ارائه شده استفاده کرده‌ایم و برای ترافیک ازدحام ناگهانی از ترافیک دانشگاه بزرگمهر در روز انتخاب واحد و در بازه‌های زمانی مختلف استفاده شده است. الگوریتم ارائه شده را با استفاده از ترافیک‌های موجود آزمایش کردیم که در ۹۲ درصد موارد، الگوریتم به درستی شناسایی کرد. همین کار را برای الگوریتم ارائه شده توسط شوئی یو و همکاران انجام دادیم که نتیجه آزمایش این بود که الگوریتم مذکور در ۶۳ درصد موارد، کار شناسایی را به درستی انجام داد. مشکل دیگری که در الگوریتم شوئی یو و همکاران وجود دارد، تعیین حد آستانه مناسب است. مقدار حد آستانه بستگی به تعداد جریان‌ها و اندازه ترافیک دارد. در هر حال تعیین حد آستانه، خود مسأله‌ای است که باید به دقت مورد بررسی قرار گیرد، چراکه اگر درست انتخاب نشود تمام نتایج تحت تأثیر این اشتباه قرار خواهند گرفت.

۱۰- ارزیابی

برای ارزیابی الگوریتم از روش 10 Folds Cross Validation که به طور گسترده در پژوهش‌ها به کار گرفته می‌شود، استفاده شده

است. این روش به منظور آموزش و تست، مجموعه داده را به ۱۰ زیرمجموعه تصادفی تقسیم‌بندی می‌کند. یک مجموعه برای تست و ۹ مجموعه دیگر برای آموزش به کار گرفته می‌شوند. فرایند آموزش و تست آن قدر تکرار می‌شود تا زمانی که تمام ۱۰ زیرمجموعه، یک‌بار به عنوان مجموعه تست و ۹ زیرمجموعه دیگر به عنوان مجموعه آموزشی، مورد استفاده قرار گیرند.

مزیت این روش در این است که تمام داده‌های موجود در مجموعه داده به عنوان داده آموزشی و تست، در نظر گرفته می‌شوند. مجموعه داده، به دو کلاس ترافیک بات‌نت نظیر به نظیر و ترافیک معمولی شبکه طبقه‌بندی می‌شوند، سپس عملکرد الگوریتم، مورد ارزیابی قرار می‌گیرد. پارامترهایی که برای ارزیابی مورد استفاده قرار می‌گیرند، پارامترهای F-Measure، Precision، Recall و Accuracy است که در ادامه شرح داده شده‌اند. جدول (۱۳) پارامترهای ارزیابی الگوریتم را نشان می‌دهد.

جدول (۱۳). پارامترهای ارزیابی الگوریتم

ترافیک معمولی شبکه	ترافیک بات‌نت نظیر به نظیر	کلاس ترافیک پارامتر ارزیابی
Recall _{Normal}	Recall _{Bot}	Recall
Precision _{Normal}	Precision _{Bot}	Precision
F_Measure _{Normal}	F-Measure _{Bot}	F_Measure
Accuracy		Accuracy

برای تعریف پارامترهای ارزیابی، ابتدا لازم است ماتریس درهم‌ریختگی معرفی شود. با این ماتریس می‌توان نتایج ارزیابی را بهتر نشان داد. ماتریس درهم‌ریختگی در جدول (۱۴) نشان داده شده است و پارامترهای آن به صورت زیر به دست می‌آیند:

- ۱- مثبت درست (TP): تعداد ترافیک‌های باتی که به عنوان ترافیک بات شناسایی شده‌اند.
- ۲- مثبت کاذب (FP): تعداد ترافیک‌های نرمالی که به عنوان ترافیک بات شناسایی شده‌اند.
- ۳- منفی درست (TN): تعداد ترافیک‌های نرمالی که به عنوان ترافیک نرمال شناسایی شده‌اند.
- ۴- منفی کاذب (FN): تعداد ترافیک‌های باتی که به عنوان ترافیک نرمال شناسایی شده‌اند.

