

## Presenting a fraud detection model in online banking systems based on credit card transactions using multiple weighted random forest and quadratic model

Farzaneh Rahmani  | Changiz Valmohammadi \* | Kiomars Fathi 

---

**Article Info****Keywords:**

Fraud detection, online banking, credit card transactions, multiple weighted random forest, quadratic model

**ABSTRACT**

With the increasing growth of online banking, banks and financial institutions are more and more inclined to use this technology and its services. Due to the high volume of transactions, it is practically impossible to manage them by human resources. For this purpose, today, approaches based on data mining have come online with the help of banking. In this article, an efficient model for identifying fraudsters in bank card transactions is presented. The proposed method uses the adjacency matrix, placement of non-valued features using weighting, and random forest aggregation algorithm, in each branch of which, by calculating the weight of each branch, the best branch of the decision maker is selected by calculating the cost of the selection model. It can be It also selects the best forest for decision-making using the multiple quadratic model. Thus, we have tested this method on two data sets, the first one had ۱۴ features and the second one had ۲۰ features, and it has been observed that the model of this research compared to the decision tree, support vector machine, neural network, and normal random forest, which is currently the highest The results have shown improvements over any method. Also, the tests show that none of the mentioned methods were able to predict the OOB error and the normal random forest which is able to predict this error performed much weaker than the proposed model..

---

---



## ارائه مدل شناسایی متقلبین در سیستم‌های بانکداری آنلاین بر مبنای تراکنش‌های کارت‌های اعتباری با استفاده از جنگل تصادفی وزن‌دار چندگانه و مدل کوادراتیک

فرزانه رحمانی

گروه مدیریت فناوری اطلاعات، دانشگاه آزاد اسلامی، واحد تهران جنوب، تهران، ایران [rahmani.f.it@gmail.com](mailto:rahmani.f.it@gmail.com)

چنگیز والمحمدی

گروه مدیریت صنعتی، دانشکده مدیریت، دانشگاه آزاد اسلامی، واحد تهران جنوب، تهران، ایران

[ch\\_valmohammadi@azad.ac.ir](mailto:ch_valmohammadi@azad.ac.ir)

کیامرث فتحی دانشیار

دانشگاه آزاد اسلامی واحد تهران جنوب [fathikiamars@yahoo.com](mailto:fathikiamars@yahoo.com)

### چکیده

با رشد روزافزون بانکداری برخط ۱ بانک‌ها و مؤسسات مالی روزبه‌روز بیشتر به سمت استفاده از این فناوری و خدمات آن سوق پیدا می‌کنند. با توجه به حجم بالای تراکنش‌ها امکان مدیریت آن‌ها توسط نیروی انسانی عملاً غیرممکن است. به همین منظور امروزه رویکردهای مبتنی بر داده‌کاوی به کمک بانکداری برخط آمده است. در این مقاله یک مدل کارآمد برای شناسایی متقلبین در تراکنش‌های کارت‌های بانکی ارائه می‌گردد. روش پیشنهادی از ماتریس مجاورت، جای‌گذاری ویژگی‌های بدون مقدار با استفاده از وزن‌دهی و الگوریتم تجمیعی جنگل تصادفی استفاده می‌کند که در هر انشعاب آن با محاسبه وزن هر انشعاب، بهترین انشعاب تصمیم‌گیرنده با استفاده از محاسبه هزینه مدل انتخاب می‌شود. همچنین با استفاده از مدل کوادراتیک چندگانه بهترین جنگل را برای تصمیم‌گیری انتخاب می‌نماید. بدین ترتیب این روش را بر روی دو مجموعه داده که اولی ۱۴ ویژگی و دومی ۲۰ ویژگی داشته است تست کرده‌ایم و مشاهده شده است که مدل این تحقیق در مقایسه با درخت تصمیم و ماشین بردار پشتیبان و شبکه عصبی و جنگل تصادفی معمولی که در حال حاضر بالاترین نتایج را نسبت به هر روشی از خود نشان داده‌اند نیز بهبودهایی داشته است. همچنین آزمایشات نشان می‌دهد که هیچ‌یک از روش‌های مذکور قادر به پیش‌بینی خطای OOB نبوده و جنگل تصادفی معمولی که قادر به پیش‌بینی این خطا می‌باشد بسیار ضعیف‌تر از مدل پیشنهادی عمل نموده است. فقط روش پیشنهادی است که می‌تواند این مقدار را محاسبه نماید و مقدار مناسبی برای آن پیش‌بینی نماید. همچنین در آزمایش روی مجموعه داده‌ی دوم نیز در همین حدود بهبودهایی داشته ایم که به تفصیل در مقاله ذکر شده است.

**کلیدواژه‌ها:** شناسایی متقلبین، بانکداری آنلاین، تراکنش‌های کارت‌های اعتباری، ماتریس مجاورت، الگوریتم تجمیعی، جنگل تصادفی وزن‌دار چندگانه، مدل کوادراتیک

## مقدمه

مزایای بانکداری آنلاین را می‌توان از دو جنبه مشتریان و مؤسسات مالی مورد توجه قرار داد. از دید مشتریان می‌توان به صرفه‌جویی در هزینه‌ها، صرفه‌جویی در زمان و دسترسی به کانال‌های متعدد برای انجام عملیات بانکی نام برد. از دید مؤسسات مالی می‌توان به ویژگی‌هایی چون ایجاد و افزایش شهرت بانک‌ها در ارائه نوآوری، حفظ مشتریان علی‌رغم تغییرات مکانی بانک‌ها، ایجاد فرصت برای جستجوی مشتریان جدید در بازارهای هدف، گسترش محدوده جغرافیایی و فعالیت و برقراری شرایط رقابت کامل را نام برد (ویشال و همکاران ۲۰۲۱).

امروزه بانک‌ها و مؤسسات مالی و اعتباری برای خدمات‌رسانی مؤثر، ناگزیر از مهاجرت از بانکداری سنتی به بانکداری مدرن و برخط شده‌اند. هر چند استفاده از این سامانه‌ها باعث مدیریت بهتر فرایندهای مالی و افزایش کارایی و سرعت خدمات‌رسانی به مشتریان این مؤسسات شده، اما تقلب و سوءاستفاده‌های مالی یکی از مشکلاتی است که این سازمان‌ها در پی پیشگیری از آن‌ها و کاهش اثرات آن‌ها بوده‌اند. ایجاد سیستم‌های کارا برای شناسایی مشتریان بانک و کنترل فعالیت‌های مالی آنها بهترین روش برای مبارزه با مجرمان در جریان انجام عملیات بانکی است (فانگ و همکاران ۲۰۲۱).

امروزه روش‌های داده‌کاوی به‌عنوان بهترین راهکار برای شناسایی خودکار تقلب در حوزه‌های مختلف شناخته شده‌اند. داده‌کاوی به‌عنوان فرایند کشف و استخراج الگوهای پنهان از حجم بالایی از داده‌ها تعریف می‌شود. در سامانه‌های بسیاری از روش‌های داده‌کاوی برای شناسایی و کشف تقلب و سوءاستفاده مالی استفاده شده است (پتیدار و شارما ۲۰۱۱).

بخش عمده‌ای از فعالیت‌های متقلبانه، معطوف به تراکنش با کارت‌های اعتباری است. از این رو ایجاد سیستمی که ناظر بر عملکرد نظام‌های پرداخت باشد، به‌منظور شناسایی تقلب در تراکنش‌های موجود در کارت‌های اعتباری بانکی، ضروری به نظر می‌رسد. تاکنون مطالعات مختلفی برای شناسایی تقلب در تراکنش‌ها انجام شده است که هر کدام دارای مزایا و معایب خاص خود هستند. در این مقاله روشی ارائه می‌شود که علاوه بر تعیین ویژگی‌های مؤثر در دسته‌بندی تراکنش‌ها، از دقت و سرعت عملکرد بهتری نسبت به سایر روش‌ها، بهره‌مند گشته است (کارتا و همکاران در ۲۰۲۰). از آنجایی که روش‌های داده‌کاوی به‌صورت پویا با محیط‌های در حال تغییر، سازگار می‌شوند، مدت زمانی را که صرف کشف الگو می‌شود، تا حد قابل توجهی نسبت به روش‌های غیر خودکار کاهش می‌دهند (وانگ و همکاران ۲۰۲۱).

