

# بهره‌گیری از تکنیک‌های خوشه بندی داده‌کاو محور جهت افزایش تشخیص نفوذ در سیستم‌های اطلاعاتی ابری

فاطمه اسحق<sup>۱</sup>، مهدی خلیلی<sup>۲</sup>، مجید انجیدنی<sup>۳</sup>

دانشگاه پیام نور مرکز تهران شمال<sup>۱</sup>، دانشگاه پیام نور مرکز تهران شمال<sup>۲</sup>، دانشگاه پیام نور مرکز نیشابور<sup>۳</sup>

مسئول مکاتبات: فاطمه اسحق

## چکیده

با گسترش روز افزون تبادل اطلاعات و استفاده از سیستم‌های بر خط، میزان حملات و نفوذ در سیستم‌های اطلاعاتی افزایش یافته است. رایانش ابری گام بعدی تکامل سرویس های فناوری اطلاعات برحسب تقاضا می‌باشد. فایروال ها که توسط بسیاری از سازمان ها، به عنوان سیستم های تشخیص نفوذ به منظور حفاظت از امنیت سیستم های اطلاعاتی به کار گرفته می شوند اغلب در تشخیص حملاتی که از درون این سازمان ها اتفاق می افتد با شکست مواجه می‌شوند.

در این مقاله، به منظور غلبه بر این مشکل در فایروال‌ها، استفاده از تکنیک‌های خوشه‌بندی داده‌کاو محور شامل تکنیک‌های خوشه‌بندی Y-Means، K-Means و Fuzzy C-Means جهت برعهده گرفتن مسئولیت رسیدگی به نفوذ از درون سازمان‌ها، ارتقاء سیستم‌های تشخیص نفوذ و افزایش سطح امنیت اطلاعات در ابرها پیشنهاد می‌گردد تا نقشی حیاتی در تشخیص نفوذ با استفاده از تجزیه و تحلیل حجم بزرگ داده‌های شبکه و دسته‌بندی آن‌ها به صورت عادی و یا غیرعادی در فایروال‌ها بازی کند. همچنین با توجه به اینکه پیش‌تر تکنیک‌های داده‌کاوی با موفقیت برای تشخیص نفوذ در حوزه‌های کاربردی مختلف از جمله بیوانفورماتیک، بازار سهام، تجزیه و تحلیل وب و غیره مورد استفاده قرار گرفته است؛ از این روش استخراج روابط قبلی و ناشناخته در پایگاه داده‌های بزرگ، الگوبرداری نموده و سپس از الگوهای استخراج شده به عنوان پایه‌ای برای شناسایی حملات جدید استفاده خواهیم نمود.

## کلمات کلیدی

رایانش ابری- حملات امنیتی به ابر، سیستم‌های تشخیص نفوذ، داده‌کاوی، سیستم‌های تشخیص نفوذ داده‌کاو محور، الگوریتم K-Means، الگوریتم Y-Means، الگوریتم Fuzzy C-Means

## ۱- مقدمه

با پیشرفت فناوری اطلاعات نیاز به انجام کارهای محاسباتی در همه جا و همه زمان و همچنین نیاز به این که افراد بتوانند کارهای محاسباتی سنگین خود را بدون داشتن سخت‌افزارها و نرم‌افزارهای گران، از طریق خدماتی انجام دهند، به وجود آمده است که رایانش ابری آخرین پاسخ فناوری به این نیازها می‌باشد (آپارنا اس. وارد، ۲۰۰۹). با تکیه بر تکنولوژی ابر، میلیون‌ها نفر از کاربران داده‌های خود را در یک ابر که دارای فضای زیادی است ذخیره می‌کنند. ذخیره سازی در ابرها خطرات بسیاری نظیر دسترسی غیرمجاز، از دست دادن داده‌ها و غیره را به همراه دارد براساس بررسی های انجام شده در سال ۲۰۱۰ میلادی، امنیت به عنوان مهمترین چالش رایانش ابری شناخته شده است (دی. جی. براون، بی. سوکو و تی. وانگ، ۲۰۱۰). جدیدترین بررسی ها در سال ۲۰۱۱ حاکی از کاهش ۳۳/۵ درصدی چالش امنیت و رسیدن این درصد به عدد ۵۵ می‌باشد (جی. هویس منز، بی. بیسنز، دی. مارتینز، کی. دنیس و جی. وانسین، ۲۰۱۱).

حفظ حریم خصوصی داده‌ها عمده نگرانی افرادی است که از خدمات ابرهای عمومی استفاده می‌کنند. تهدیدهای امنیتی بر روی کاربران ابر به دو دسته داخلی و خارجی تقسیم بندی می‌شود. تهدیدهای خارجی شامل تهدید مراکز داده بزرگ می‌باشد که این نگرانی امنیتی در میان کاربران ابر و فراهم آورندگان (اشخاص ثالث) در حصول اطمینان از نرم افزارهای امن موجود، امکان پذیر است. علاوه بر مسائل مربوط به امنیت خارجی، ابر دارای برخی از مسائل مربوط به امنیت داخلی نیز می‌باشد که در آن کاربران باید در مقابل حملات از طرف یکدیگر، محافظت شوند. (انجمن کامپیوتر IEEE، ۲۰۱۰).