جدول (۱۴). ماتریس درهم‌ریختگی

نرمال شناسایی شده	بات شناسایی شده	
منفی کاذب (FN)	مثبت درست (TP)	ترافیک‌های بات
منفی درست (TN)	مثبت کاذب (FP)	ترافیک‌های نرمال

- 2- True Positive
- 3- False Positive
- 4- True Negative
- 5- False Negative

$$\text{Precision}_{\text{Bot}} = \frac{TP}{TP+FP} \quad (3-10)$$

۴- Precision Normal: نشان‌دهنده این است که چند درصد ترافیک‌هایی که توسط الگوریتم به‌عنوان ترافیک نرمال شناسایی شده‌اند، واقعاً نرمال هستند. این پارامتر از نسبت تعداد ترافیک‌های نرمالی که توسط الگوریتم به‌درستی به‌عنوان ترافیک نرمال شناسایی شده‌اند، به تعداد کل ترافیک‌هایی که توسط الگوریتم به‌عنوان ترافیک نرمال شناسایی شده‌اند، به‌دست می‌آید و از طریق ذابطه (۴-۱۰) محاسبه می‌شود.

$$\text{Precision}_{\text{Normal}} = \frac{TN}{TN+FN} \quad (4-10)$$

۵- F_Measure: این پارامتر، ترکیبی از دو پارامتر Precision و Recall است و از طریق رابطه‌های (۵-۱۰) و (۶-۱۰) برای دو کلاس ترافیک بات‌نت نظیربه‌نظیر و ترافیک نرمال شبکه محاسبه می‌شود.

$$F_Measure_{\text{Bot}} = 2 * \frac{\text{Recall}_{\text{Bot}} * \text{Precision}_{\text{Bot}}}{\text{Recall}_{\text{Bot}} + \text{Precision}_{\text{Bot}}} \quad (5-10)$$

$$F_Measure_{\text{Normal}} = 2 * \frac{\text{Recall}_{\text{Normal}} * \text{Precision}_{\text{Normal}}}{\text{Recall}_{\text{Normal}} + \text{Precision}_{\text{Normal}}} \quad (6-10)$$

۶- Accuracy: این پارامتر نشان‌دهنده دقت کلی الگوریتم است و از نسبت تعداد کل ترافیک‌هایی که توسط الگوریتم به‌درستی شناسایی شده‌اند، به تعداد کل ترافیک‌ها به‌دست می‌آید. این پارامتر از طریق رابطه (۷-۱۰) محاسبه می‌شود.

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} \quad (7-10)$$

برای ارزیابی الگوریتم ارائه‌شده توسط روش ارزیابی فوق ابتدا مقادیر ماتریس درهم‌ریختگی که توسط الگوریتم به‌دست آمده است، مورد بررسی قرار می‌گیرد. سپس عملکرد الگوریتم در شناسایی ترافیک بات‌نت‌های نظیربه‌نظیر و ترافیک نرمال شبکه براساس پارامترهای ارزیابی، موردبررسی و ارزیابی قرار می‌گیرد.

ماتریس درهم‌ریختگی نشان‌دهنده نتایج ارزیابی الگوریتم است. جدول (۱۵) مقادیر ماتریس درهم‌ریختگی را برای الگوریتم

جدول (۱۵). مقادیر ماتریس درهم‌ریختگی برای الگوریتم

نرمال شناسایی شده	بات شناسایی شده	
۶	۲۴	ترافیک‌های بات
۲۷	۳	ترافیک‌های نرمال

جدول (۱۶) نتایج ارزیابی عملکرد الگوریتم را در شناسایی ترافیک بات‌نت نظیربه‌نظیر براساس پارامترهای تعریف‌شده، نشان می‌دهد.

ماتریس درهم‌ریختگی نشان می‌دهد که:

- هرچه پارامتر مثبت درست بیشتر باشد یا به‌عبارتی پارامتر منفی کاذب کم‌تر باشد، الگوریتم توانسته است تعداد ترافیک بات بیشتر تری را به‌درستی شناسایی کند.
 - هرچه پارامتر منفی درست بیشتر باشد یا به‌عبارتی پارامتر مثبت کاذب کم‌تر باشد، الگوریتم توانسته است تعداد ترافیک نرمال بیشتر تری را به‌درستی شناسایی کند.
 - باتوجه به این که در این پژوهش از روش 10 Folds Cross Validation برای آموزش و تست الگوریتم استفاده شده است، تمام داده‌های موجود در مجموعه داده، به‌عنوان داده تست در نظر گرفته می‌شوند، در نتیجه روابط زیر برقرار است:
 - مثبت درست به‌علاوه منفی کاذب نشان‌دهنده تعداد کل ترافیک بات‌نت است.
 - منفی درست به‌علاوه مثبت کاذب نشان‌دهنده تعداد کل ترافیک‌های نرمال است.
- باتوجه به ماتریس درهم‌ریختگی، پارامترهای ارزیابی به‌صورت زیر تعریف می‌شوند:

۱- Recall Bot: نشان‌دهنده این است که الگوریتم توانسته است چند درصد ترافیک بات را به‌درستی شناسایی کند. این پارامتر از نسبت تعداد ترافیک باتی که توسط الگوریتم به‌درستی به‌عنوان ترافیک بات شناسایی شده‌اند، به تعداد کل ترافیک‌های بات به‌دست می‌آید و از طریق رابطه (۱-۱۰) محاسبه می‌شود:

$$\text{Recall}_{\text{Bot}} = \frac{TP}{TP+FN} \quad (1-10)$$

۲- Recall Normal: نشان‌دهنده این است که الگوریتم توانسته است چند درصد ترافیک‌های نرمال را شناسایی کند. این پارامتر از نسبت تعداد ترافیک‌های نرمالی که توسط الگوریتم به‌درستی به‌عنوان ترافیک نرمال شناسایی شده‌اند، به تعداد کل ترافیک‌های نرمال به‌دست می‌آید و از طریق رابطه (۲-۱۰) محاسبه می‌شود:

$$\text{Recall}_{\text{Normal}} = \frac{TN}{TN+FP} \quad (2-10)$$

۳- Precision Bot: نشان‌دهنده این است که چند درصد ترافیک‌هایی که توسط الگوریتم به‌عنوان ترافیک بات شناسایی شده‌اند، واقعاً بات هستند. این پارامتر از نسبت تعداد ترافیک‌های باتی که توسط الگوریتم به‌درستی به‌عنوان ترافیک بات شناسایی شده‌اند، به تعداد کل ترافیک‌هایی که توسط الگوریتم به‌عنوان ترافیک بات شناسایی شده‌اند، به‌دست می‌آید و از طریق رابطه (۳-۱۰) محاسبه می‌شود.

جدول (۱۶). نتایج ارزیابی الگوریتم در شناسایی ترافیک‌های بات نت نظریه نظیر

۸۳/۸۰٪	F_Measure _{Bot}	۸۱٪	Precision _{Normal}
۸۵/۲۶٪	F_Measure _{Normal}	۸۸٪	Precision _{Bot}
۸۵٪	Accuracy	۹۰٪	Recall _{Normal}
		۸۰٪	Recall _{Bot}

نتایج ارزیابی نشان می‌دهد که عملکرد الگوریتم باتوجه به پارامترهای مختلف، متفاوت است. دو پارامتر Recall و Precision از اهمیت یکسانی برخوردار نیستند. Precision نشان‌دهنده این است که چند درصد ترافیک‌های باتی که الگوریتم به‌عنوان ترافیک بات شناسایی کرده است، واقعاً بات هستند. پارامتر Recall نیز مشخص می‌کند که الگوریتم توانسته است چند درصد ترافیک‌های بات را شناسایی کند. باتوجه به توضیحات زیر، پارامتر Recall از پارامتر Precision مهم‌تر است: اگر Recall کم باشد، یعنی الگوریتم توانسته است درصد کمی از ترافیک‌های بات را شناسایی کند. به‌عبارت دیگر، تعداد زیادی از ترافیک‌های بات به‌عنوان ترافیک نرمال شناسایی شده‌اند. این برای امنیت سیستم بسیار خطرناک است، زیرا الگوریتم نتوانسته است چنین ترافیک‌های باتی را تشخیص دهد. اگر Precision کم باشد، یعنی الگوریتم نتوانسته است درصد کمی از ترافیک‌های باتی را که واقعاً بات هستند، شناسایی کند. به‌عبارت دیگر، الگوریتم تعداد زیادی ترافیک نرمال را به‌عنوان ترافیک بات شناسایی کرده‌است. این برای امنیت سیستم خطر چندانی ایجاد نمی‌کند، فقط ممکن است موجب ناراحتی کاربر شود، زیرا ترافیک‌های نرمال به‌عنوان ترافیک بات شناسایی شده‌اند و ممکن است یک برنامه عادی، غیرفعال شود.