در این تحقیق به دنبال پاسخ به این سؤال هستیم که استفاده از ماتریس مجاورت و الگوریتم تجمیعی جنگل تصادفی، کارایی کشف تقلب در سیستم بانکی را تا چه حد بهبود می‌بخشد؟

## مبانی نظری و پیشینه تحقیق

در مقاله قلی‌پور و همکاران به‌منظور شناسایی تقلب در تراکنش‌های موجود در کارت‌های اعتباری بانکی، از روش یادگیری دسته‌جمعی استفاده گردیده است. در این روش علاوه بر تعیین ویژگی‌های مؤثر در دسته‌بندی تراکنش‌ها، دقت و صحت دسته‌بندی افزایش یافته است (قلی‌پور و همکاران ۱۴۰۰).

در مقاله وثوق و همکاران به روش‌های کشف تقلب پرداخته شده است. در تقلب‌های برخط، تراکنش‌ها از راه دور انجام شده و تنها به جزئیات کارت نیاز است نه لزوماً خود کارت. آن‌ها دریافته‌اند که به علت استفاده‌ی گسترده از اینترنت، کاربران می‌توانند موقعیت و هویت تراکنش اینترنتی خود را پنهان کنند (وثوق و همکاران ۱۳۹۸).

در مقاله حاتمی راد و همکاران به موضوع تقلب و اینکه در تقلب قصد شخصی و فریب دیگران مطرح است می‌پردازد. روش مطرح شده در این مقاله به استخراج ویژگی‌های پنهان پرداخته و سپس از بین ویژگی‌های استخراجی بهترین ویژگی‌ها جهت تصمیم‌گیری در مورد یک تراکنش را انتخاب می‌نماید (حاتمی راد و همکاران ۱۳۹۷).

بنایی و همکاران در مقاله خود مطرح نمودند که با توجه به حجم بالای تراکنش‌ها باید با رویکردهای مبتنی بر داده‌کاوی به کمک بانکداری برخط رفت تا مشکلات مربوط به تقلب را شناسایی نمود. آن‌ها از روش پردازش داده‌های حجیم بهره بردند (بنایی و همکاران ۱۳۹۶).

نظیر و همکارانش در ۲۰۲۳ مقاله‌ای در راستای تشخیص تقلب در کارت‌های اعتباری با استفاده از دو تکنیک یادگیری عمیق و رمزنگاری اطلاعات منتشر کردند. آن‌ها در لایه ابتدایی داده‌ها را رمزگذاری و در لایه انتهایی داده‌ها را رمزگشایی نمودند. اگرچه این روش پیچیدگی محاسباتی را افزایش داده است اما استفاده از رمزنگاری در کنار یادگیری عمیق چندلایه سبب افزایش دقت و قدرت تشخیص تقلب و همچنین کاهش احتمال وقوع تقلب شده است.

موریرا و همکاران در ۲۰۲۲ روشی مبتنی بر رگرسیون لجستیک، بیزین و نزدیک‌ترین همسایه برای شناسایی تقلب در سیستم بانکی مطرح نمودند. آن‌ها ابتدا به استخراج ویژگی‌های مؤثر پرداختند و فضای ویژگی ایجاد شده را با روش‌های نام برده شده دسته‌بندی نمودند. از نظر نویسندگان این مقاله، روش‌های یادگیری ماشین قادر است با دقت بالایی تراکنش‌های بانکی را دسته‌بندی نماید و علاوه بر آن می‌تواند صحت دسته‌بندی تراکنش‌ها را نیز افزایش دهد. نتایج به دست آمده نیز بیانگر همین موضوع است.

در مقاله دیگر مجموعه داده آموزش جدید، برای تبدیل داده بانک به داده‌های مناسب برای الگوریتم CLOPE که در خصوص تکنیک خوشه‌بندی برای داده اسمی (مقادیر رشته‌ای) می‌باشد، به منظور تشخیص موارد تقلب، ارائه گردیده است. نتایج آزمایشی نشان می‌دهد که CLOPE یک الگوریتم مناسب برای تشخیص موارد متقلبانه می‌باشد. اما این سیستم نمی‌تواند به تنهایی، به طور کامل، اجرا گردد و باید از توانائی تحلیل گران در تجزیه و تحلیل داده‌ها، و ارائه مجموعه‌ای از قوانین (معیارهای تعیین شده) برای اعتباربخشی به خوشه‌ها پس از عمل خوشه‌بندی، استفاده نمود (ویشال و همکاران ۲۰۲۱).

مقاله (دزاسکس و همکاران ۲۰۲۱) یک رویکرد تشخیص ناهنجاری هیبریدی است که استفاده از خوشه‌بندی برای ایجاد رفتارهای نرمال مشتریان و استفاده از تکنیک‌های آماری برای تعیین انحراف معامله خاص از رفتار گروه مربوطه است. این رویکرد بر روی یک مجموعه داده واقعی که شامل ۸,۲ میلیون معاملات انجام می‌شود، مورد آزمایش قرار گرفته و نتایج نشان می‌دهد که TEART به خوبی از لحاظ پارامترهایی که در مقایسه با الگوریتم K-mean سنتی به دست می‌آید، خوب است.

در (آتا و همکاران ۲۰۲۱) نویسندگان برای تشخیص رفتارهای غیرمعمول مشتری ماشین بردار پشتیبان (هند و همکاران ۲۰۱۹) را توسعه داده‌اند. آن‌ها ترکیبی از الگوریتم‌های نظارت شده و بدون ناظر ماشین بردار پشتیبان را ارائه داده‌اند. مزیت این روش آن است که می‌تواند با مجموعه داده‌های ناهمگون کار کند. با این حال ارزیابی عملکرد آن بر اساس مجموعه داده‌های شبیه‌سازی شده برای موارد مشکوک می‌باشد.

نویسندگان در (فوا ۲۰۲۰) بیان داشته‌اند، درخت تصمیم یکی از پرکاربردترین روش‌های استنباط استقرایی از سال ۱۹۶۰ تاکنون است. در این مقاله روش درخت تصمیم برای تعیین میزان ریسک تشخیص تقلب و پول‌شویی، بر اساس مشخصات مشتری بکار گرفته شده است

در طرح (یونکوک و همکاران ۲۰۱۹)، یک چارچوب تشخیص تقلب آنلاین بانکی مؤثر ارائه شده که به ترکیب منابع مربوط و به کارگیری تکنیک‌های داده کاوی پیشرفته می‌پردازد. در با ایجاد یک بردار تقابل ۱ برای هر تراکنش، بر اساس توالی رفتار مشتری در طول زمان، نرخ تمایز تراکنش جاری مشتری را با رفتار رایج وی نشان داده شده است. در این تحقیق، یک الگوریتم جدید برای استخراج مؤثر الگوهای تضاد و تشخیص رفتار جعلی از رفتار واقعی، برگرفته از یک الگوی انتخاب مؤثر و رتبه‌بندی ریسک که پیش‌بینی‌ها از مدل‌های مختلف را ترکیب می‌کند، معرفی شده است.

در سال‌های اخیر (سالچنبرگر و همکاران ۲۰۱۸)، شبکه عصبی مصنوعی (ANN) را به دلیل زمینه‌های کاربرد گسترده مطرح نمودند. در بسیاری از این برنامه‌های کاربردی تمرکز بر یادگیری حساس به هزینه وجود دارد به طوری که هزینه‌های مختلفی برای انواع مختلف طبقه‌بندی نادرست وجود دارد. هزینه طبقه‌بندی نادرست یک مثال از یک زمینه متفاوت است. در بسیاری از طبقه‌بندی‌های دودویی حساس به هزینه مانند مشکلات تشخیص، دو طبقه‌بندی نادرست مختلف وجود دارد و هر یک از آنها دارای هزینه است. با این حال، در کسب و کار مشکلاتی از قبیل تشخیص جعل کارت اعتباری و بازاریابی مستقیم هر مشاهده طبقه‌بندی شده، می‌تواند هزینه متفاوتی داشته باشد و علاوه بر این ممکن است سودی برای درستی طبقه‌بندی هر یک وجود داشته باشد؛ بنابراین، در چنین مواردی، ضرورتی برای توسعه یک مدل طبقه‌بندی وجود دارد که به منافع و هزینه‌های شخصی رسیدگی کند.

در مقاله‌ای دیگر (واهلر و همکاران ۲۰۱۵) نشان داده شده است که استفاده از قواعد انجمنی و به کارگیری آن‌ها بر روی مجموعه داده‌گان تراکنش‌های بانکی نتایج مطلوبی را حاصل نموده است. همچنین در رویکرد دیگری (پوزولو و همکاران ۲۰۱۴ و باتاچاریان و همکاران ۲۰۱۱) ترکیب این روش در کنار تکنیک یادگیری ماشین توانسته است دقت و صحت تشخیص را افزایش دهد. نویسندگان در (باتاچاریا و همکاران ۲۰۱۱)، به بررسی برخی مدل‌های پیش‌بینی‌کننده معروف داده کاوی برای شناسایی تقلب پرداخته‌اند. در این مطالعه از دو تکنیک داده کاوی رگرسیون لجستیک و ماشین بردار پشتیبان برای شناسایی تقلب استفاده شده است. در ادامه ابتدا به معرفی این دو تکنیک می‌پردازیم.

