



## کاربرد داده کاوی جهت کشف رفتارهای مشکوک در تراکنش های کارت های بانکی

لیلا جدیدی،<sup>۱</sup> پروین پهاران

۱- عضو هیأت علمی دانشگاه آزاد اهواز مؤسسه فرهنگی آموزشی سماء خوزستان، اهواز، ایران

۲- دانشجوی کارشناسی نرم افزار مؤسسه غیرانتفاعی جهاد دانشگاهی خوزستان، اهواز، ایران

Leila.jadidi@yahoo.com

### چکیده

تکنیک های داده کاوی به صورت گسترده ای در دنیا در حوزه کشف تقلب مورد استفاده قرار گرفته شده است، دلیل این موضوع اطباق زیاد این فرآیند با نیازمندی های مساله کشف تقلب می باشد. در این مقاله سعی بر پیاده سازی قدم به قدم فرآیند داده کاوی بر داده های مربوط به تراکنش های کارت بانکی فردی خاص، جهت کشف رفتار مشکوک در این تراکنش ها شده است. در مرحله مدل سازی این فرآیند از تکنیک تحلیل خوشجهت آموزش مدل هایی با پارامترها و تنظیمات متنوع (برای مثال تعداد خوشنه های مختلف، تعداد نرون های مختلف...) بر داده های بدون برقسپ، استفاده شده است. تکنیک خوشجهت بندی انتخاب گردیده در این مقاله نگاشت خود سازمانده یا همان شبکه عصبی کوهنن می باشد، بعد از خوشجهت بندی داده ها از مدل طبقه بندی، داده ها از مدل طبقه بندی، جهت ارزیابی مدل های غیر نظراتی و همچنین انتخاب بهینه ترین مجموعه مدل ترکیبی (مدل های نظراتی و غیر نظراتی) جهت کشف رفتار مشکوک تراکنش های آتنی، استفاده شده است. تکنیک طبقه بندی انتخاب شده، مدل مجموعه قوانین با استفاده از الگوریتم C5 می باشد. در انتهای تحقیق نیز، قواعد مدل طبقه بندی بهینه، جهت استخراج دانش هایی مورد تفسیر قرار گرفته است.

واژه های کلیدی: داده کاوی، رفتار مشکوک، مدل های غیر نظراتی، تحلیل خوشجهت، شبکه عصبی کوهنن

### مقدمه

افزایش تنوع کارت های اعتباری و بانکی ارائه شده از سوی بانک ها و موسسات مالی و اختصاص سهم بالای از حجم مبادلات مالی افراد حقیقی و حقوقی در دنیا توسط این کارت ها، احتمال ایجاد تخلف و سوء استفاده از این کارت ها و ایجاد مبالغه بالای ضرر و زیان نیز برای مشتریان و موسسات مالی رو به افزایش می دهد، برای مثال لونوارد هزینه های مربوط به تقلب برای کارت های ویزا و مستر کارت در کانادا را به ترتیب برای سال های ۱۹۸۹ و ۱۹۹۰ و ۱۹۹۱ و ۴۶ میلیون دلار برآورد نموده است. همچنین منبع آماری ایکس روند رشد تقلب کارت های اعتباری در کشور انگلستان در سال های ۱۱۹۷ و ۲۰۰۰ را به ترتیب برای مبالغه ۱۲۲ و ۱۳۵ و ۱۸۸ و ۲۹۳ میلیون یورو برآورد نموده است.[۱]

بیشتر مسائل کشف تقلب در حوزه کارت های بانک، شامل حجم عظیمی از مجموعه داده می باشد، برای مثال شرکت کارت های اعتباری بارکلای کارد تقریباً به تنها بیان دارای ۳۵۰ میلیون تراکنش در سال در کشور انگلستان می باشد. به دلیل حجم بالای داده در پایگاه اطلاعاتی مربوط به تراکنش های کارت های بانکی و همچنین وجود اطلاعاتی بصورت عددی و اسمی (سمبولیک) یا ترکیبی از این دو، پرداز بر روی این مجموعه داده ها برای جستجو تراکنش های متقابلانه، نیاز به استفاده از الگوریتم های کارا و سریع دارد. با معرفی و رشد دانش داده کاوی در دنیا و همچنین برآورده نمودن نیازمندی های مساله کشف تقلب، یعنی امکان استفاده از این ابزار بر روی حجم بالای اطلاعات و کشف الگوهای پنهان در این داده ها بصورت خودکار و سریع، شاهد توسعه بکار گیری این رویکرد در حل مسائل کشف تقلب در حوزه کارت های بانکی در سال های اخیر در دنیا بوده ایم.