پارامتر Accuracy دقت کلی الگوریتم را در شناسایی ترافیک بات‌نت نظریه‌نظیر و ترافیک نرمال شبکه، نشان می‌دهد. این پارامتر مشخص می‌کند که الگوریتم توانسته است چند درصد ترافیک‌های بات‌نت نظریه‌نظیر و ترافیک‌های نرمال شبکه را به‌درستی شناسایی کند.

۱۱- نتیجه‌گیری

در این مقاله، ما سیستم تشخیص خود را که مبتنی بر تجزیه و تحلیل رفتار ترافیک و مبتنی بر جریان است، را ارائه دادیم. نکته‌ای که از آن برای تشخیص بات‌نت استفاده شد این است که جریان‌های حمله جاری در مقایسه با جریان‌های ازدحام ناگهانی، شباهت بیش‌تری به هم دارند؛ و طبق همین اصل، الگوریتمی ارائه‌شده که از ضریب هم‌بستگی جریان، به‌عنوان یک پارامتر شباهت برای جریان‌های مشکوک استفاده می‌شود. راه‌حل‌های دیگری برای مقابله با چنین حملاتی استفاده می‌شود که در آن از جورچین‌های گرافیکی برای تفاوت

قائل شدن بین انسان و بات‌ها استفاده می‌کنند. بدی این روش این است که استفاده از پاسخ‌های انسانی، می‌تواند برای کاربر آزاردهنده باشد. پژوهش ارائه‌شده مشکلات فوق را برطرف می‌کند. در این پژوهش، شناسایی بات‌نت نظریه‌نظیر در مرحله فرمان-کنترل از چرخه حیات، براساس تجزیه و تحلیل ترافیک انجام می‌شود؛ و از الگوریتم K-Means برای خوشه‌بندی ضرایب هم‌بستگی استفاده شده است و درنهایت براساس شباهت ترافیک جدید با جریان‌های بات‌نت می‌توان مخرب‌بودن ترافیک جدید را تشخیص داد. برای ارزیابی الگوریتم از روش 10 Folds Cross Validation که به‌طور گسترده در پژوهش‌ها به‌کار گرفته می‌شود، استفاده شده است. پارامتر Accuracy دقت کلی الگوریتم را در شناسایی ترافیک بات‌نت نظریه‌نظیر و ترافیک نرمال شبکه، نشان می‌دهد. این پارامتر مشخص می‌کند که الگوریتم توانسته است چند درصد ترافیک‌های بات‌نت نظریه‌نظیر و ترافیک‌های نرمال شبکه را به‌درستی شناسایی کند که مقدار این پارامتر برای الگوریتم ارائه شده ۸۵ درصد است. اما همان‌طور که گفته شد، داده‌های اولیه که برای انجام آزمایش مورد نیاز است، بسیار محدود است، بنابراین نرخ تشخیصی که از آزمایش این مجموعه داده به‌دست می‌آید نرخ قابل اطمینانی نیست. و برای دقیق‌تر شدن نتایج باید آزمایش را روی تعداد بیش‌تری از مجموعه داده انجام داد. که مطمئناً اگر مجموعه داده مورد آزمایش تنوع و تعداد بیش‌تری داشته باشد نرخ تشخیص نیز بالاتر خواهد رفت. اما با این حال، مزیتی که این الگوریتم نسبت به الگوریتم‌های مشابه دارد، تشخیص ترافیک بات‌نت از ترافیک ازدحام ناگهانی است. که این ویژگی خاص باعث می‌شود الگوریتم پیشنهادی نسبت به الگوریتم‌های مشابه اهمیت بیش‌تری داشته باشد.