هنگام برخورد با ساختارهای گرافی، ناهنجاری می‌تواند طبق مشخصات گراف به خوبی طبقه‌بندی شود. در تحقیق (ابرله ۲۰۰۷) ناهنجاری‌ها را بر اساس آنچه گفته شد به سه گروه تقسیم‌بندی کرده‌اند. انجام درج با وجود یک رأس یا یک لبه غیرمنتظره در گراف. انجام اصلاح با حضور یک برجسب غیره منتظره روی یک رأس یا یک لبه. حذف شامل عدم وجود یک رأس یا لبه‌ی مورد انتظار می‌باشد. گاهی اوقات، حتی شامل مفاهیم، مرتبط با لبه نیز می‌شود به عنوان مثال حذف یک رأس خاصی از کل لبه‌های مجاور که ممکن است حتی حذف شده باشد. این امر سبب بهبود تشکیل گراف و کشف ناهنجاری در یال‌های آن می‌شود. این یال‌ها هر یک می‌تواند یک تراکنش را شامل شود.

نویسندگان در (شن و همکاران ۲۰۰۷)، برای شناسایی تقلب در کارت‌های اعتباری، مقایسه‌ای بین روش‌های داده کاوی مختلف ارائه داده‌اند. در این مطالعه، سه روش پرکاربرد نظیر درخت تصمیم، شبکه‌ی عصبی و رگرسیون لجستیک استفاده شده است.

### کارکردهای مدیریتی:

با استفاده از روش ارائه شده در این مقاله و زمانی که تراکنش‌های مشکوک شناسایی شوند مدیریت تراکنش‌ها بسیار ساده تر خواهد شد و هزینه‌های مدیریتی کاهش پیدا خواهد کرد.

## شکاف تحقیقاتی

### شکاف عملی

تراکنش‌های تقلبی هزینه‌های هنگفتی را به بانک‌ها تحمیل می‌کنند. بنابراین مدلی مناسب و سودده است که هزینه‌ها را تا حد امکان کاهش دهد. داده‌های بانک‌ها حجم بالایی دارند و پردازش آن‌ها زمان‌بر است در نتیجه آن‌ها ملزم به استخدام نیروی انسانی زیادی هستند که زمان‌بر است و هزینه بالایی با به بانک تحمیل می‌کند.

### شکاف نظری

باتوجه به بررسی مطالعات انجام شده، شکاف تحقیقاتی موجود، عدم وجود مدلی است که بتواند اطلاعات حجیم و غیرنرمال را کاوش نموده و روشی مؤثر برای استخراج اطلاعات ارزشمند در داده‌های بزرگ و داده‌های ناهمخوان باتوجه به هزینه مدل باشد. مدلی که قادر باشد هم‌زمان با فرایند آموزش، متغیرهای مهمی که بیش‌ترین تأثیر را در طبقه‌بندی داده‌ها دارند، استخراج کند و سرعت پاسخگویی را افزایش دهد و بهترین مدل را با احتمال و دقت بالا تشخیص دهد. همچنین بتواند داده‌های بدون مقدار یا دارای مقدار نامتعارف و پرت را شناسایی نموده و برای آن‌ها مقدار مناسب جایگزین نماید.

### روش تحقیق

تحقیق حاضر یک تحقیق کاربردی در حوزه بانکداری آنلاین می‌باشد. مدل مطرح شده در این تحقیق می‌تواند به بهبود کیفیت خدمات بانکی و امنیت تراکنش‌ها کمک نماید.

در این مقاله برای شناسایی متقلبین، از ماتریس مجاورت، مدل تجمیعی جنگل تصادفی وزن‌دار و الگوریتم مدل احتمالاتی استفاده کرده‌ایم. در مدل پیشنهادی ویژگی‌های بدون مقدار و دارای مقدار غیرنرمال مقداردهی مناسب می‌گردد. همچنین با استفاده از جنگل تصادفی وزن‌دار که برای هر انشعاب آن وزن و هزینه محاسبه می‌شود و برای هر دسته داده بهترین انشعاب جهت تصمیم‌گیری انتخاب می‌گردد، بهترین و کم‌هزینه‌ترین تصمیم برای هر نوع داده اتخاذ می‌شود. همچنین با استفاده از الگوریتم محاسبه احتمال، بهترین مدل بر روی هر انشعاب ایجاد می‌گردد. در واقع این مدل یک روش یادگیری دسته‌جمعی محسوب شده و برای یادگیری از تعداد زیادی درخت تصمیم استفاده می‌کند که در هر درخت محاسبه وزن و هزینه برای هر نوع داده و هر ویژگی لحاظ می‌شود. پس از پایان آموزش، برای دسته‌بندی یک نمونه‌ی جدید، بین درختان رأی‌گیری بر اساس بهترین مدل برای آن داده انجام شده و کلاس با بیش‌ترین رأی و کم‌ترین هزینه، برای نمونه‌ی جدید انتخاب می‌شود.

در این مدل، دو پارامتر توسط کاربر تعیین می‌شود. پارامتر اول تعداد درختان تصمیم است که در جنگل ساخته خواهد شد. این پارامتر بسته به تعداد داده‌های آموزشی و تعداد ویژگی‌های موجود در هر مسئله‌ای متغیر است.

پارامتر دوم با نام  $m$  شناخته می‌شود. در زمان انتخاب ویژگی شکست در هر گره‌ی درخت، باید  $m$  ویژگی به صورت تصادفی انتخاب شده و از بین این  $m$  ویژگی، با کمک معیارهای کارایی بهترین ویژگی برای بخش‌بندی داده‌ها انتخاب شود.

این پارامتر در کل فرایند آموزش مدل و در بین کلیه‌ی درختان ثابت در نظر گرفته می‌شود. مقادیر معمولی که برای این پارامتر در نظر گرفته می‌شود، عبارت‌اند از:  $\sqrt{nVariable}$ ،  $\log(nVariable)+1$  و منظور از  $nVariable$  تعداد ویژگی‌های موجود در مجموعه داده‌ها یا همان متغیرهای مسئله است.

همچنین برای آموزش هر درخت، از مجموعه داده‌های آموزشی اولیه، به طور تصادفی و با جای‌گذاری به تعداد  $N$  نمونه‌ی آموزشی انتخاب می‌شود. پارامتر  $N$  را معمولاً به اندازه‌ی کل مجموعه داده‌های اولیه‌ای که در دسترس است، در نظر می‌گیرند.

بنابراین، ممکن است در بین زیرمجموعه‌های آموزشی ایجاد شده اشتراکاتی وجود داشته باشد. نکته‌ی قابل توجه اینجاست که این مدل می‌تواند پدیده‌ی بیش‌برازش<sup>۱</sup> را به‌خوبی مدیریت کرده و عمومیت بیش‌تری را در فضای مسئله‌اش داشته باشد. منظور از عمومیت این است که مدل در مواجهه با داده‌های جدید، بتواند به‌خوبی با آن‌ها تطبیق پیدا کرده و کمترین خطای تشخیص یا پیش‌بینی را در پی داشته باشد. یعنی تنها به داده‌های آموزش محدود نبوده (overfit نشده) و عمومیت داشته باشد. در واقع به دلیل همین سیاست تصادفی عمل کردن در انتخاب زیرمجموعه‌های آموزشی و به‌خصوص در انتخاب  $m$  ویژگی، به‌خوبی از پدیده‌ی بیش‌برازش جلوگیری می‌شود.

### روند ساخت مدل

مدل پیشنهادی از ۳ بخش اصلی تشکیل می‌شود.

در بخش اول بارگذاری و پیش‌پردازش داده‌ها انجام می‌شود. مجموعه‌داده‌های در دسترس شامل ویژگی‌های عددی پیوسته و ویژگی‌های رده‌ای می‌باشند.

پس از بارگذاری داده‌ها، اگر ویژگی‌های رده‌ای با مقادیر رشته‌ای مقاردهی شده باشند، باید این مقادیر رشته‌ای به عدد تبدیل شوند. بخش دوم، پارامترهای موردنیاز الگوریتم مقاردهی اولیه می‌شوند. این پارامترها عبارت‌اند از: تعداد متغیرهای مسئله یا همان ویژگی‌های موجود در داده‌ها، تعداد نمونه‌های موجود در مجموعه‌داده‌ها، تعداد درختان موردنیاز، تعیین پارامتر  $N$  که تعداد نمونه‌های آموزشی را برای هر درخت تعیین می‌کند و تعیین پارامتر  $m$  برای مشخص شدن تعداد ویژگی‌هایی که در هر گره باید بررسی شده و بهترین آن برای بخش‌بندی داده‌ها انتخاب شود.