تقسیم بندی مدل های کشف تقلب در حوزه کارت های بانکی

استراتژی های کلان مسائل کشف تقلب در حوزه کارت های بانکی را نیز می توان، منطبق با استراتژی های فرآیند داده کاوی دانست.

دو استراتژی کلان برای فرآیند داده کاوی قابل تصور می باشد که عبارتند از:

یادگیری نظراتی (هدايت شده)

یادگیری غیر نظراتی (غیر هدايتی)

انتخاب یکی از این دو رویکرد و استفاده از مدل ها و روش های قابل استفاده مرتبط با آن ها نیز بستگی به ماهیت اطلاعات و داده های تاریخی موجود دارد. زمانی که داده های موجود بصورت برقسپ خورده باشند می توان از رویکرد نظراتی استفاده نمود و در غیر این صورت رویکرد غیر نظراتی (هدايت نشده) رویکرد مناسب تری جهت این موضوع می باشد. تکنیک های مدل های مورد استفاده در خصوص مساله کشف تقلب در سال ۲۰۰۲ توسط بلتن و هند بررسی موری گردید، آن ها مدل های کشف تقلب در حوزه کارت های اعتباری را تحت دو رویکرد کلی نظراتی و غیر نظراتی طبقه بندی نمودند.[۱] همچنین فوالي، اسمیت و گایلر نیز در سال ۲۰۰۵ تکنیک های مرتبط با حوزه کارت های اعتباری را در چهار گروه ۱- گروه های نظراتی ۲ - غیر نظراتی ۳ - ترکیب تکنیک های نظراتی ۴ - ترکیب تکنیک های نظراتی تقسیم بندی نمودند.[۲]

### بیان مساله و رویکرد اصلی تحقیق

امکان انجام پروژه های داده کاوی مستلزم استفاده از داده های مرتبط در خصوص مساله می باشد، به همین منظور داده های تراکنش های کارت های بانکی به عنوان اصلی ترین ورودی مدل های کشف تقلب، از یکی از بانک های موجود در کشور برای انجام این تحقیق جمع آوری شده است و از آنجا که امکان دسته بندی داده های مربوط به تراکنش های بیگر بر چسب زنی داده ها در بیشتر بانک ها و موسسات مالی و از جمله بانک مورد نظر، ز متدالوگی های داده کاوی به ویژه استفاده از مدل های غیر نظراتی، بر داده

# همایش ملی الکترونیکی دستاوردهای نوین در علوم مهندسی و پایه

National e-Conference on Advances in Basic Sciences and Engineering

WWW.AEBSCONF.IR



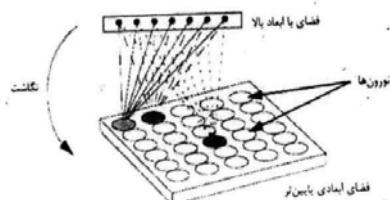
های تراکنش‌های کارت‌های بانکی مربوط به فردی خاص در این بانک می‌باشد، این امر به جهت طراحی سرویسی ویژه در بانک مورد نظر برای مشتریان بانک در خصوص امکان ایجاد کنترل و نظارت بر تراکنش‌های صورت گرفته از طریق کارت‌های متعلق به آن‌ها می‌باشد.

## مرور ادبیات مرتبط با تحقیق

تکنیک مورد استفاده در این تحقیق نگاشت خود سازمانده (شبکه عصبی کوهن) جهت تحلیل خوش داده‌ها با رویکرد غیرنظری می‌باشد. در سال ۲۰۰۶ زاسلاویسکی و استریزاک اقدام به راهه مدلی برگرفته از الگوریتم نگاشت خود سازمانده جهت تشخیص تراکنش‌های مشکوک کارت‌های اعتباری با رویکرد غیرنظری، برای سیستم پرداخت بانک‌ها نمودند.<sup>[۳]</sup> همچنین کواه و اسریگانش در سال ۲۰۰۸ با توسعه رویکرد قبلی و با تاکید بر حل مساله کشف تقلب در زمان حقیقی با حداقل هزینه و زمان، اقدام به طراحی مدلی ب سه لایه نمودند. لایه اصلی و هسته‌ای ابتدا با استفاده از نگاشت خود سازمانده اقدام به غربالگری و دسته بندی اولیه داده‌ها به دو گروه تراکنش‌های مشکوک قانونی نموده و سپس با بکارگیری شبکه عصبی پیش‌خور با الگوریتم پس انتشار، اقدام به تحلیل تراکنش‌های مشکوک بر اساس الگوهای تقلب می‌پردازد و تقلب تراکنش را در نهایت تایید یا رد می‌نماید.<sup>[۴]</sup>