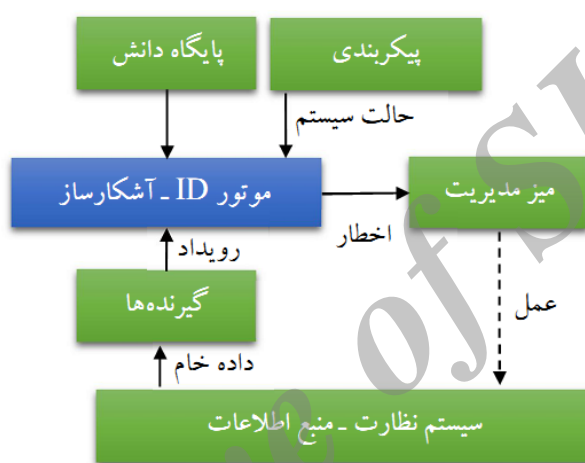
شکست در امنیت رایانش ابری به دلایل: الف- سخت افزاری و در لایه زیرساخت به عنوان سرویس ابر، ب- نفوذ کدهای مخرب در نرم‌افزار و در لایه نرم‌افزار کاربردی و ج- نفوذ کدهای مخرب در حال اجرا، توسط برنامه کاربردی کاربر یا تزریق اطلاعات ساختگی به برنامه توسط شخص ثالث؛ رخ می‌دهد که راه حل‌های: ۱- تنظیم تغییر رمز عبور کاربران به صورت اجباری در زمان‌های مشخص. ۲- پشتیبان‌گیری از داده‌ها در فواصل منظم ۳- تست نفوذ در سیستم در بازه‌های زمانی منظم ۴- نصب فایروال و ضدویروس‌های بروز شده ۵- استفاده از ناوگانی از سرورها که هر سرور جهت انجام کار خاصی در نظر گرفته شده برای رفع این مشکلات مورد استفاده قرار می‌گیرد (دیپسی کی. دانش و آیتا جان، ۲۰۱۳).

در این مقاله و به منظور غلبه بر این مشکل در فایروال‌ها، استفاده از تکنیک‌های خوشه بندی داده‌کاو محور شامل تکنیک‌های خوشه‌بندی K-Means، Y-Means، Fuzzy C-Means جهت برعهده گرفتن مسئولیت رسیدگی به نفوذ از درون سازمان‌ها، ارتقاء سیستم‌های تشخیص

نفوذ و افزایش سطح امنیت اطلاعات در ابرها پیشنهاد می گردد تا نقشی حیاتی در تشخیص نفوذ با استفاده از تجزیه و تحلیل حجم بزرگ داده های شبکه و دسته بندی آن ها به صورت عادی و یا غیرعادی در فایروال ها بازی کند. همچنین با توجه به اینکه پیش تر تکنیک های داده کاوی با موفقیت برای تشخیص نفوذ در حوزه های کاربردی مختلف از جمله بیوانفورماتیک، بازار سهام، تجزیه و تحلیل وب و غیره مورد استفاده قرار گرفته است؛ از این روش استخراج روابط قبلی و ناشناخته در پایگاه داده های بزرگ، الگوبرداری نموده و سپس از الگوهای استخراج شده به عنوان پایه ای برای شناسایی حملات جدید استفاده خواهیم نمود.

## ۲- نفوذ و سیستم تشخیص نفوذ

یک سیستم تشخیص نفوذ (IDS)، یک دستگاه یا برنامه نرم افزاری است که بر شبکه یا فعالیت های سیستم جهت فعالیت های مخرب یا نقض سیاست ها نظارت می کند و تولید گزارش ها برای یک ایستگاه مدیریت را بر عهده دارد که در شکل ۱ نشان داده شده است. (تشخیص نفوذ، ۱۹۹۳)



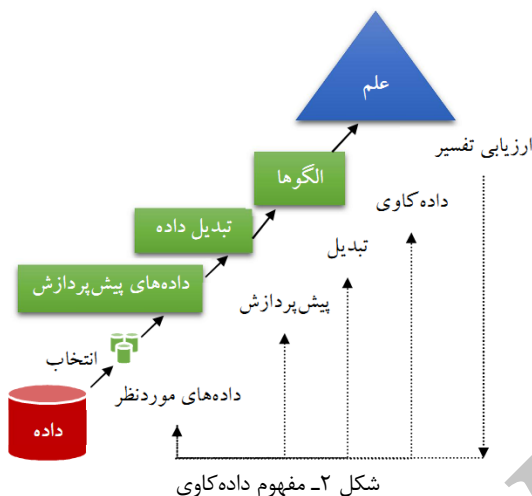
شکل ۱- مدل IDS

روش های متعددی در پیاده سازی یک ID براساس روش IDS (سیستم تشخیص نفوذ) وجود دارند:

- تشخیص ناهنجاری: این روش براساس تشخیص ناهنجاری ترافیکی است.
- سوءاستفاده/تشخیص امضاء: این روش الگوها و امضاهای حملات شناخته شده پیشین در ترافیک شبکه را جستجو می کند. سبک این روش با تشخیص نفوذ سر و کار دارد که شبیه به روش نرم افزار ضد ویروس عمل می کند.
- تشخیص ترکیبی: این روش ترکیبی از هر دو روش تشخیص ناهنجاری و تشخیص امضاء را بکار می برد.

## ۳- داده کاوی و سیستم تشخیص نفوذ

همانگونه که در شکل ۲ نشان داده شده است؛ داده کاوی، فرایند استخراج الگوها از داده ها می باشد. داده کاوی به عنوان یک ابزار مهم رو به افزایش برای کسب و کار مدرن که اطلاعات را به آگاهی شغلی تبدیل می کند، یک برتری اطلاعاتی را نمایش می دهد. در حال حاضر داده کاوی در طیف گسترده ای از شیوه های پروفایل، مانند بازاریابی، نظارت، تشخیص تقلب و اکتشاف علمی به کار می رود. (سی. یو. پی. مجموعه داده ها، ۲۰۱۰) دلیل اولیه برای استفاده از داده کاوی کمک به تجزیه و تحلیل مجموعه مشاهدات رفتار است. یک واقعیت غیرقابل اجتناب داده کاوی این است که تجزیه و تحلیل زیرمجموعه داده ها ممکن است نماینده کل دامنه و بنابراین شامل نمونه هایی از روابط بحرانی و رفتاری خاص که در دیگر بخش های دامنه وجود دارد، نباشد. (بان فاستر، یونگ ژائو، ایوان رایکو و شیانگ لو، ۲۰۰۸).