بخش سوم، مدل ساخته می‌شود. برای ساخت مدل، به‌ازای هر درخت، به تعداد  $N$  نمونه به‌صورت تصادفی و با جای‌گذاری از مجموعه‌داده‌ی اولیه انتخاب و به‌عنوان trainingSample ذخیره می‌شود.

روند کلی ساخت مدل در ادامه شرح داده شده است.

به تعداد درختان تصمیم‌موردنظر، از داده‌های آموزشی اولیه، زیرمجموعه‌های آموزشی را نمونه‌گیری می‌کنیم.

- در هر زیرمجموعه‌ی آموزشی، یک درخت تصمیم را بدون هرس کردن و تا انتهای فرایند آموزش توسعه می‌دهیم. برای هر درخت با توجه به تعداد ویژگی‌ها، یک وزن اولیه در نظر می‌گیریم. در هر گره، به‌جای انتخاب بهترین ویژگی شکست از بین کلیه‌ی ویژگی‌ها،  $m$  ویژگی را به‌طور تصادفی انتخاب کرده و از بین این  $m$  ویژگی، بهترین را به‌عنوان ویژگی شکست انتخاب می‌کنیم. وزن انشعاب را با توجه به اهمیت ویژگی انتخاب شده و هزینه محاسبه آن برای نمونه‌های موجود در انشعاب به‌روز می‌نماییم. یادگیرنده‌های ضعیف در حین اضافه‌شدن به مجموعه، وزن‌دهی می‌شوند که این وزن‌دهی بر اساس میزان دقت در طبقه‌بندی نمونه‌هاست. پس از اضافه‌شدن هر طبقه‌بند، نمونه‌های موجود (داده‌ها) نیز وزن‌دهی می‌گردند (وزن‌شان اصلاح می‌گردد). وزن‌دهی نمونه‌ها به‌صورتی است که در هر مرحله، وزن نمونه‌هایی که به‌صورت صحیح طبقه‌بندی می‌شوند کاهش یافته و وزن نمونه‌هایی که به‌درستی طبقه‌بندی نشده‌اند، بیشتر می‌شود تا در مراحل بعدی (توسط یادگیرنده‌های جدید) بیشتر موردتوجه بوده و با دقت بیشتری طبقه‌بندی گردند؛ بنابراین تمرکز یادگیرنده‌های ضعیف جدید، بیشتر بر روی داده‌های خواهد بود که سیستم در مراحل قبلی قادر به طبقه‌بندی صحیح آنها نبوده است.

<sup>۱</sup>overfit



۲. ماتریس هزینه داده‌ها را محاسبه می‌نمایم.

در جدول ۱ ماتریس هزینه محاسبه شده نشان داده شده است.

۳. جدول ۱: ماتریس هزینه برای دسته بندی

|                 |                |                |   |   |
|-----------------|----------------|----------------|---|---|
|                 |                |                |   |   |
| پیش بینی نادرست | TN یا $c(i,i)$ | FN یا $c(i,i)$ | ۱ | ۲ |
| پیش‌بینی درست   | FP یا $c(i,i)$ | TP یا $c(i,i)$ | ۱ | ۰ |

لازم به ذکر است که  $c(i,i)$  که همان (TN و TP) می‌باشد هنگامی که نمونه به درستی پیش‌بینی شود معمولاً به عنوان سود در نظر گرفته می‌شود. همچنین اقلیت یا طبقه نادر به عنوان کلاس مثبت در نظر گرفته می‌شود. مشخص است که هزینه‌ی دسته‌بندی نادرست یک نمونه‌ی اقلیت بیشتر از هزینه‌ی دسته‌بندی نادرست یک نمونه‌ی اکثریت است به همین دلیل ارزش FN معمولاً بیشتر از FP است. هزینه مورد انتظار نمونه‌ها محاسبه می‌شود.

با توجه به ماتریس هزینه یک نمونه باید در رده‌ای که کمترین هزینه از آن انتظار می‌رود دسته‌بندی شود. هزینه مورد انتظار  $R(i|x)$  از دسته‌بندی نمونه  $x$  در رده  $i$  ام به شرح ذیل بیان می‌گردد.

فرمول ۱:

$$R(i|x) = \sum_j P(j|x)C(i,j) \quad . ۱$$

که در آن  $P(j|x)$  احتمال تعلق نمونه  $x$  به رده  $j$  را مشخص می‌کند. در ابتدا هر دسته برای هر مجموعه به طور مستقل استفاده می‌شود. سپس نتایج به منظور تولید با هم ترکیب می‌شوند. روش مبتنی بر ترکیب دسته‌بندی‌ها معمولاً دارای دقت بالاتر، همچنین FP کمتر نسبت به دسته‌بندی‌های جداگانه دارد. روند کلی روش پیشنهادی در ادامه شرح داده شده است.

هر انشعاب را توسعه می‌دهیم تا زمانی که آموزش آن به پایان رسد.

داده‌ی جدید را به کلیه‌ی درختان اعمال کرده و برجسب نهایی داده، رأی اکثریت درختان و در نظر گرفتن کمترین هزینه خواهد بود. در مسئله‌ی رگرسیون، مقدار پیش‌بینی شده‌ی نهایی، میانگین مقادیر پیش‌بینی شده توسط همه‌ی درختان خواهد بود. روش‌های مختلفی برای ترکیب نتایج دسته‌بندی‌کننده‌ها وجود دارد، متداول‌ترین روش‌ها میانگین‌گیری و یا استفاده از رأی اکثریت هستند. به منظور نشان‌دادن رویکرد ترکیبی از دو تکنیک استفاده می‌شود. اولین روش قانون رأی‌گیری اکثریت ۱ می‌باشد. در این روش اظهار نظر هر دسته‌بند در مورد کلاس الگوی ورودی، به عنوان یک رأی محسوب می‌شود و تصمیم‌گیری نهایی بر اساس آرای اخذ شده از دسته‌بندی‌های مختلف صورت می‌گیرد. در این تکنیک تمامی داده‌ها دارای وزن یکسان هستند و حساسیتی نسبت به هزینه داده‌ها وجود ندارد.

۱. Majority voting rule.

۴. بهترین مدل برای داده محاسبه می گردد.

فرض ابتدایی در تعیین بهترین مدل این است که داده‌ها به صورت نرمال توزیع شده‌اند. در این روش نیازی به برابر بودن کوواریانس بین دو گروه و انشعاب نیست. ابتدا مجموعه‌ای از مشاهدات به نام  $x$  به طبقه‌بند داده می‌شوند که از نوع خانواده  $y$  هستند. به این قسمت مرحله تعلیم گفته می‌شود که متناسب با آن طبقه‌بند، سطح فرضی را برای جداسازی دو دسته از ویژگی‌ها بکار می‌برد. سپس پس از تعلیم طبقه‌بند، داده‌ها تست به آن داده می‌شود تا دو دسته ویژگی‌های مختلف را از یکدیگر جدا کند.

در روش پیشنهادی برای حساسیت به هزینه از تابع باور و برای محاسبه احتمال از تابع کوادراتیک استفاده می‌شود. به این ترتیب که زمانی می‌توانیم یک تراکنش را قانونی بدانیم که هر دو دسته‌بندی‌کننده آن را قانونی تشخیص دهند. با این کار می‌توان معیار FN (تراکنش ورودی تقلبی می‌باشد و سیستم آن را به اشتباه قانونی تشخیص داده است) که بیشترین هزینه را برای ما دارد به حداقل رساند. همچنین برای ایجاد وزن برای نمونه‌ها از تکنیک میانگین قانون ۱ استفاده می‌نماییم. به این صورت هر نمونه را با توجه به الگوی ورودی، با یک میانگین احتمال خلفی به یک کلاس تخصیص می‌دهیم. با استفاده از الگوریتم کوادراتیک، ویژگی‌های حاصل شده برای هر مدل و انشعاب طبقه‌بندی می‌شوند.

## ماتریس مجاورت ۲

یک ماتریس  $N \times N$  است و  $N$  تعداد کل مجموعه داده‌ی آموزشی اولیه است (بنسال و شرما ۲۰۲۱). زمانی که یک درخت تصمیم ساخته شد، داده‌های آموزش متعلق به خودش را به درخت اعمال می‌کنیم. اگر نمونه‌ی  $i$  ام با نمونه‌ی  $j$  ام در یک نود پایانی مشابه قرار گرفتند، عنصر  $(i, j)$  ماتریس مجاورت را یکی اضافه می‌کنیم. در نهایت، درایه‌های ماتریس را با تقسیم کردن بر تعداد کل درختان، نرمال می‌کنیم.