## شبکه عصبی کوهن (نگاهداشت خود سازمانده SOM) [۵]

نگاشت‌های خود سازمانده نوع دیگری از شبکه‌های عصبی رقابتی بدون ناظر است. نگاشت‌های خود سازمانده (SOM) مدلی غیرنظری بر شبکه‌های عصبی است که توسط کوهن نیز مانند سایر شبکه‌های عصبی دارای ارتباطات کامل و بصورت پیش‌خور می‌باشد و پیشخور بودن آن به این معنی که این گونه از شبکه‌ها اجزاء ایجاد حلقه را ندارد. همچنین ارتباط کامل شبکه کوهن نیز به این معنی که بین هر یک از عصب‌های (گره یا نرون) هر لایه و هریک از عصب‌های لایه بعد یک ارتباط وجود دارد به هر ارتباطی بین دو لایه عددی، وزنی که مقدار آن بین صفر و یک می‌باشد، بصورت تصادفی تخصیص داده می‌شود، که در واقع تعديل این وزنها همان مکانیسم یادگیری شبکه‌های عصبی می‌باشد. اما علی‌رغم سایر انواع شبکه‌های عصبی این نوع شبکه دارای لایه پنهان نمی‌باشد و داده‌ها بصورت مستقیم از لایه ورودی به لایه خروجی ارائه می‌شود. لایه خروجی معمولاً به شکل یک یا دو می‌باشد، که معمولاً بصورت مستطیل شکل بوده و گاه‌ها شش گوش نیز می‌باشند. شبکه کوهن می‌تواند برای به تصویر کشیدن داده‌های ابعادی بال به کار گرفته شود. این نگاشت سطوحی را که در فضای ورودی به هم نزدیک‌اند به فضایی که در آن داده‌ها در فضای خروجی نزدیک‌اند، تصویر می‌کند. شکل (۱) نشان می‌دهد که چگونه یک بردار ورودی هفت بعدی به یک فضای اصلاح شده دو بعدی تصویر شده است.



شکل ۱: کاهش ابعادی در نگاشت خود سازمانده

اگرچه شبکه کوهن یکی از روش‌های کاهش بعد نیز به حساب می‌آید اما کاربرد اصلی این شبکه در تحلیل خوش‌خواه (خوش بندی) می‌باشد. بطوری که این شبکه قابلیت تفکیک داده‌های ورودی را به گروه‌هایی جدا از هم، که بر اساس شباهت‌ها و الگوهای پنهان بین این داده‌ها تقسیم‌بندی می‌شوند را دارا می‌باشد. بردار ورودی کوهن را بصورت مجموعه ای از مقادیر  $M$  فیلد برای  $N$  این رکورد بصورت زیر در نظر می‌گیریم  $\dots, \dots, xn1, xn2, \dots, xnM = w1, w2, \dots, wM$  در نظر گرفته می‌شود.

الگوریتم شبکه کوهن طی مراحل زیر صورت می‌پذیرد.

باایستی برای هر بردار ورودی  $X$  مراحل زیر را انجام می‌دهیم:

رقبات: برای هر گره خروجی  $J$  ام مقداری برای  $D(wj, Xn)$  به عنوانتابع ارزیابی محاسبه می‌کنیم برای فرض فاصله اقلیدسی  $D$  بصورت زیر بدست می‌آید.

$$D(wj, Xn) = \sqrt{\sum_i w_{ji}^2}$$

گره ای بعنوان گره برنده معرفی می‌گردد که مینیمم فاصله را از بین گره‌ها ورودی از گره  $J$  ام داشته باشد.

همکاری: در این مرحله اقدام به شناسایی همه خروجی گره  $J$  ام که در همسایگی تعریف شده از گره  $J$  ام داشته باشد، اقدام می‌نماییم و سپس برای همه این گره‌ها عملیات بالا را تکرار می‌کنیم تا تعدیلی در اوزان به وجود آید.