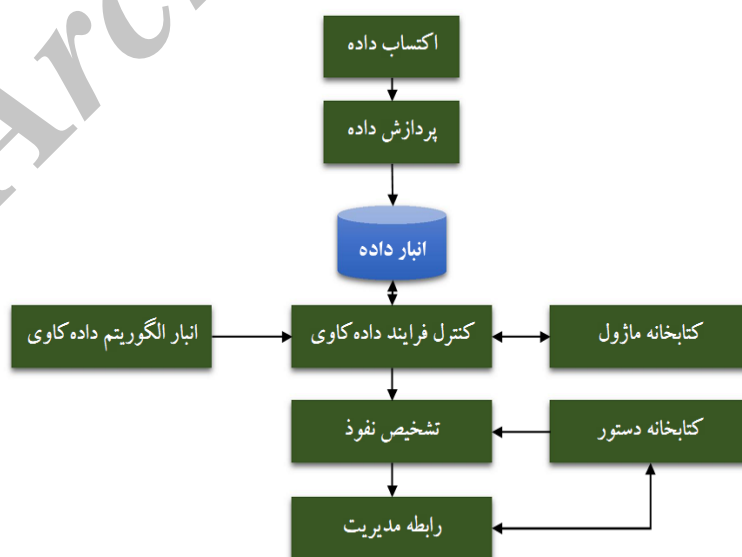


شکل ۲- مفهوم داده کاوی

چارچوب داده کاوی به طور خودکار الگوهای موجود در مجموعه داده‌های ما را تشخیص می‌دهد و از این الگوها برای پیدا کردن مجموعه‌ای از فایل‌های دودویی مخرب استفاده می‌کند. تکنیک‌های داده کاوی می‌تواند الگوهای مقادیر زیادی از داده‌ها را تشخیص دهد (مانند بایت کُد) و از این الگوها برای تشخیص داده‌های مشابه در نمونه‌های آینده استفاده کند. در سیستم تشخیص نفوذ، اطلاعات از منابع مختلف مانند: داده‌های میزبان، اطلاعات ورود به سیستم شبکه، پیام‌های هشدار و غیره به دست می‌آید. (داده کاوی"، کارگاه آموزشی بین المللی IEEE دوم در فناوری آموزش و پرورش و سیستم های تشخیص در پایگاه داده اوراکل ۱۰، ع، ۱۹۹۸)

تکنولوژی داده کاوی قابلیت استخراج پایگاه داده‌های بزرگ را دارد و از اهمیت زیادی برای استفاده در تشخیص نفوذ برخوردار است. با به کار بردن تکنولوژی داده کاوی، سیستم تشخیص نفوذ می‌تواند به طور گسترده‌ای داده‌ها را برای به دست آوردن یک مدل بررسی کند، بدین گونه برای مقایسه بین الگوی غیرطبیعی و الگوی رفتار طبیعی کمک می‌کند. تجزیه و تحلیل دستی برای این روش لازم نیست (اوما شانکار، کانیکا لاکانی و منیش موندرا، ۲۰۱۰).

یک مسأله مهم در تشخیص نفوذ این است که چگونه می‌توان به طور مؤثر الگوهای حمله و الگوهای طبیعی را از تعداد زیادی از داده‌های شبکه جدا کرد و چگونه می‌توان به طور مؤثر قوانین نفوذ خودکار را پس از جمع آوری داده‌های خام شبکه به وجود آورد. برای انجام این کار، از تکنیک‌های داده کاوی مختلفی مانند: طبقه‌بندی، دسته‌بندی و استخراج قانون رابطه استفاده می‌شود. در شکل ۳، یک سیستم که براساس الگوی تشخیص نفوذ طراحی گردیده است را نشان می‌دهد. (دیپسی کی. دانشز و آنیتا جان، ۲۰۱۲).



شکل ۳- داده کاوی براساس سیستم تشخیص نفوذ

داده کاوی براساس سیستم تشخیص نفوذ به دانش کمتر کارشناس نیاز دارد با این وجود اجرای خوب و امنیت را تأمین می کند. این سیستم ها به صورت حملات شناخته شده به اندازه حملات ناشناخته شبکه قابل تشخیص هستند. تکنیک های مختلف داده کاوی می تواند شبیه رده بندی، خوشه بندی و کاوش قواعد وابستگی، برای تحلیل ترافیک شبکه و در نتیجه تشخیص نفوذها استفاده شود. از میان موارد فوق، الگوریتم های خوشه بندی در بسیاری از جاها برای تشخیص نفوذ استفاده می شوند چون آنها به رده بندی داده آموزش به صورت دستی نیاز ندارند (میشل ام. جی، ۲۰۱۳).

اکثر نفوذها برای حمله به سیستم های مقصد از راه شبکه با استفاده از پروتکل های شبکه اتفاق می افتد. این نوع ارتباطات به عنوان ارتباطات غیرعادی (ناهنجار) و سایر ارتباطات، به عنوان ارتباطات عادی طبقه بندی می شوند. کلاً چهار نوع حمله به شرح ذیل وجود دارد (یه چینگ، وو شیائوپینگ و هوانگ گاوونگ، ۲۰۱۲):

- ۱) **DOS - عدم قبولی سرویس:** حمله کننده سعی می کند مانع دسترسی کاربران مجاز به خدمات در رایانه شود. برای مثال: فرمان مرگ، طغیان SYN و غیره
- ۲) **بررسی - نظارت و کاوش:** حمله کننده یک شبکه را برای کشف آسیب پذیری های شناخته شده رایانه آزمایش می کند. بررسی های این شبکه به طور معمول برای حمله کننده ای که یک حمله را در آینده برنامه ریزی می کند، بارز است. برای مثال: درگاه اسکن، فرمان پاکسازی و غیره
- ۳) **R2L - کنترل از راه دور به محلی:** حمله کننده های غیرمجاز، دسترسی محلی رایانه را از یک ماشین کنترل از راه دور به دست می آورند و سپس از آسیب پذیری های ماشین استفاده می کنند. برای مثال: حدس رمز ورود به سیستم.
- ۴) **U2R - کاربر با ریشه:** به ماشین اخیراً حمله شده، ولی حمله کننده تلاش می کند تا دسترسی با امتیازات کاربر بسیاری را به دست آورد. به عنوان مثال، حافظه موقت از حملات پُر می شود.