ماتریس مجاورت می‌تواند در تعریف ساختار داده‌ها و یا یادگیری غیر نظارت شده به کار گرفته شود (بنسال و شرما ۲۰۲۱).

## ویژگی‌های مدل پیشنهادی

در نهایت، ویژگی‌های مدل پیشنهادی را به صورت زیر خلاصه می‌کنیم:

- یک الگوریتم قدرتمند با ماهیت دسته‌جمعی بودن، محسوب شده که قدرت یادگیری و تعمیم خوبی را فراهم می‌آورد.
- این الگوریتم به صورت کارا بر روی مجموعه داده‌های بسیار بزرگ اجرا می‌شود.
- بدون حذف ویژگی (یا متغیر) می‌تواند با بیش از هزار متغیر در مسئله کار کند.
- برآوردی از میزان اهمیت هر یک از متغیرها را نشان می‌دهد.
- در برابر پدیده‌ی بیش‌برازش مقاوم است.
- می‌تواند یک تخمین بدون بایاس داخلی از خطای عمومی‌سازی را در حین فرایند آموزش، ارائه دهد.
- برای افزایش سرعت و کاهش زمان فرایند آموزش، می‌توان به صورت موازی درختان را آموزش داد.

- دارای روشی مؤثر برای برآورد داده‌های گم شده است و بدون کاهش دقت می‌تواند در مجموعه داده‌هایی که مقادیر گم شده‌ی زیادی دارند، به خوبی کار کند.
- مدل‌های ساخته شده می‌توانند برای استفاده‌های بعدی بر روی داده‌های دیگر، ذخیره شوند.
- این الگوریتم میزان مجاورت (یا شباهت) بین هر جفت از داده‌ها را محاسبه می‌کند. این امکان می‌تواند در خوشه‌بندی، برآورد داده‌های گم شده و همچنین ایجاد دید کلی در مورد داده‌ها بسیار مؤثر باشد.

### نتایج روش پیشنهادی

برای بررسی تأثیر پارامترهای مختلف مدل پیشنهادی آزمایشاتی را طراحی و اجرا کرده‌ایم. با این آزمایش‌ها مقادیر قابل قبولی برای پارامترهای آزاد مدل تعیین می‌شود. در ادامه نیز برای نمایش قدرت و کارایی مدل پیشنهادی، آن را با سایر روش‌های موجود مقایسه و نتایج به دست آمده ارائه شده است.

### داده‌ها و نرم‌افزار مورد استفاده

در این تحقیق از دو مجموعه داده‌ی استاندارد، متعلق به کارت‌های اعتباری موجود در کشور آلمان و کشور استرالیا استفاده می‌نمایم. در هر دو مجموعه داده برای حفظ امنیت، نام و مقادیر فیلدها به صورت یکتا، گمنام‌سازی شده‌اند. داده‌های فیلتر شده شامل نام و نام خانوادگی مشتریان و اطلاعات شخصی آن‌ها می‌باشد.

### پارامترهای منحصر به فرد در مدل پیشنهادی

به‌طور کلی در هر الگوریتم تعدادی پارامتر آزاد و مؤثر در کارایی الگوریتم وجود دارد. الگوریتم پیشنهادی نیز از این قضیه مستثنی نبوده و تعدادی پارامتر دارد که تنظیم بهینه‌ی آن‌ها موجب افزایش دقت و قدرت الگوریتم در شناسایی تقلب خواهد شد. در این بخش در قالب آزمایشاتی تأثیر این پارامترها را نشان داده‌ایم. کلیه‌ی آزمایش‌ها انجام شده در این بخش بر روی دو مجموعه داده‌ی استاندارد است که در بخش قبل معرفی شد، انجام شده و نتایج ارائه شده‌اند.

### تعداد انشعابات

یکی از پارامترهای مهم و مؤثر در دقت شناسایی، تعداد انشعابات در مدل پیشنهادی است. هر یک از انشعابات موجود در مدل به‌تنهایی کارایی زیادی نداشته و در واقع قدرت تشخیص مدل به برآیند قدرت کلیه‌ی آن‌ها وابسته است. در واقع این ویژگی روش‌های دسته-جمعی است که از نوعی سیستم رأی‌گیری بین اعضا برای تصمیم‌گیری نهایی استفاده می‌کند. هر چه تعداد این انشعاب‌ها در مدل بیش‌تر باشد دقت تشخیص بهبود می‌یابد.

برای مشاهده‌ی تأثیر این پارامتر در کارایی الگوریتم و انتخاب تعداد مناسب آن در مدل، آزمایش زیر انجام شده است:

### آزمایش ۱: تأثیر پارامتر تعداد انشعابات

در این آزمایش، الگوریتم پیشنهادی را با شرایط زیر اجرا کرده‌ایم:

۱. معیار انتخاب نقطه‌ی شکست = شاخص جینی

۲.  $m = n \text{Variable}$  به معنی آن که پارامتر  $m$  برابر با تعداد کل ویژگی‌های موجود است. با توجه به دو مجموعه داده‌ی در دسترس، با به کارگیری مجموعه داده‌ی اول این پارامتر مقدار ۱۴ و با به کارگیری مجموعه داده‌ی دوم، مقدار ۲۰ را خواهد داشت.

۳. تعداد انشعابات = متغیر بوده و از مقدار ۱ با گام افزایشی ۱۰ و تا مقدار ۲۰۰ مقداردهی خواهد شد.

معیارهای ارزیابی الگوریتم‌های دسته‌بندی که در بخش قبل به آن‌ها اشاره شد، هم در فرایند اجرای الگوریتم محاسبه و در جدول ۲ نشان داده شده است.

جدول (۲) معیارهای ارزیابی محاسبه شده بر حسب تعداد انشعابات در مدل پیشنهادی بر مجموعه داده‌ی استرالیایی با ۱۴ ویژگی

| F-measure | Preci si on | Sensi ti vi ty | Accuracy | تعداد انشعابات |
|-----------|-------------|----------------|----------|----------------|
| ۰,۸۴      | ۰,۸۶        | ۰,۸۲           | ۰,۸۴     | ۱              |
| ۰,۹۱      | ۰,۹         | ۰,۹۲           | ۰,۹      | ۱۰             |
| ۰,۹       | ۰,۹         | ۰,۹            | ۰,۸۹     | ۲۰             |
| ۰,۹۲      | ۰,۹۲        | ۰,۹۱           | ۰,۹۱     | ۳۰             |
| ۰,۹۲      | ۰,۹۲        | ۰,۹۳           | ۰,۹۲     | ۴۰             |
| ۰,۹۳      | ۰,۹۳        | ۰,۹۲           | ۰,۹۲     | ۵۰             |
| ۰,۹۲      | ۰,۹۳        | ۰,۹۲           | ۰,۹۲     | ۶۰             |
| ۰,۹۴      | ۰,۹۵        | ۰,۹۴           | ۰,۹۴     | ۷۰             |
| ۰,۹۳      | ۰,۹۴        | ۰,۹۲           | ۰,۹۲     | ۸۰             |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۳           | ۰,۹۳     | ۹۰             |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۳           | ۰,۹۳     | ۱۰۰            |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۳           | ۰,۹۳     | ۱۱۰            |
| ۰,۹۴      | ۰,۹۵        | ۰,۹۳           | ۰,۹۴     | ۱۲۰            |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۴           | ۰,۹۳     | ۱۳۰            |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۴           | ۰,۹۳     | ۱۴۰            |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۴           | ۰,۹۳     | ۱۵۰            |
| ۰,۹۳      | ۰,۹۴        | ۰,۹۲           | ۰,۹۲     | ۱۶۰            |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۴           | ۰,۹۴     | ۱۷۰            |
| ۰,۹۴      | ۰,۹۵        | ۰,۹۴           | ۰,۹۴     | ۱۸۰            |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۴           | ۰,۹۳     | ۱۹۰            |
| ۰,۹۴      | ۰,۹۵        | ۰,۹۳           | ۰,۹۳     | ۲۰۰            |