Anatipac: تنظیم نرخ یادگیری و اندازه همسایگی در صورت نیاز

هنگامی که شرایط خاتمه مهیا گردد توقف می‌نماییم.

# همایش ملی الکترونیکی دستاوردهای نوین در علوم مهندسی و پایه

National e-Conference on Advances in Basic Sciences and Engineering



WWW.AEBSCONF.IR

## جمع آوری اولیه داده

داده های اولیه جهت بررسی و پیاده سازی فرایند داده کاوی مربوط به تراکنشهای صادره بانک ملت می باشد، بعد از بررسی های انجام شده در بانک مورد نظر، داده های مربوط به کارت اکسیس کارت این بانک مرود انتخاب قرار گرفت، دلیل انتخاب این کارت بیان این موضوع می باشد، که به هر دارنده صاحب کارت یک شماره حساب بانکی نیز تخصیص داده شده و همچنین این کارتها قابلیت انجام تراکنش با مبالغ بالا و از طریق مختلف (اینترنت، خودپرداز و پوز) را دارا می باشد، برای مثال بعضی از کارتاهای بانکی سیستم شتاب دارای این ویژگی ها نبوده و سقف تراکنش های انجام شده روزانه برای آنها حداقل مبلغ ۱۲۰۰۰۰ تومان می باشد. به علت حجم بالای کلیه داده های تراکنشهای ذخیره شده در پایگاه اطلاعاتی بانک و همچنین تصمیم مدیران ارشد بانک در خصوص ارائه این سرویس ویژه برای افرادی خاص که دارای درجه اهمیت بالاتری نسبت به سایر دارندگان کارتاهای بانکی می باشند، داده های مربوط به ۱۰۰ نفر اول از لحاظ حجم میادلات ریالی صورت گرفته از طریق کارت های آنها در ایران، جمع آوری گردیده شد.

## توصیف داده

بعد از بارگذاری داده ها در نرم افزار بانک اطلاعاتی تعداد تراکنشهای ثبت گردیده (تعداد رکوردها) برای ۱۰۰ مشتری انتخاب شده حدود ۷۴۰۰۰ رکورد مشخص گردید همچنین فیلدهای اطلاعاتی موجود در پایگاه داده ابتدایی، بصورت فیلدهای نام برده شده در جدول شماره ۱ مشخص گردیده شد.

جدول ۱: مشخصات پایگاه داده تراکنش ها

نام ستون (Column Name)	نوع داده (Data Type)	طول رشته (String length)	تشریح فیلد (Description)
PAN	Varchar	۲۰	شماره کارت
TRANSNO	Varchar	۲۰	شماره تراکنش
TRANSDATE	Numeric	Integer	تاریخ تراکنش
DEVICENO	Numeric	Integer	کد دستگاه
DEVICETYPE	TEXT	۵	نوع تراکنش
AMOUNT	Numeric	Long Integer	مقدار تراکنش
RESULTTRANS	TEXT	۵	نتیجه تراکنش

## اکتساف داده

بعد از بارگزاری داده ها در نرم افزار داده کاوی و بررسی و تحلیل داده ها در قالب گرافهای توزیع مختلف بر اساس فیلدهای اطلاعاتی، متوجه وجود ارتباط معنی داری بین مبلغ تراکنش و دستگاه مورد استفاده و نوع تراکنش انجام شده، شدیم، به اینصورت که مبالغی که به تعداد بالا تکرار شده بودند و یا دارای مقادیر بالایی بودند، بیشتر از طریق پوز و اینترنت انجام گردیده شده اند.

## انتخاب داده

با توجه به اینکه فرض اولیه مساله در نظر گرفتن داده های یک فرد خاص به عنوان رفتاری قانونی بدون داشتن آگاهی قبلی نسبت به احتمال وقوع تخلف در تراکنشهای او می باشد، در این مرحله اقدام به انتخاب یک نفر با رفتاری نرمال تر و با توزیع متنوع تر نسبت به بقیه از بین صد نفر، شد، به این معنی که حدالملکان در تراکنش های آن قدر خاص، سه نوع تراکنش های اینترنت و پوز و دستگاه های خود پرداز و همچنین تراکنش رد و قبول نیز به نسبتهای نزدیک به هم، موجود باشد. در نهایت فردی دارای حدوداً ۳۲۰ تراکنش انجام شده در بازه زمانی دو سال انتخاب شد.