#### ۴- تکنیک های خوشه بندی داده کاومحور و سیستم تشخیص نفوذ

در این بخش، چند الگوریتم خوشه بندی برای تشخیص نفوذ استفاده شده است. همه این الگوریتم ها میزان خطای واقعی (مثبت های کاذب) را کاهش و میزان تشخیص نفوذ را افزایش می دهند. میزان تشخیص مانند تعداد نمونه های تشخیص نفوذ تعریف می شود تا سیستم تقسیم شده توسط جمع کل تعداد نمونه های تشخیص نفوذ در مجموعه داده را نشان دهد. میزان خطای واقعی به صورت تعداد کل نمونه های عادی که به اشتباه به صورت نفوذهای تعریف شده با تعداد کل نمونه های عادی، تعریف می شود. تعدادی از تکنیک های خوشه بندی نظیر الگوریتم های خوشه بندی داده کاو محور k-Means، Y-Means و Fuzzy C-Means در ذیل بحث می شوند.

#### ۴-۱ الگوریتم خوشه بندی k-Means

الگوریتم k-Means الگوریتمی است که به سختی قسمت بندی شده و در بسیاری جاها به سادگی و با سرعت به کار می رود. این الگوریتم فاصله اقلیدسی را با اندازه مشابه استفاده می کند. مفهوم خوشه بندی سخت این است که یک آیتم در یک مجموعه داده می تواند متعلق به یکی و فقط یک خوشه در یک زمان باشد. این الگوریتم، یک تحلیل خوشه بندی است به طوری که آیتم های گروه ها براساس مقادیر مشخص در داخل K به خوشه هایی که آیتم ها در همان خوشه ویژگی های مشابه دارند ملحق نمی شوند چون آن ها در خوشه های متفاوت، مشخصه های مختلفی دارند. عملکرد فاصله اقلیدسی با محاسبه فاصله بین دو آیتم ارائه شده به صورت:

$$d(p, q) = \sqrt{(p_1 - q_1)^2 + (q_2 - q_2)^2}$$

استفاده می شود.

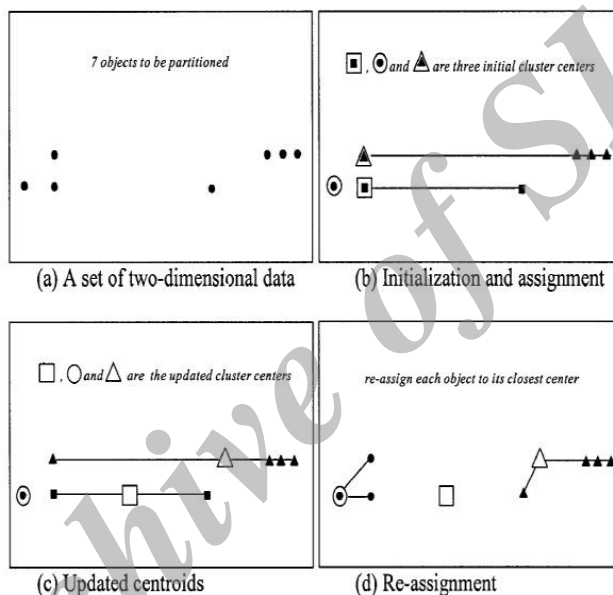
در اینجا؛  $P=(P_1, P_2, \dots, P_m)$  و  $q=(q_1, q_2, \dots, q_m)$  دو بردار ورودی با مشخصه های کمیت  $m$  هستند. الگوریتم اجرا می شود تا مجموعه های آموزش داده که ممکن است شامل ترافیک عادی و یا غیرعادی شود بدون اینکه قبلاً برچسب خورده باشد.

ایده اصلی این دیدگاه براساس فرضیه ای است که ترافیک عادی و غیرعادی از خوشه های مختلف می باشد. همچنین داده ممکن است شامل بخش های مجزایی شود که کدام آیتم های داده از آیتم های دیگر در خوشه خیلی متفاوت هستند و کدام متعلق به هیچ خوشه ای نیستند (جیبیان کی. پال، ۲۰۱۱). یک بخش مجزا با مقایسه شعاع آیتم های داده درک می شود به طوری که اگر شعاع یک آیتم داده بزرگتر از یک آستانه ارائه شده باشد، پس به صورت یک بخش مجزا دیده می شود. اما این موضوع، پردازش خوشه بندی k-means را از وقتی که تعداد بخش مجزا کوچک است، برهم نمی زند (ونکاتادری ام. لوکاناسا سی. ردی، ۲۰۱۳).

الگوریتم خوشه بندی k-means به صورت زیر است:

۱. تعریف تعداد خوشه های K. برای مثال، اگر  $K=2$  باشد ما در آموزش داده از دو خوشه متفاوت، ترافیک عادی و غیرعادی را فرض می کنیم.
۲. مقداردهی اولیه مراکز ثقل خوشه K. این کار با انتخاب اتفاقی آیتم های داده K از مجموعه داده انجام می شود.

۳. محاسبه فاصله هر آیتم تا مراکز ثقل همه خوشه‌ها با به‌کارگیری فاصله اقلیدسی متریک که برای یافتن شباهت بین آیتم‌ها در مجموعه داده به کار می‌رود.
  ۴. اختصاص هر آیتم با نزدیک‌ترین مرکز ثقل خوشه. در این روش همه آیتم‌ها به خوشه‌های مختلف اختصاص خواهند یافت، به طوری که هر خوشه آیتم‌هایی با ویژگی‌های مشابه خواهد داشت.
  ۵. پس از اختصاص یافتن همه آیتم‌ها به خوشه‌های مختلف، میانگین خوشه‌های تغییر یافته مجدداً محاسبه می‌شود. میانگین اخیراً محاسبه شده به عنوان مرکز ثقل جدید اختصاص می‌یابد.
  ۶. تکرار مرحله ۳ تا زمانی که مرکز ثقل خوشه تغییر نکند.
  ۷. برچسب زدن خوشه‌ها به صورت عادی و غیرعادی، به تعداد آیتم‌های داده هر خوشه بستگی دارد.
- یک مشکل اساسی روش خوشه‌بندی K-means، تعیین بخش اولیه و تعداد مناسب خوشه‌های K می‌باشد. همچنین گاهی اوقات منجر به هم‌ترازی می‌شود که میانگین کدام پردازش خوشه‌بندی ممکن است با چند خوشه خالی به پایان برسد. در شکل ۴، تولید خوشه‌های خالی را نشان داده شده است.