جدول (۳) معیارهای ارزیابی محاسبه شده بر حسب تعداد انشعابات در مدل پیشنهادی بر مجموعه داده‌ی آلمانی با ۲۰ ویژگی

| F-measure | Preci si on | Sensi ti vi ty | Accuracy | تعداد انشعابات |
|-----------|-------------|----------------|----------|----------------|
| ۰,۸۶      | ۰,۸۱        | ۰,۹۲           | ۰,۸      | ۱              |
| ۰,۸۵      | ۰,۸۹        | ۰,۸۱           | ۰,۸      | ۱۰             |
| ۰,۸۷      | ۰,۸۸        | ۰,۸۶           | ۰,۸۳     | ۲۰             |
| ۰,۸۹      | ۰,۸۹        | ۰,۸۸           | ۰,۸۵     | ۳۰             |
| ۰,۹       | ۰,۹۱        | ۰,۸۹           | ۰,۸۶     | ۴۰             |
| ۰,۹۱      | ۰,۹۳        | ۰,۸۹           | ۰,۸۸     | ۵۰             |
| ۰,۹۲      | ۰,۹۲        | ۰,۹۱           | ۰,۸۹     | ۶۰             |
| ۰,۹۲      | ۰,۹۱        | ۰,۹۳           | ۰,۸۹     | ۷۰             |
| ۰,۹۲      | ۰,۹۲        | ۰,۹۲           | ۰,۸۹     | ۸۰             |
| ۰,۹۳      | ۰,۹۳        | ۰,۹۳           | ۰,۹۱     | ۹۰             |
| ۰,۹۲      | ۰,۹۳        | ۰,۹۲           | ۰,۹      | ۱۰۰            |
| ۰,۹۳      | ۰,۹۱        | ۰,۹۴           | ۰,۹      | ۱۱۰            |
| ۰,۹۳      | ۰,۹۱        | ۰,۹۴           | ۰,۹      | ۱۲۰            |
| ۰,۹۳      | ۰,۹۲        | ۰,۹۵           | ۰,۹۱     | ۱۳۰            |
| ۰,۹۳      | ۰,۹۱        | ۰,۹۵           | ۰,۹      | ۱۴۰            |
| ۰,۹۴      | ۰,۹۳        | ۰,۹۵           | ۰,۹۱     | ۱۵۰            |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۴           | ۰,۹۲     | ۱۶۰            |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۴           | ۰,۹۲     | ۱۷۰            |
| ۰,۹۴      | ۰,۹۳        | ۰,۹۵           | ۰,۹۱     | ۱۸۰            |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۴           | ۰,۹۲     | ۱۹۰            |
| ۰,۹۴      | ۰,۹۴        | ۰,۹۴           | ۰,۹۲     | ۲۰۰            |

همان‌گونه که مشخص است با افزایش تعداد انشعاب‌ها معیارهای ارزیابی نیز روند صعودی دارند. پس قدرت الگوریتم در شناسایی نیز افزایش یافته است. باتوجه به جدول ۱ که مربوط به انجام آزمایش بر روی مجموعه داده‌ی استرالیای با ۱۴ ویژگی است، مشاهده می‌کنیم زمانی که تعداد انشعاب‌ها به ۹۰ می‌رسد تقریباً رفتار مدل نیز قابل پیش‌بینی بوده و تغییرات چندانی در معیارها مشاهده نمی‌شود. همچنین بهترین نتایج با تعداد انشعاب ۱۲۰، به دست آمده است. علاوه بر این، باتوجه به نتایج به دست آمده در جدول ۲ که مربوط به انجام آزمایش بر روی مجموعه داده‌ی آلمانی با ۲۰ ویژگی است، همین نکته را دریافتیم که با افزایش تعداد انشعابات در مدل دقت تشخیص الگوریتم افزایش می‌یابد. نکته‌ی قابل توجه دیگر آن است با افزایش تعداد ویژگی‌های مسئله به تعداد انشعاب‌ها بیش‌تری نیاز است تا به دقت مطلوب دست یابیم. چنان‌که مشاهده می‌کنیم زمانی که تعداد انشعاب به ۱۶۰ می‌رسد، رفتار مدل تقریباً ثابت می‌شود درحالی‌که در مجموعه داده‌ی قبلی که ۱۴ ویژگی داشت، با داشتن ۹۰ انشعاب به کارایی مطلوب دست می‌یابیم.

### درجه کوادراتیک

یکی از پارامترهای مهم در تعیین بهترین مدل در فضای ویژگی، درجه کرنل کوادراتیک در مدل پیشنهادی است. هر چه تعداد ویژگی‌ها و درجه کرنل افزایش می‌یابد، انتخاب ویژگی مناسب‌تر و در نتیجه دقت تشخیص بهبود می‌یابد. برای مشاهده تأثیر این پارامتر در کارایی الگوریتم و انتخاب درجه مناسب آن در مدل، آزمایش زیر انجام شده است:

#### آزمایش ۲: تأثیر پارامتر درجه کرنل کوادراتیک

در این آزمایش، الگوریتم پیشنهادی را با شرایط زیر اجرا کرده‌ایم:

۱. معیار انتخاب نقطه‌ی شکست = شاخص جینی
۲.  $m = nVariable$  به معنی آن که پارامتر  $m$  برابر با تعداد کل ویژگی‌های موجود است. با توجه به دو مجموعه داده‌ی در دسترس، با به کارگیری مجموعه داده‌ی اول این پارامتر مقدار ۱۴ و با به کارگیری مجموعه داده‌ی دوم، مقدار ۲۰ را خواهد داشت.
۳. درجه کرنل = متغیر بوده و از مقدار ۱ تا مقدار ۳ مقداردهی خواهد شد.

معیارهای ارزیابی الگوریتم‌های دسته‌بندی که در بخش قبل به آن‌ها اشاره شد، هم در فرایند اجرای الگوریتم محاسبه و در جدول ۴ نشان داده شده است.

جدول (۴) معیارهای ارزیابی محاسبه شده بر حسب کرنل کوادراتیک در مدل پیشنهادی بر مجموعه داده‌ی استرالیایی با ۱۴ ویژگی

| F-measure | Precision | Sensitivity | Accuracy | کرنل کوادراتیک بر روی فضای ویژگی انشعاب‌ها |
|-----------|-----------|-------------|----------|--|
| ۰,۸۳      | ۰,۸۲      | ۰,۸۱        | ۰,۸۲     | تک هسته‌ای                                 |
| ۰,۹۰      | ۰,۸۹      | ۰,۸۸        | ۰,۹۱     | دو هسته‌ای                                 |
| ۰,۹۴      | ۰,۹۵      | ۰,۹۵        | ۰,۹۴     | سه هسته‌ای                                 |

جدول (۵) معیارهای ارزیابی محاسبه شده بر حسب کرنل کوادراتیک در مدل پیشنهادی بر مجموعه داده‌ی آلمانی با ۲۰ ویژگی

| F-measure | Precision | Sensitivity | Accuracy | کرنل کوادراتیک بر روی فضای ویژگی انشعاب‌ها |
|-----------|-----------|-------------|----------|--|
| ۰,۸۲      | ۰,۸۴      | ۰,۸۲        | ۰,۸۴     | تک هسته‌ای                                 |
| ۰,۹۲      | ۰,۹۲      | ۰,۹۱        | ۰,۹۳     | دو هسته‌ای                                 |
| ۰,۹۵      | ۰,۹۶      | ۰,۹۶        | ۰,۹۶     | سه هسته‌ای                                 |

همان گونه که مشخص است با افزایش تعداد ابعاد کرنل معیارهای ارزیابی نیز روند صعودی دارند. پس قدرت الگوریتم در شناسایی نیز افزایش یافته است. با توجه به جداول فوق که مربوط به انجام آزمایش بر روی مجموعه داده‌ی استرالیای با ۱۴ ویژگی و مربوط به انجام آزمایش بر روی مجموعه داده‌ی آلمانی با ۲۰ ویژگی است، دریافتیم که با افزایش ابعاد کرنل از خطی به درجه سه در مدل دقت تشخیص الگوریتم افزایش می‌یابد. نکته‌ی قابل توجه دیگر آن است با افزایش تعداد ویژگی‌های مسئله استفاده از کوادراتیک نتایج مطلوب‌تری به همراه داشته است.