## پاکسازی داده ها

بعد از بررسی های صورت گرفته بر روی پایگاه اطلاعاتی موجود، به دلیل فرایند ذخیره سازی کاملاً خودکار داده ها در بانک، رکوردهایی با مقادیر گمشده ملاحظه نگردید. اما دو رکورد به عنوان نویز با مقادیر اشتباه برای داده ها مشخص گردید (کد دستگاه ۱۱۷۱۱ اشتباها در دو تراکنش ۱۱۷۱۱ و کد دستگاه ۱۰۰۰۱۰۰۰ ثبت شده بود) که بصورت دستی اقدام به اصلاح و یا به تعبیر دیگر هموارسازی داده ها نمودیم. از طرفی ۷ رکورد نیز به عنوان داده های پرت شناسایی شد که به جهت احتمال تشخیص آنها توسط مدل مورد نظر، جهت کشف رفتار مشکوک، اقدامی برای خروج آنها از داده ها صورت نپذیرفت.

## ساخت داده

در این مرحله خصوص عدم قابلیت استفاده اگوریتم خوش بندی شبکه عصبی کوهنن از فیلدهای اطلاعاتی کد دستگاه به دلیل حروفی بودن نوع قابل بندی و تنوع زیاد مقادیر این متغیر ورودی و همچنین فیلد اسطلایاتی تاریخ انجام تراکنش به دلیل تنوع زیاد آن و قالبندی عددی برای آن در پایگاه داده موجود، با استفاده از نرم افزار بانک

نمد.

# همایش ملی الکترونیکی دستاوردهای نوین در علوم مهندسی و پایه

National e-Conference on Advances in Basic Sciences and Engineering

WWW.AEBSCONF.IR



مقادیر شمارش شده از تکرار هر دستگاه مشخص در کل رکوردها به جای عدد کد دستگاه در همه رکوردهای پایگاه داده، جایگزین شده و فیلد اطلاعاتی جدیدی با نام Count-Deviceno ایجاد گردید.

هر تراکنش بر اساس تاریخ انجام تراکنش در یکی از شش گروه تاریخی تراکنشهای پنج روزه اول تا پنج روزه ششم ماه طبقه بندی شده و مقادیر مربوط به آن در فیلد اطلاعاتی جدیدی با نام fiveday قرار گرفت.

تقسیم بندی داده ها

بعد از انتخاب و آماده سازی داده ها، داده ها بصورت تصادفی به سه مجموعه داده مجزا با نام آموزشی، آزمون و اعتبار سنجی تقسیم بندی شدند. کاربرد این مجموعه داده ها را می توان بترتیب جهت ساخت اولیه یک مجموعه مدل، مشخص نمودن عملکرد مدل ساخته شده بروی داده ها دیده نشده و انتخاب بهترین و کارترین مدل از میان کلیه مدلها، دانست. میزان درصد حجم تراکنشهای مجموعه داده آموزشی ۷۵٪، آزمایشی ۲۰٪ و اعتبار سنجی ۵٪ می باشد.

آماده سازی داده های مشکوک

چون از ابتدای تحقیق پایگاه داده مربوط به تراکنشهای کارتهای بانکی برای هر فرد بعنوان رفتار قانونی آن فرد، فرض شده است، بنابراین به دلیل ارزیابی مدل در خصوص شناسایی رفتار مشکوک در این مرحله ۲۰ رکورد داده جدید با رفتاری مشکوک در پایگاه داده قبلی اضافه گردیده شد. ساز و کار تهیه تراکنشهای مشکوک بصورت دستی و با اعمال تغییراتی بصورت معنی وار بر روی داده ها بود است

تعیین متغیرهای ورودی

در این مرحله فیلدهای اطلاعاتی نتیجه تراکنش(X1)، مقدار تراکنش(X2)، نوع تراکنش(X3)، شمارش کد دستگاه (بهای کد دستگاه) X4 و پنج روزه تاریخ انجام تراکنش (به جای روز انجام تراکنش) X5 به عنوان پنج متغیر ورودی جهت ساخت مدل های خوشبندی مورد انتخاب قرار گرفت.

مدل سازی و ارزیابی جهت انتخاب مدل بهینه

ساخت مدل های خوشبندی

در این مرحله با استفاده نرم افزار داده کاوی اقدام به طراحی مدل های خوشبندی مختلفی با تکنیک شبکه عصبی کو亨ن شد، ۳ مدل با استفاده از همان ۵ متغیر ورودی تشریح شده (X1, X2,...,X5) و مجموعه داده های ورودی آموزشی و ارائه تنظیمات مختلف بر روی این مدلها مطابق با جدول ۲ طراحی شدند.