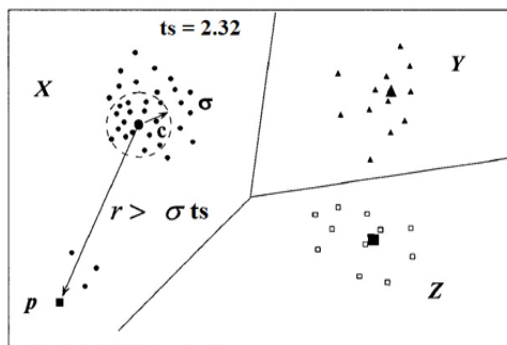
شکل ۴- تولید خوشه‌های خالی

## ۲-۴ الگوریتم خوشه بندی Y-Means

Y-Means، الگوریتم دیگری از خوشه‌بندی است که برای تشخیص نفوذ استفاده می‌شود. این تکنیک به طور خودکار یک مجموعه داده را درون تعداد قابل قبولی از خوشه‌ها قسمت‌بندی می‌کند به طوری که آیتم‌های داده درون خوشه‌های عادی و غیرعادی رده‌بندی می‌شوند. امتیاز اصلی الگوریتم خوشه‌بندی Y-Means آن است که به سه کمبود الگوریتم K-Means یعنی وابستگی مراکز ثقل اولیه، تعداد خوشه‌ها و هم‌ترازی برتری دارد. خوشه‌بندی Y-Means موانع خوشه‌های خالی را رفع می‌کند.

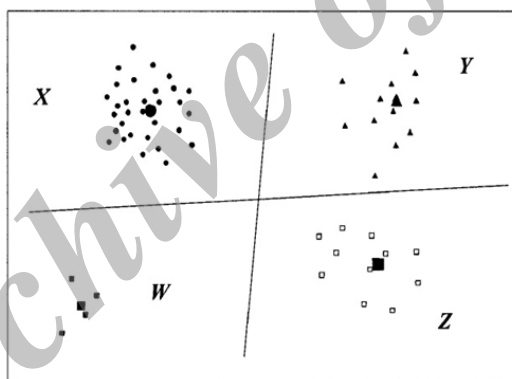
تفاوت اصلی بین خوشه‌بندی‌های Y-Means و K-Means تعداد خوشه‌های (k) در Y-Means است (یک خودتعریف متغیر به جای کاربرتعریف ثابت). اگر میزان K خیلی کوچک باشد، Y-Means تعداد خوشه‌ها را با چنددسته‌ای کردن آن‌ها افزایش می‌دهد. از طرف دیگر، اگر میزان K خیلی بزرگ باشد، Y-Means تعداد خوشه‌ها را با ادغام کردن در نزدیکی آن‌ها کاهش می‌دهد. Y-Means میزان مناسبی از K را به وسیله چنددسته‌ای کردن و ارتباط بین خوشه‌ها حتی بدون داشتن هیچ دانشی از آیتم توزیع شده تعیین می‌کند. این کار باعث می‌شود Y-Means یک تکنیک خوشه بندی موثر برای تشخیص نفوذ باشد چون ثبت وقایع شبکه داده به صورت اتفاقی توزیع شده و به دست آوردن میزان K به صورت دستی کار مشکلی است. Y-Means فاصله اقلیدسی را برای ارزیابی تشابه بین دو آیتم در مجموعه داده استفاده می‌کند. خوشه‌بندی سه مرحله اصلی دارد:

- **مرحله اول- تخصیص آیتم‌ها به خوشه‌های K:** وابستگی میزان K مشخص شده توسط کاربر، آیتم‌هایی را در یک مجموعه داده به نزدیک‌ترین خوشه‌های وابسته به فاصله بین آیتم و مرکز ثقل هر خوشه اختصاص می‌دهد.



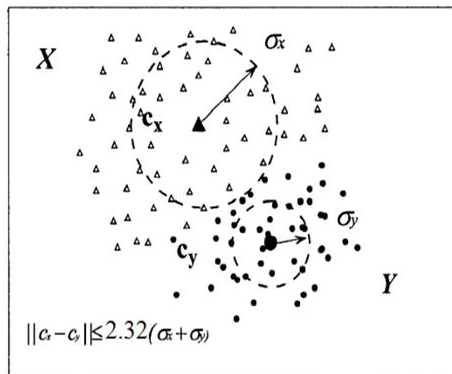
شکل ۵- تخصیص آیتم‌ها به خوشه‌های K

- **مرحله دوم - چنددسته‌ای کردن خوشه‌ها:** برطبق عملکرد توزیع عادی جمع‌ی استاندارد شده، ۹۹٪ نمونه‌های خوشه در داخل شعاع دایره  $\sigma$  قرار دارد جایی که  $\sigma$  انحراف استاندارد داده می‌باشد. بنابراین آستانه  $t = 2.32 \sigma$  انتخاب می‌شود. منطقه درون دایره، منطقه مطمئن خوشه نامیده می‌شود. بدین ترتیب، همه نقاط در خوشه که خارج از منطقه مطمئن قرار دارند، به صورت بخش‌های مجزا در نظر گرفته می‌شوند. سپس این بخش‌های مجزا از خوشه‌های فعلی دور شده و به مراکز ثقل جدید اختصاص می‌یابند. مراکز ثقل جدیداً شکل یافته ممکن است چند آیتم را از خوشه‌های مجاور و به موجب آن از خوشه‌های جدید جذب کنند. چنددسته‌ای کردن خوشه‌ها تا وقتی که هیچ بخش مجزایی وجود نداشته باشد ادامه پیدا می‌کند. چند دسته ای کردن، خوشه‌ها را به شاخه‌های کوچک‌تر تبدیل، تعداد خوشه‌ها را افزایش و سپس آیتم‌های درون همان خوشه را به یکدیگر شبیه‌تر می‌سازد.