### پارامتر m

یکی دیگر از پارامترهایی که در مدل پیشنهادی قابل تنظیم می‌باشد، با نام m شناخته می‌شود. یادآور می‌شویم که این پارامتر به معنی تعداد ویژگی‌هایی است که در هر گره به طور تصادفی از مجموعه‌ی کل ویژگی‌های مسئله استخراج شده و با استفاده از این زیرمجموعه‌ی انتخاب شده بهترین ویژگی برای شکست انتخاب می‌شود. در طول زمان ساخت مدل، این پارامتر برای تمامی انشعابات ثابت در نظر گرفته می‌شود و مقادیر معمولی که برای آن انتخاب می‌گردد،  $\sqrt{nVariable}$  و یا  $\log(nVariable)$  خواهد بود. برای مشاهده‌ی تأثیر این پارامتر و انتخاب بهترین مقدار ممکن برای آن، آزمایش زیر را اجرا کرده‌ایم:

### آزمایش ۳: تأثیر پارامتر m

در این آزمایش مدل پیشنهادی را با شرایط زیر اجرا کرده‌ایم:

۱. تعداد انشعاب‌ها = ۱۲۰
  ۲. معیار تعیین ویژگی شکست = شاخص جینی
  ۳. پارامتر m متغیر بوده و سه مقدار مجزا به خود می‌گیرد. ۱. جذر تعداد کل ویژگی‌ها در هر مجموعه داده. ۲. لگاریتم تعداد کل ویژگی‌ها در هر مجموعه داده. ۳. تعداد کل ویژگی‌ها در هر مجموعه داده.
- در جداول ۵ و ۶ نتایج ارزیابی الگوریتم در این سه حالت مختلف و بر روی دو مجموعه داده‌ی موجود، داده نشان داده شده است.

جدول (۵) معیارهای ارزیابی مدل پیشنهادی بر حسب مقادیر مختلف پارامتر m، مجموعه داده‌ی استرالیایی با ۱۴ ویژگی

| F-measure | Sensitivity | Accuracy | پارامتر m          |
|-----------|-------------|----------|--------------------|
| ۰,۹۳      | ۰,۹۳        | ۰,۹۳     | $\sqrt{nVariable}$ |
| ۰,۹۳      | ۰,۹۴        | ۰,۹۲     | $\log(nVariable)$  |
| ۰,۹۳      | ۰,۹۳        | ۰,۹۳     | nVariable          |

جدول (۶) معیارهای ارزیابی مدل پیشنهادی بر حسب مقادیر مختلف پارامتر m، مجموعه داده‌ی آلمانی با ۲۰ ویژگی

| F-measure | Sensitivity | Accuracy | پارامتر m          |
|-----------|-------------|----------|--------------------|
| ۰,۹۲      | ۰,۹۳        | ۰,۸۹     | $\sqrt{nVariable}$ |
| ۰,۹۲      | ۰,۹۲        | ۰,۸۹     | $\log(nVariable)$  |

|      |      |     |           |
|------|------|-----|-----------|
| ۰,۹۳ | ۰,۹۳ | ۰,۹ | nVariable |
|------|------|-----|-----------|

باتوجه به نتایج فوق می توان به این نکته پی برد که در مسئلهی حاضر و با مجموعه داده های در دسترس کاهش تعداد ویژگی ها به مقادیر  $\sqrt{\text{nVariable}}$  و  $\log(\text{nVariable})$  نتایج امیدبخشی را در پی ندارد.

### معیار تعیین نقطه ی شکست

یکی از مهم ترین پارامترهای موجود در الگوریتم پیشنهادی و کلیه ی روش هایی که بر مبنای درخت تصمیم عمل می کنند، انتخاب معیاری مناسب برای تعیین نقطه ی شکست در زمان ساخت مدل خواهد بود. از آنجایی که هسته ی اصلی مدل پیشنهادی، درختان تصمیم موجود در آن است، پس انتخاب بهترین معیار برای تعیین نقطه یا همان ویژگی برای بخش بندی داده ها نکته ی مهمی محسوب می شود. همان طور که ذکر شد، معیارهای مختلفی برای این کار وجود دارد که به تفصیل توصیف شدند. برای مشاهده ی تأثیر این پارامتر و انتخاب بهترین معیار انتخاب نقطه ی شکست، آزمایش سوم به شکل زیر انجام شده است:

آزمایش ۴: تأثیر پارامتر تعیین ویژگی شکست

در این آزمایش مدل پیشنهادی را با شرایط زیر اجرا کرده ایم:

- تعداد انشعابات = ۱۲۰

۴.  $m = \text{nVariable}$  به معنی آن که پارامتر  $m$  برابر با تعداد کل ویژگی های موجود است. باتوجه به دو مجموعه داده ی در دسترس، با به کارگیری مجموعه داده ی اول این پارامتر مقدار ۱۴ و با به کارگیری مجموعه داده ی دوم، مقدار ۲۰ را خواهد داشت.

- معیار انتخاب نقطه ی شکست = متغیر بوده و هر بار یکی از معیارهای بهره ی اطلاعاتی، نسبت بهره و شاخص جینی انتخاب و آزمایش می شود.

در جداول ۷ تا ۱۱ نتایج ارزیابی الگوریتم در این سه حالت مختلف و بر روی دو مجموعه داده ی موجود، داده نشان داده شده است.

جدول (۷) معیارهای ارزیابی مدل پیشنهادی بر حسب معیارهای مختلف تعیین ویژگی بخش بندی داده ها، مجموعه داده ی استرالیایی با

۱۴ ویژگی

| F-measure | Sensitivity | Accuracy |               |
|-----------|-------------|----------|---------------|
| ۰,۷۲      | ۰,۶         | ۰,۹      | بهره اطلاعاتی |
| ۰,۵۵      | ۰,۴         | ۰,۷۲     | نسبت بهره     |
| ۰,۹۳      | ۰,۹۵        | ۰,۹۳     | شاخص جینی     |

جدول (۸) معیارهای ارزیابی مدل پیشنهادی بر حسب معیارهای مختلف تعیین ویژگی بخش بندی داده ها، مجموعه داده ی آلمانی با ۲۰

ویژگی



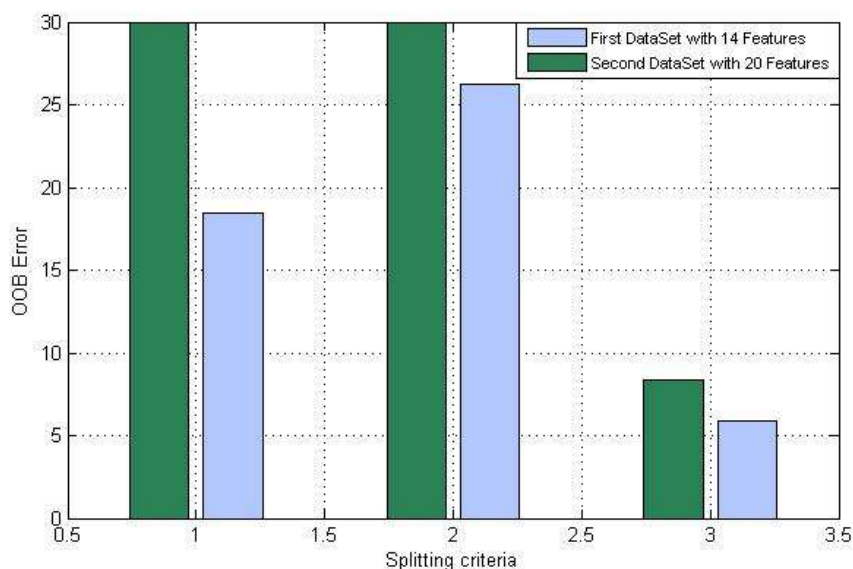
| F-measure | Sensitivity | Accuracy |               |
|-----------|-------------|----------|---------------|
| ۰,۸۵      | ۰,۸۶        | ۰,۷۵     | بهره اطلاعاتی |
| ۰,۷۶      | ۰,۸         | ۰,۷      | نسبت بهره     |
| ۰,۹۵      | ۰,۹۳        | ۰,۹۶     | شاخص جینی     |

جدول (۹) معیارهای ارزیابی مدل پیشنهادی بر حسب معیارهای مختلف تعیین ویژگی بخش‌بندی داده‌ها، مجموعه‌داده‌ی استرالیایی با ۱۴ ویژگی

| Precision | خطای OOB |               |
|-----------|----------|---------------|
| ۰,۹۳      | ۱۸,۴٪    | بهره اطلاعاتی |
| ۰,۹۳      | ۲۶,۲٪    | نسبت بهره     |
| ۰,۹۵      | ۴,۶٪     | شاخص جینی     |

جدول (۱۰) معیارهای ارزیابی مدل پیشنهادی بر حسب معیارهای مختلف تعیین ویژگی بخش‌بندی داده‌ها، مجموعه‌داده‌ی آلمانی با ۲۰ ویژگی

| Precision | خطای OOB |               |
|-----------|----------|---------------|
| ۰,۹       | ۳۰٪      | بهره اطلاعاتی |
| ۰,۸۸      | ۳۰٪      | نسبت بهره     |
| ۰,۹۵      | ۴,۲٪     | شاخص جینی     |



شکل (۱) خطای OOB بر حسب معیارهای مختلف تعیین ویژگی بخش‌بندی داده‌ها

باتوجه به نتایج به دست آمده از این آزمایش در مسئله‌ی حاضر و باتوجه به مجموعه داده‌های در دسترس، مشاهده می‌شود که استفاده از معیار شاخص جینی میزان دقت بیش تر شده است و نتایج بهتری را در سایر معیارهای ارزیابی نسبت به دو معیار دیگر فراهم آورده است. همچنین مشاهده می‌شود که با استفاده از معیار شاخص جینی میزان صحت بیش تر و خطای OOB کمتر شده است.