جدول ۲: مشخصات مدل های تولید شده شبکه عصبی کو亨ن

شماره مدل	تعداد خوشبندی	لایه ورودی(خروجی)	لایه خروجی(خروجی)	تنظیمات(طول×عرض)	نوع مدل
۱	هشت	۱۲	۹	۳×۳	خبره
۲	نه	۱۲	۹	۳×۳	ساده
۳	یازده	۱۲	۱۲	۳×۴	ساده

تولید یک طرح آزمون جهت ارزیابی مدلها و انتخاب مدل بهینه

جهت ارزیابی نتایج پروژه های داده کاوی، روش های مشخص و متنوعی جهت ارزیابی مدل های نظریتی به ویژه مدل های طبقه بندی، موجود می باشد، این موضوع برای مدل های غیر نظریتی صدق نمی کند، در مقابل یکی از مهم ترین روش های ارزیابی مدل های غیر نظریتی، تبدیل این گونه مسائل به یک مدل نظریتی و استفاده از روش های ارزیابی مرتبط با این گونه از مدل های می باشد.<sup>[۸ و ۹]</sup> در تحقیق حاضر نیز از این مفهوم، جهت تبدیل هر مدل خوش بندی به یک مدل طبقه بندی و ارزیابی از طریق روش های ارزیابی مدل های طبقه بندی، استفاده گردیده شد.

در همین خصوص گامهای زیر جهت ارائه یک الگوریتم ارزیابی، پیشنهاد می گردد:

بعد از خوشبندی داده ها، هر خوشبندی داده در این مدل را بعنوان یک کلاس در نظر گرفته و در صورت عدم انتخاب خود کار اسم برای آن، هر کلاس را با اسمی مانند... Cluster1, Cluster2 نامگذاری کرده و مجموعه تراکنشهای عضو آن خوش بندی را متعلق به کلاس نام گذاری شده متناظر آن نیز می دانیم.

در این مرحله توسط همان داده های آموزشی ابتدایی اقدام به آموخته متناظر مدل مجموعه قوانین C5 برای هر یک از مدل های خوش بندی می گردد. لازم به ذکر است فیلدهای اطلاعاتی ورودی مدل های خوش بندی (X1,X2,...,X5) به عنوان فیلدهای ورودی مدل های طبقه بندی و اسمی کلاس های تخصیص داده شده به عنوان متغیر خروجی در مدل های طبقه بندی می باشد.

در این مرحله با عبور مجموعه داده های آزمون و ارزیابی از مجموعه مدل های خوش بندی و طبقه بندی مجموعه قوانین، میزان قابلیت اطمینان مدل ها محاسبه می گردد.

با ارزیابی مقادیر شاخص های قابلیت اطمینان مدل ها مطابق با جدول ۳، مدل بهینه انتخاب می گردد.

# همایش ملی الکترونیکی دستاوردهای نوین در علوم مهندسی و پایه

National e-Conference on Advances in Basic Sciences and Engineering

WWW.AEBSCONF.IR



جدول(۳): معیارهای ارزیابی مدل‌ها

ردیف	تشریح معیار	شاخص
۱	قابلیت اطمینان کل مدل طبقه بندی برای مجموعه داده‌های آموزشی	تعداد کل طبقه بندی صحیح
۲	قابلیت اطمینان کل مدل طبقه بندی برای مجموعه داده‌های آزمایشی	تعداد کل داده‌های ارائه گردیده
۳	قابلیت اطمینان کل مدل طبقه بندی برای مجموعه داده‌های ارزیابی	تعداد کل طبقه بندی صحیح
۴	بیبود در روند شاخص‌های قابلیت اطمینان برای مجموعه داده‌های ارزیابی نسبت به داده‌های آزمون نسبت به داده‌های آموزشی	تعداد کل داده‌های ارائه گردیده میزان روند رو به رشد

## ارزیابی نهایی و انتخاب مدل پیشنهادی

با توجه به نتایج قابلیت اطمینان کل مدل‌ها که در جدول(۴) مشخص می‌باشد، به دلیل مقدار کم و روند نزولی شاخص قابلیت اطمینان کل، مدل شماره ۱ را به عنوان مدل نهایی جهت کشف رفتار مشکوک انتخاب می‌کنیم.