شکل ۶- چنددسته‌ای کردن خوشه‌ها

- **مرحله سوم - پیوند بین خوشه‌ها:** وقتی دو خوشه نزدیک به هم یک اشتراک دارند، آنها می‌توانند درون یک خوشه بزرگ‌تر ادغام شوند. ادغام آستانه، مجموعه‌ای است با  $2.32 \sigma$  یعنی هر جایی که آن‌ها هستند هیچ آیتم داده‌ای در یک منطقه مطمئن نیست، همچنین اگر در منطقه مطمئن خوشه دیگری قرار گیرد، دو خوشه می‌توانند ادغام گردند. وقتی دو خوشه ادغام شدند مراکز ثقلشان کامل و سالم مانده و مرکز ثقل جدیدی ایجاد نمی‌شود یعنی خوشه جدید دو مرکز ثقل دارد. مزیت آن است که خوشه‌ها می‌توانند از شکل‌های دلخواه نظیر شکل‌های زنجیره مانند باشند.



شکل ۷- پیوند بین خوشه‌ها

الگوریتم خوشه‌بندی Y-means به صورت زیر است:

۱. تعریف تعداد خوشه‌های K
۲. مقداردهی اولیه مراکز ثقل خوشه K. این کار توسط انتخاب اتفاقی آیت‌های K از مجموعه داده انجام می‌شود.
۳. محاسبه فاصله هر آیت تا مراکز ثقل همه خوشه‌ها با استفاده از فاصله اقلیدسی متریک که برای یافتن شباهت بین آیت‌ها در مجموعه داده به کار می‌رود.
۴. اختصاص هر آیت به نزدیک‌ترین مرکز ثقل خوشه. در این روش همه آیت‌ها به خوشه‌های مختلف اختصاص خواهند یافت، به طوری که هر خوشه آیت‌هایی با ویژگی‌های مشابه خواهد داشت.
۵. پس از اختصاص یافتن همه آیت‌ها به خوشه‌های مختلف، میانگین خوشه‌های تغییر یافته مجدداً محاسبه می‌شود. میانگین اخیراً محاسبه شده به عنوان مرکز ثقل جدید اختصاص می‌یابد.
۶. کنترل کنید آیا هم ترازوی وجود دارد؟ اگر پاسخ بله است خوشه‌های خالی را حذف کرده و به مرحله ۷ بروید.
۷. اگر پاسخ خیر است، بخش‌های مجزا را پیدا کنید. اگر در هیچ یک از خوشه‌ها پیدا نشدند پس خوشه‌ها را چند دسته ای کنید.
۸. کنترل کنید آیا هم ترازوی وجود دارد؟ اگر پاسخ بله است خوشه‌های خالی را حذف کرده و به مرحله ۹ بروید.
۹. اگر پاسخ خیر است، خوشه‌های دارای اشتراک را لینک کنید.
۱۰. خوشه‌ها را به صورت عادی و غیرعادی برچسب بزنید.

#### ۴-۳ الگوریتم خوشه بندی Fuzzy C-Means

FCM یک الگوریتم بازبینی‌نشده خوشه بندی براساس مجموعه تئوری fuzzy است که اجازه می‌دهد یک عامل متعلق به بیش از یک خوشه باشد. درجه ارتباط هر آیت داده با خوشه، محاسبه شده و تصمیم می‌گیرد که کدام خوشه با کدام آیت داده وابسته شود. برای هر آیت، یک عامل مشترک داریم که مشخصه‌های درجه عضویت (uij) از بودن در خوشه Kth به صورت

$$u_{ij} = \sum_{k=1}^n (d_{ij}/d_{ik})^{(2/m-1)}$$

می‌باشد. در اینجا؛ d - فاصله آیت ith از خوشه jth، dik - فاصله آیت ith از خوشه kth و m - عامل نادقیق‌سازی.

وجود یک آیت داده در بیش از یک خوشه وابسته به مقدار m یعنی نادقیق‌سازی تعریف شده توسط کاربر در حد [0,1] می‌باشد که درجه نادقیقی را در خوشه تعیین می‌کند. بدین ترتیب، آیت‌های روی گوشه یک خوشه ممکن است در خوشه‌ای با درجه کمتر از آیت‌ها در مرکز خوشه باشد. وقتی m به مقدار ۱ برسد الگوریتم شبیه یک الگوریتم قسمت‌بندی حلقوی کار می‌کند و برای مقادیر بزرگ‌تر از m خوشه‌های متمایل به بیشتر بودن به اشتراک گذاشته می‌شوند.

هدف اصلی الگوریتم خوشه‌بندی fuzzy قسمت‌بندی داده درون خوشه‌هاست به طوری که تشابه آیت‌های داده درون هر خوشه افزایش یافته و تشابه آیت‌های داده در خوشه‌های مختلف کم می‌شود. به علاوه، این خوشه‌بندی، کیفیت قسمت‌بندی را طوری اندازه‌گیری می‌کند که یک مجموعه داده درون خوشه‌های C تقسیم شوند. الگوریتم FCM به کاهش میزان عملکرد عینی به شرح  $\sum_{i=1}^n \sum_{j=1}^m (u_{ij})^m \cdot ||x_i - v_j||^2$  توجه می‌کند. در اینجا؛ m - هر شماره واقعی بیش از ۱ باشد، uij - درجه ارتباط xi در خوشه j، xi - ith از داده اندازه‌گیری شده ابعدی، vj - فاصله مرکز ابعد از خوشه و  $||x_i - v_j||^2$  - هر مقیاس که بیان‌کننده تشابه بین هر داده اندازه‌گیری شده و مرکز است.



الگوریتم FCM به شرح زیر است:

۱. انتخاب به صورت اتفاقی مراکز خوشه  $c_{ij}$  -

۲. مقداردهی اولیه عضویت fuzzy ماتریکس  $u_{ij}$  با استفاده از:

$$u_{ij} = \sum_{k=1}^n (d_{ij}/d_{ik})^{(2/m-1)}$$

۳. محاسبه مراکز  $v_j$  در fuzzy با استفاده از:

$$v_j = \left( \sum_{i=1}^n (u_{ij})^m x_i \right) / \sum_{i=1}^n (u_{ij})^m, \quad \forall j = 1, 2, \dots, c$$

۴. به روز کردن ماتریس عضویت یعنی به مرحله ۲ بروید.