۱۲- پیشنهادها

باتوجه به این‌که در شناسایی بات‌نت‌ها نرخ خطا در تشخیص وجود دارد، نیاز است برای تکمیل کار رویکردی به‌منظور شناسایی بات‌نت‌های نظریه‌نظیر در مرحله حمله از چرخه حیات ارائه شود. برخی بات‌نت‌ها ممکن است در مرحله فرماندهی و کنترل از چرخه حیات شناسایی نشوند، بنابراین نیاز است رویکردی وجود داشته باشد تا بتواند آن‌ها را در مرحله حمله از چرخه حیات شناسایی کند. باتوجه به این موضوع، در تحقیقات آینده قصد داریم رویکردی را برای شناسایی بات‌نت‌های نظریه‌نظیر در تمامی مراحل چرخه حیات ارائه دهیم. این رویکرد بایستی بتواند علاوه بر شناسایی بات‌نت‌های نظریه‌نظیر در مرحله فرماندهی و کنترل از چرخه حیات، آن‌ها را در مرحله حمله نیز شناسایی کند. امروزه استفاده از دستگاه‌های سیار مانند موبایل در میان

- and Security, 2009 ARES'09, International Conference on, 2009.
- [11] S. Garg, A. K. Sarje, and S. K. Pedd, "Improved Detection of P2P Botnets through Network Behavior Analysis," *Recent Trends in Computer Networks and Distributed Systems Security*, Springer Berlin Heidelberg, pp. 334-345, 2014.
- [12] J. Han, K. Micheline, and P. Jian, "Data mining," southeast asia edition: Concepts and techniques, Morgan kaufmann, 2006.
- [13] L. Xinying and W. Peizhi, "Data Mining Technology and its Application in Electronic Commerce," In *Wireless Communications, Networking and Mobile Computing*, 2008 WiCOM'08, 4th International Conference on, 2008.
- [14] J. R. Binkley and S. Singh, "An algorithm for anomaly-based botnet detection," *Proceedings of USENIX Steps to Reducing Unwanted Traffic on the Internet Workshop (SRUTI)*, pp. 43-48, 2006.
- [15] H. Husna, S. Phithakkitnukoon, and S. Pa, "Behavior analysis of spam botnets," *Communication Systems Software and Middleware and Workshops*, 2008 COMSWARE, 3rd International Conference on, 2008.
- [16] N. Pratik, S. Ray, C. Hota, and V. Venkatakrishnan, "Peershark: detecting peer-to-peer botnets by tracking conversations," In *Security and Privacy Workshops (SPW)*, 2014 IEEE, 2014.
- [17] D. Zhao, I. Traore, B. Sayed, W. Lu, S. Saad, A. Ghorbani, and D. Garant, "Botnet detection based on traffic behavior analysis and flow intervals," *Computers & Security*, vol. 39, pp. 2-16, 2013.
- [18] S. Sherif, I. Traore, A. Ghorbani, S. Bassam, D. Zhao, W. Lu, J. Felix, and P. Hakimian, "Detecting P2P botnets through network behavior analysis and machine learning," In *Privacy, Security and Trust (PST)*, 2011 Ninth Annual International Conference on, 2011.
- [19] W. Tarng, L.-Z. Den, K.-L. Ou, and M. Chen, "The analysis and identification of P2P botnet's traffic flows," *International Journal of Communication Networks and Information Security (IJCNIS)*, vol. 2, no. 3, 2011.
- [20] F. Chen, M. Wang, Y. Fu, and J. Zeng, "New detection of peer-to-peer controlled bots on the host," *Wireless Communications, Networking and Mobile Computing WiCom'09*, 5th International Conference, Beijing, 2009.
- [21] M. Stevanovic and J. M. Pedersen, "An efficient flow-based botnet detection using supervised machine learning," In *Computing, Networking and Communications (ICNC)*, 2014 International Conference on, 2014.
- [22] S. Yu, W. Zhou, W. Jia, S. Guo, Y. Xiang, and F. Tang, "Discriminating DDoS attacks from flash crowds using flow correlation coefficient," *Parallel and Distributed Systems*, IEEE Transactions, vol. 23, no. 6, 2012.
- [23] F. Kovács, C. Legány, and A. Babos, "Cluster validity measurement techniques," 6th International symposium of hungarian researchers on computational intelligence, 2005.