جدول (۴): مقایسه قابلیت اطمینان کل ۳ مدل با روش شبکه عصبی کوهنن

شماره مدل	شاخص قابلیت اطمینان با استفاده از داده‌های آزمون(درصد)	شاخص قابلیت اطمینان با استفاده از داده‌های ارزیابی(درصد)	شاخص قابلیت اطمینان با استفاده
۱ مدل	۱۰۰	۱۰۰	۱۰۰
۲ مدل	۹۶.۹۲	۹۶.۹۲	۹۹.۵۶
۳ مدل	۹۲.۵۹	۱۰۰	۹۹.۵۶

## شناسایی رفتار مشکوک

در مدل‌های طراحی شده اعدادی در خروجی مدل‌های طبقه بندی مجموعه قوانین C5 به عنوان مقادیر احتمال قابلیت اطمینان تخصیص یک رکورد به کلاس تخصیص داده شده قابل محاسبه و استفاده می‌باشد، به این صورت فیلد جدیدی که با عبارت \$CC\$ شروع می‌گردد و حاوی این مقادیر می‌باشد، در پایگاه واده موجود ایجاد می‌شود حال جهت شناسایی رفتار مشکوک از این مقادیر استفاده می‌کنیم. مینیمم مقدار آن اعداد را در داده‌های آموزشی پیدا کرده و در صورت اینکه مقادیر این شاخص برای داده‌های جدید(برای مثال داده‌های آزمایشی) کمتر از آن مقدار مینیمم باشد، آن تراکنش(رکورد) را به عنوان تراکنش مشکوک(رفتار مشکوک) در نظر می‌گیریم. مدل حاضر ۱۷ رفتار را به عنوان رفتار مشکوک شناسایی کرد که ۱۴ عدد از تراکنشها شناسایی شده از همان ۲۰ تراکنش مشکوک تهیه شده قبلی به روش دستی می‌باشد، این مقدار خود نیز نشان دهنده احتمال ۷۰ درصدی شناسایی رفتار مشکوک توسط این مدل می‌باشد.

## استخراج قواعد، تحلیل و نتیجه‌گیری

در این مرحله ۱۳ قانون مربوط به مدل طبقه بندی مجموعه قوانین C5 متناظر با شبکه کوهنن با تعداد ۸ خوشه را استخراج نموده‌ایم و با بررسی های بیشتر بر روی مکانیزم تولید هر یک از این قوانین، رفتارهای مشکوک شناسایی شده توسط این مدل و همچنین نکات دیگر بدست آمده در حین تحقیق، بصورت کلی نتایج زیر حاصل گردید. نتایج فرد گرایانه همچون: اگر شخصی با کارت فرد مورد نظر از تاریخ ششم ماه تا تاریخ ۱۰ ماه تراکنش از طریق پوز انجام دهد این تراکنش مشکوک به نظر می‌رسد چون معمولاً در مدل رفتاری این فرد هیچگاه چنین تراکنش انجام نشده است، و یا اینکه فرد مذکور همیشه از طریق یک ترمینال خاص از طریق اینترنت تراکنش مالی انجام رفتارهای مشکوک قرار گیرد و موارد دیگر...

آن دسته از خصوصیات متغیرهای ورودی که در کل داده‌های تاریخی، دارای سهم(درصدی) پایینی می‌باشند، دارای احتمال بالای برای تخصیص نسبتهای بالا به خود، در کل رفتار مشکوک شناسایی شده هستند. برای مثال تعداد تراکنشهای POS که در کل داده‌های آموزشی سهم ۳ درصدی دارا بوده است، در کل رفتار مشکوک شناسایی شده توسط مدل، سهمی معادل با ۵۳ درصد را دارا می‌باشد. عدم وجود یک فیلد در متغیرهای ایجاد کننده شروط در قوانین، دو معنی می‌تواند داشته باشد، معنی اول تنوع زیاد آن متغیر ورودی و عدم ایجاد الگو توسط مدل برای آن به دلیل تنوع زیاد می‌باشد و معنی دوم وجود رفتاری یکسان و الگوی بیش از اندازه شفاف و بدیهی برای آن متغیر می‌باشد.