۵. اگر  $\|u^{(k+1)} - u^{(k)}\|$  کوچک تر از آستانه باشد پس پایان کار و اگر نباشد به مرحله ۳ برگردید. جایی که  $k$  مرحله تکراری است.

#### ۵- مقایسه الگوریتم های خوشه بندی داده کاومحور

این بخش مقایسه سه تکنیک خوشه بندی یعنی K-Means, Y-Means و Fuzzy C-Means را نشان می دهد. مقایسه با محاسبه معیار مختلف شبیه اجرا، کارایی، میزان تشخیص، میزان خطای واقعی، خالص بودن و ... انجام می شود. هر تکنیک تعدادی خصوصیت خوب برای رفع موانع سایر تکنیک ها دارد.

جدول ۱. مقایسه تکنیک های خوشه بندی های Fuzzy C-Means و Y-Means, K-Means

|   |   |   |                                  |
|---|---|---|----------------------------------|
| تعداد خوشه های $C$ در صورتی که $C$ کوچکتر از $m$ باشد.<br>مجموعه آیتم های داده $(X_1, X_2, \dots, X_m)$<br>مجموعه مراکز خوشه $(V_1, V_2, \dots, V_C)$ | تعداد خوشه های $Y$ در صورتی که $Y$ کوچکتر از $m$ باشد.<br>مجموعه آیتم های داده $(X_1, X_2, \dots, X_m)$ | تعداد خوشه های $K$ در صورتی که $K$ کوچکتر از $m$ باشد.<br>مجموعه آیتم های داده $(X_1, X_2, \dots, X_m)$ | ورودی                            |
| مجموعه ای از خوشه های $C$ ، جایی که هر خوشه آیتم های مشابه بیشتری دارد.   | مجموعه ای از خوشه های غیر خالی، جایی که هر خوشه آیتم های مشابه دارد.                                    | مجموعه خوشه های $K$ ، جایی که هر خوشه آیتم های مشابه دارد.  | خروجی                            |
| یک میزان عضویت به صورت "uiz" دارد.  | وجود ندارد.   | وجود ندارد.   | میزان عضویت                      |
| محاسبه چند فرمول را شامل می شود، به طوری که به زمان بیشتری نیاز دارد.   | چنددسته ای شدن و ارتباط خوشه ها را شامل می شود، به طوری که به زمان بیشتری نیاز دارد.                    | ساده و مستقیم به جلو به طوری که به زمان کمتری نیاز دارد.  | زمان محاسبه                      |
| بالا  | بالا  | پایین   | خالص بودن خوشه                   |
| خیر   | خیر   | ممکن است تولید داشته باشد شاید هم خیر   | تولید خوشه خالی                  |
| برای مجموعه داده های کوچک به اندازه مجموعه داده های بزرگ خوب کار می کند.  | برای مجموعه داده های کوچک به اندازه مجموعه داده های بزرگ خوب کار می کند.                                | برای مجموعه داده کوچک خوب کار می کند.   | کارایی/بازده                     |
| یک یا بیشتر از یک خوشه  | یک  | یک  | تعداد خوشه های وابسته به یک آیتم |



|                  |   |   |
|------------------|---|---|
| اجرای کامل       | به تعداد اولیه خوشه های K بستگی دارد.   | به تعداد اولیه خوشه های C بستگی دارد.   |
| شکل خوشه         | برای فشردن و کروی نمودن خوشه ها خوب کار می کند.   | هم برای کروی کردن و هم غیر کروی کردن خوشه ها خوب کار می کند.                                      |
| میزان تشخیص      | بالاترین  | بالا  |
| میزان خطای واقعی | پایین ترین  | پایین تر  |
| امتیازات         | ساده و سریع کارها برای فشردن و کروی نمودن خوشه ها خوب انجام می شود.   | هم ترازی نیست. به تعداد اولیه خوشه های K بستگی ندارد. با خوشه های کروی و غیر کروی خوب کار می کند. |
| معایب            | نیاز به تعیین تعداد مناسبی از خوشه ها دارد. هم ترازی با خوشه های غیر کروی خوب کار نمی کند. کیفیت خوشه ها اندازه گیری نمی شود. | اجرا به تعداد اولیه خوشه ها بستگی دارد. زمان بر بودن  |

## ۶- بحث و نتیجه گیری

همان طور که می دانید، روش های تشخیص سوءاستفاده برای شناسایی حملات ناشناخته کافی نیست لذا برای تشخیص نفوذ ناشناخته، ما باید به سمت تشخیص ناهنجاری برویم. تکنیک های مختلف داده کاوی شبیه رده بندی، خوشه بندی و کاوش قواعد وابستگی در تحلیل شبکه داده بسیار مفید است و از آن جایی که میزان وسیعی از ترافیک شبکه برای تشخیص نفوذ نیاز به جمع آوری دارد، پس خوشه بندی مناسب تر از رده بندی در دامنه تشخیص نفوذ می باشد، به طوری که جهت برچسب زدن مجموعه داده به صورت دستی، به کم کردن تلاش نیاز ندارد (کیزی، آنورهین، سوزان و وبسکی، ۲۰۰۹).

تکنیک های داده کاوی می تواند حملات شناخته شده را به خوبی حملات ناشناخته تشخیص دهند. فناوری داده کاوی به درک رفتار عادی در یک طرف داده و استفاده از این دانش برای تشخیص نفوذهای ناشناخته کمک می کند. در این مقاله، سه الگوریتم خوشه بندی داده کاومحور یعنی K-Means، Y-Means و Fuzzy C-Means بحث و بررسی شده اند. هر یک از آنها هم امتیازات و هم معایبی دارند و همدیگر را اصلاح می کنند:

الف- الگوریتم خوشه بندی K-Means در هم ترازی نتیجه می دهد و برای پایگاه های داده بزرگ مناسب نیست.