افراد افزایش یافته است. دستگاه‌های سیار قابلیت اتصال به اینترنت را دارند و افراد ممکن است بیش‌تر وقت خود را با استفاده از این دستگاه‌ها در اینترنت بگذرانند. همان‌طور که فن‌آوری در این زمینه پیشرفت کرده است و امروزه شاهد این هستیم که بیش‌تر افراد کارهای خود را با استفاده از گوشی‌های هوشمند انجام می‌دهند، مهاجمین و توسعه‌دهندگان بات‌نت از این فرصت استفاده کرده‌اند و فعالیت در محیط سیار را هدف قرار داده‌اند. بات‌نت‌هایی هستند که در محیط سیار فعالیت‌های مخرب خود را انجام می‌دهند. این بات‌نت‌ها به دستگاه‌های سیار افراد حمله می‌کنند تا به اهداف خود از قبیل سرقت اطلاعات مهم، دست یابند. باتوجه به این موضوع، نیاز است تا رویکردی برای شناسایی این بات‌نت‌ها نیز ارائه شود، بنابراین مسیر دیگر ادامه این پژوهش، کار بر روی شناسایی بات‌نت‌ها در محیط سیار است.

۱۳- مراجع

- [1] A. Cole, M. Michael, and D. Noyes, "Botnets: The rise of the machines," In *Proceedings on the 6th Annual Security Conference*, 2007.
- [2] B. Assadhan, M. José, and D. Lapsley, "Periodic Behavior in Botnet Command and Control Channels Traffic," *IEEE*, 2009.
- [3] H. Choi, H. Lee, and H. Kim, "BotGAD: detecting botnets by capturing group activities in network traffic," In *Proceedings of the Fourth International ICST Conference on communication System software and middleware*, 2009.
- [4] A. Karasaridis, B. Rexroad, and D. Hoeflin, "Wide-scale botnet detection and characterization," *Proceedings of the first conference on First Workshop on Hot Topics in Understanding Botnets*, vol. 7, 2007.
- [5] J. Govil and J. Govil, "Criminology of botnets and their detection and defense methods," *Electro/Information Technology*, 2007 IEEE International Conference, 2007.
- [6] S. S. Silva, R. M. Silva, and R. C. Pinto, "Botnets: A survey," *Computer Networks*, vol. 57, no. 2, pp. 372-403, 2013.
- [7] H. R. Zeidanloo, M. Safar, M. Zamani, P. Vahdani Amoli, and M. J. Z. Shooshtari, "A taxonomy of botnet detection techniques," *Computer Science and Information Technology (ICCSIT)*, 3rd IEEE International Conference on, vol. 2, 2010.
- [8] K.-S. Han and G. I. Eul, "A survey on p2p botnet detection," *Proceedings of the International Conference on IT Convergence and Security*, Springer Netherlands, pp. 589-593, 2012.
- [9] A. H. Lashkari and S. G. Ghalebani, "A Wide Survey on Botnet," *Digital Information and Communication Technology and Its Applications*, Springer Berlin Heidelberg, pp. 445-454, 2011.
- [10] J. Leonard, S. Xu, and R. Sandhu, "A framework for understanding botnets," In *Availability, Reliability*

Botnet Detection for Peer to Peer Networks

R. Azmi, M. Gholinezhad*, M. Saberi

* Alzahra University, Engineering College

(Received: 07/07/2015, Accepted: 12/01/2016)

ABSTRACT

Botnets are the latest types of internet-scale malware in recent years that has been the greatest threats to the web servers. Bot is an infected computer by a malware that are controlled remotely by one or more human factors without the user's knowledge. This controller agent called "bot master" and sometimes the infected system is called "victim". Peer to peer botnet is one type of botnets that use peer to peer protocols and detection of this type of botnet is more difficult than other types. Our suggested approach is a new method to detect such botnets. This approach uses network flow analysis and clustering method in data mining to detect peer to peer botnets. This approach is flow-based and compares the similarity between flows and K-Means clustering algorithm and eventually determines that the new traffic is an attack or not. This approach has good performance in detection of botnets in flash crowd traffic and this characteristic is distinction of the suggested algorithm and similar algorithms. Finally, the suggested approach has been tested with different traffic.

Keywords: Botnet, Botnet Detection, Bot, Peer to Peer Botnet, Network Flow Analysis.