# همایش ملی الکترونیکی دستاوردهای نوین در علوم مهندسی و پایه

National e-Conference on Advances in Basic Sciences and Engineering

WWW.AEBSCONF.IR



از تکنیک شبکه عصبی کوهن می‌توان جهت غربال‌گری اولیه داده‌ها در سیستمهای کشف رفتار مشکوک در زمان حقیقی استفاده نموده تعبیری می‌توان ابتدا یک سری از داده‌ها را توسط این مدل از کل تراکنشهای صورت گرفته به عنوان رفتارهای مشکوک انتخاب نموده و سپس توجه خود را بصورت دقیق تر بر این داده‌ها جهت کشف رفتار متنقلانه معطوف نمائیم، این موضوع باعث کاهش هزینه‌ها و افزایش سرعت و امکان کشف رفتار متنقلانه در زمان حقیقی می‌گردد.

مدل طراحی شده در این تحقیق فرد گرا بوده و نتایج مرتبط با آن نیز به میزان بالای ممکنی بر رفتار افراد می‌باشد و در صورت ارائه داده‌های افراد دیگر به مدل احتمال وجود نتایج متفاوت وجود دارد، اما می‌توان اینگونه بیان نمود که مدل طراحی شده فعلی قابلیت تعیین و استفاده برای مشتریانی با رفتاری شبیه به رفتار فرد مذکور را دارا می‌باشد.

یکی از نتایج مهم بدست آمده در این تحقیق وجود الگوهای مختلف در فرایند اکتشاف می‌باشد که این موضوع باعث می‌گردد که ما با نگاه از طریق مختلف به مساله همواره به دنبال استفاده از تکنیک‌های متفاوت جهت حل مساله باشیم که نهایتاً این موضوع نیز باعث هم افزایی و کشف دانش بیشتر جهت حل مساله می‌گردد.

## پیشنهاد برای تحقیقات آتی

کشف رفتار مشکوک با رویکرد غیر ناظارتی و استفاده از تکنیک‌های دیگر خوش بندی مانند Two-K-Means-Step ...

توسعه مکانیزم ارزیابی مدل‌ای غیر ناظارتی و ناظارتی با استفاده از شاخص‌ای گستره‌تر و استفاده از تکنیک‌های تصمیم‌گیری در ارزیابی مدل‌ها جهت انتخاب مدل بهینه کشف رفتار مشکوک کارت‌های بانکی.

ترکیب مدل‌های غیر ناظارتی و شبکه عصبی پیش خور جهت کشف رفتار مشکوک در تراکنش‌های کارت‌های بانکی.  
چگونگی استفاده از مدل‌های ناظارتی و غیر ناظارتی بصورت توان در طراحی سیستمهای مکانیزه کشف تقلب در بانکها، جهت ایجاد هم افزایی در این خصوص.

تحقیق در خصوص طراحی سیستم یکپارچه ثبت و ذخیره سازی اطلاعات تراکنشها در بانکها جهت تقسیم بندی و برچسب زنی داده‌های تراکنشها در دو دسته تراکنشهای قانونی و متنقلانه.

## منابع

- [1] Bolton,R.& Hand, D.Statistical Fraud Detection: A Review( With Discussion). Statistical Science 17(3)
- [2] Phua, C. , And Lee, V ., et al. A comprehensive survey of data mining-based fraud detecton reseach,2005.
- [3] Zaslavshy,v.,& Strizhak, A .Credit Card fraud detection using self-organizing maps. Information and Security. 18,48-63,2006
- [4] Quah, J. T. S., and Sriganesh, M. "Real-time credit card fraud detection using computational intelligence" Expert Systems with Applications35(4):1721-1732,2008
- [5] Larose,Daniel T., "Discovering Knowledge In Data Mining", John Wiley & Sons, Inc, Hoboken ,2005
- [6] Chapman, P., Clinton, J .,Kerber, R., Khabaza Th., Rienart, Th., Shearer.C., & Writh R., "CRISP\_DM Stepby step Sata Mining Guide"
- [7] Gupta, G. k.(2006) Introduction to Data Mining with Case Studies, PHI Learning Pvt. Ltd.
- [8] Williams,G. , Huang, Z.(1997) Mining the Knowledhe Mine: The Hot Spots Methodology for Mining Largi Real World Databases. Proc. Of the 10<sup>th</sup> Australian Joint Conference On Artificial Intelligence
- [9] Roiger, R . J. and M. W Geatz(2003) Data Mining: A tutorial-based Primer,Addison Wesley Boston.
- [10] Kantardezic, m.,(2003), " Data Mining Concepts, Models,Mothosd and algorithms" John wiley & Sons. Edited by jj, Hoboken.