ب- الگوریتم خوشه بندی Y-Means انجام یک تغییر بر روی K-Means است که خوشه بندی های خالی را حذف می کند.

ج- الگوریتم خوشه بندی Fuzzy C-Means بر پایه یک منطق نامعلوم است که اجازه می دهد یک آیتم به بیش از یک خوشه بند وابسته باشد و روی به حداقل رساندن تابع هدف که کیفیت قسمت بندی را بازرسی می کند، تمرکز می کند. اجرا و کارایی خوشه بندی Fuzzy C-Means در دوره های تشخیص نفوذ بهتر از دو تکنیک دیگر است. در انتها، سه الگوریتم خوشه بندی را مقایسه می کنیم.

از میان خوشه بندی ها، Fuzzy C-Means را می توان یک الگوریتم مؤثر برای تشخیص فرض کرد چون اجازه می دهد یک آیتم متعلق به بیش از یک خوشه باشد و همچنین کیفیت قسمت بندی را اندازه می گیرد. تکنیک می تواند برای مجموعه داده بزرگ به خوبی مجموعه های داده که آیتم های اشتراکی دارند، استفاده شود. به علاوه، هیچ خوشه خالی را ایجاد نمی کند و در هنگام ایجاد خوشه بالاترین خلوص را دارد.

مزیت اصلی خوشه بندی Fuzzy C-Means برای تشخیص نفوذ، پیشنهاد میزان تشخیص بالا و میزان خطای واقعی پایین تر است. اگرچه Fuzzy C-Means یک تکنیک مؤثر است ولی زمان بر می باشد. اجرای سیستم های تشخیص نفوذ می تواند با ترکیب ویژگی های تکنیک خوشه بندی Fuzzy C-Means با چند تکنیک دیگر پیشرفت کند به طوری که زمان مورد نیاز با Fuzzy C-Means را برای پردازش خوشه بندی کاهش داده، میزان تشخیص را افزایش و همچنین میزان خطای واقعی را کم کند. در نتیجه سیستم تشخیص نفوذ را با دقت تر و موثرتر سازد.

ترکیب دسته بندی K-Means و درخت تصمیم C4.5 بهترین نرخ تشخیص (۹۹/۶) و کمترین نرخ خطای واقعی (۰/۱) را ارائه می دهد، اما حملات را به انواع مختلف طبقه بندی نمی کند. در حالی که آبشار K-Means و دو طبقه بندی KNN و ساده بیز؛ نرخ تشخیص ۹۸/۱۸ و نرخ خطای

واقعی ۰/۸۳۰ را ارائه می دهد و همچنین حملات مختلف را به صورت عادی، DOS، U2R، R2L و کاوشگر طبقه بندی می کند. داده کاوی روش مدرنی برای تشخیص نفوذ شبکه است. الگوریتم های داده کاوی ساخته شده - آماده شده، در دسترس هستند. مقدار زیادی از اطلاعات را می توان با فن آوری داده کاوی اداره کرد. این روش هنوز در حال توسعه است و می تواند به طور مؤثرتری به سرعت در حال رشد باشد. وظیفه اصلی ما، رسیدن به نرخ صحیح تر تشخیص نفوذ برای کاهش نرخ خطای واقعی است. داده کاوی هنوز هم در حال توسعه است بنابراین مطالعه و تحقیقات بیشتری باید انجام شود.

## ۷- منابع مورد استفاده

- Aparna S. Varde "Challenging research issues in data mining, databases and information retrieval" ACM SIGKDD Explorations Newsletter Volume 11 Issue 1, June 2009 Pages 49-52.
- CUP 99 Data Set", Proceedings of the 2009 IEEE Symposium on Computational Intelligence, pp 695-698, 2010.
- D. J. BROWN, B. SUCKOW and T. WANG, "A Survey of Intrusion Detection Systems", and SVM", 2010.
- Data Mining", IEEE Second International Workshop on Education Technology and Detection Systems in Oracle Database 10g", Proceedings of the Fourth International detection systems: The 1998 darpa off-line intrusion detection evaluation," discex, vol.02.
- Deepthy K Denatious and Anita John, —Survey on Data Mining Techniques to Enhance Intrusion Detection||, International Conference on Computer Communication and Informatics, Jan 2012.
- Denatious, D.K., and John, A. 2013. Survey on Data Mining Techniques to Enhance Intrusion Detection.
- Ian Foster, Yong Zhao, Ioan Raicu, Shiyong Lu. "Cloud Computing and Grid Computing 360-Degree", IEEE Conference, Date of Conference: 12-16 Nov. 2008.
- IEEE Computer Society, 2010 Sixth International Conference on Semantics, Knowledge and Grids, Security and Privacy in Cloud Computing: A Survey, 2010.
- INTRUSION DETECTION", National Research Council of Canada, NRC 45842, May 4-7 items in large databases", In Proceedings of the 1993 ACM SIGMOD Conference, New Machine and Decision Tree", International Journal of Computer Applications (0975 – 8887),
- J. HUYSMANS, B. BAESENS, D. MARTENS, K. DENYS and J. VANTHIENEN "New Trends in Data Mining" Tijdschrift voor Economie en Management Vol. L, 4, 2011.
- Jiban K Pal "Usefulness and applications of data mining in extracting information from different perspectives" Annals of Library and Information Studies Vol. 58, March 2011, pp. 7-16.
- Kazi, Aunnurhain, Susun, Vrbsky, Security Attacks and Solutions in Clouds, University of Alabama, 2009.
- Michael, M. J., A view Of Cloud Computing, Communications of the ACM, April 2013.
- Uma Somani, Kanika Lakhani, Manish Mundra, Implementing Digital Signature with RSA Encryption to Enhance Data Security of Cloud in Cloud Computing, IEEE, 2010.
- Venkatadri. M, Lokanatha C. Reddy "A Review on Data mining from Past to the Future" International Journal of Computer Applications (0975 – 8887) Volume 15– No.7, February 2013-19.
- Ye Qing, Wu Xiaoping and Huang Gaofeng. An Intrusion Detection Approach based on Data Mining||, 2nd International Conference on Future Computer and Communication, pp. No. 695 – 698, 2012 IEEE.