### محدودیت‌های روش پیشنهادی

از جمله محدودیت‌های موجود افزایش تعداد تراکنش‌های بانکی و مواجهه با حجم بسیار زیادی از اطلاعات به‌روز شده در هر لحظه است.

محدودیت دیگر این است که چون عملیات بانکی نیازمند پاسخ بلادرنگ می‌باشند، زمان بسیار محدودی برای بررسی تراکنش و مسدود کردن آن در صورت مشکوک بودن می‌باشد.

یکی از محدودیت‌های دیگر در این زمینه اطلاعات غیرساختاریافته است. به همین دلیل روش شناسایی تقلب باید قادر باشد این اطلاعات را به‌درستی سازماندهی کند تا بتواند تحلیل و تشخیص درستی ارائه دهد.

### بحث و مقایسه

در بخش قبل با انجام آزمایش‌های مختلف سعی در یافتن بهترین مقادیر برای پارامترهای مختلف این الگوریتم داشتیم؛ بنابراین باتوجه به این آزمایش اول دریافتیم که مناسب‌ترین تعداد انشعاب برای مدلی که داده‌های تست را دسته‌بندی کند، در هر دو مجموعه داده‌ی در دسترس، ۱۸۰ بوده است. این ادعا باتوجه به نتایج موجود در جداول ۱ و ۲ ثابت می‌شود.

با انجام آزمایش دوم، سعی در یافتن بهترین مقدار برای پارامتر  $m$  در این الگوریتم داشتیم. باتوجه به نتایج به دست آمده در جداول ۳ و ۴، اگر این پارامتر را به تعداد ویژگی‌های موجود در مجموعه داده مقارنه کنیم، بهترین شرایط در معیارهای ارزیابی حاصل می‌شود. در واقع باتوجه به مجموعه داده‌های در دسترس و محدود بودن تعداد ویژگی در آن‌ها، کاهش مقدار این پارامتر به جذر تعداد کل ویژگی‌ها و یا لگاریتم آن، باعث کاهش کارایی الگوریتم شده و در واقع مدل حاصل به‌خوبی آموزش داده نشده و با تمامی فضای مسئله منطبق نمی‌گردد.

در نهایت، هدف از انجام آزمایش سوم تعیین بهترین معیار انتخاب ویژگی شکست در گره‌ها بود. باتوجه به نتایج نشان داده شده در جداول ۴ و ۵ دریافتیم که استفاده از معیار شاخص جینی در هر دو مجموعه داده عملکرد مناسبی دارد.

در این بخش برای نمایش قدرت و کارایی مدل پیشنهادی، این الگوریتم را با الگوریتم درخت تصمیم و ماشین بردار پشتیبان و شبکه عصبی و جنگل تصادفی معمولی که از معروف‌ترین الگوریتم‌های داده‌کاوی محسوب می‌شوند و همچنین روش ارائه شده در مقاله اوش و همکاران ارزیابی می‌کنیم. روش ارائه شده در این مطالعه یک شبکه عصبی را برای شناسایی تقلب در کارت‌های اعتباری و آزمایش آن بر روی مجموعه داده‌ی آلمانی با ۲۰ ویژگی که ما در این پروژه از آن استفاده کرده‌ایم، به کار گرفته است. در این تحقیق، شبکه‌ی عصبی با الگوریتم شبیه‌سازی ذوب فلزات آموزش داده شده است. به معنی آن که اوزان موجود در شبکه با این الگوریتم آموزش داده شده‌اند. در واقع به جای استفاده از روش آموزش گرادیان نزولی که تاکنون بسیار مورد استفاده قرار گرفته است و در اکثر مواقع در کمینه‌ی محلی گیر می‌افتد، از این الگوریتم هیوریستیک که توان فرار از کمینه‌های محلی را دارد، بهره گرفته و سعی در افزایش دقت شبکه عصبی حاصل داشته است. میزان دقت (Accuracy) به دست آمده با به کارگیری این روش، ۸۹٫۶٪ گزارش شده

است. با اجرای الگوریتم پیشنهادی با شرایط موجود در جدول ۱۱، به دقت تشخیص ۹۶٪ در مجموعه داده‌ی آلمانی با ۲۰ ویژگی و دقت تشخیص ۹۵٪ در مجموعه داده‌ی استرالیایی با ۱۴ ویژگی، دست یافته‌ایم.

جدول (۱۱) پارامترهای تنظیم شده برای ساخت مدل و ارزیابی آن

|           |                        |
|-----------|------------------------|
| مقدار     | پارامتر                |
| ۱۸۰       | تعداد انشعابات         |
| nVariable | پارامتر m              |
| شاخص جینی | معیار تعیین ویژگی شکست |

همچنین برای ساخت مدل، در هر دو مجموعه داده، از ۷۰ درصد داده‌ها برای آموزش و از ۳۰ درصد باقیمانده برای تست آن استفاده کرده‌ایم. معیار انتخاب نقطه‌ی شکست در انشعاب نیز شاخص جینی انتخاب شده است. چراکه مشاهده کردیم این معیار نسبت به بقیه نتایج بهتری در پی دارد.

برای ساخت مدل ماشین بردار پشتیبان نیز از ۷۰ درصد داده‌ها برای آموزش و از ۳۰ درصد دیگر برای تست آن استفاده کردیم. همچنین از تابع خطی به عنوان تابع kernel در این مدل استفاده کرده‌ایم. سپس هر دو الگوریتم را بر روی هر دو مجموعه داده اجرا کرده، معیارهای ارزیابی را برای هر دو مدل محاسبه و در جداول ۱۲ تا ۱۵ با الگوریتم پیشنهادی مقایسه شده است. همان‌طور که مشخص است مدل پیشنهادی این پژوهش در مقایسه با این روش‌ها کاملاً بهتر عمل کرده و امید بخش بوده و استفاده از آن دقت عملکرد بهتری را فراهم آورده است.

جدول (۱۲) معیارهای ارزیابی الگوریتم پیشنهادی و درخت تصمیم، ماشین بردار پشتیبان، شبکه عصبی و جنگل تصادفی معمولی مجموعه داده استرالیایی با ۱۴ ویژگی

| F-measure | Sensitivity | Accuracy | الگوریتم            |
|-----------|-------------|----------|---------------------|
| ۰٫۹۵      | ۰٫۹۵        | ۰٫۹۵     | مدل پیشنهادی        |
| ۰٫۹۱      | ۰٫۹         | ۰٫۹      | درخت تصمیم          |
| ۰٫۸۱      | ۰٫۸۴        | ۰٫۷۵     | ماشین بردار پشتیبان |
| ۰٫۹۱      | ۰٫۹۲        | ۰٫۹۱     | شبکه عصبی           |
| ۰٫۹۴      | ۰٫۹۳        | ۰٫۹۴     | جنگل تصادفی معمولی  |

جدول (۱۳) معیارهای ارزیابی الگوریتم پیشنهادی و درخت تصمیم، ماشین بردار پشتیبان، شبکه عصبی و جنگل تصادفی معمولی مجموعه داده آلمانی با ۲۰ ویژگی

| الگوریتم            | Accuracy | Sensitivity | F-measure |
|---------------------|----------|-------------|-----------|
| مدل پیشنهادی        | ۰,۹۶     | ۰,۹۶        | ۰,۹۵      |
| درخت تصمیم          | ۰,۷۷     | ۰,۷۵        | ۰,۸۲      |
| ماشین بردار پشتیبان | ۰,۸۶     | ۰,۹         | ۰,۹۱      |
| شبکه عصبی           | ۰,۸۷     | ۰,۹         | ۰,۸۹      |
| جنگل تصادفی معمولی  | ۰,۹۲     | ۰,۹۲        | ۰,۹۳      |

جدول (۱۴) معیارهای ارزیابی الگوریتم پیشنهادی و درخت تصمیم، ماشین بردار پشتیبان، شبکه عصبی و جنگل تصادفی معمولی مجموعه داده استرالیایی با ۱۴ ویژگی

| الگوریتم            | خطای OOB | Precision |
|---------------------|----------|-----------|
| مدل پیشنهادی        | ۴,۲٪     | ۰,۹۶      |
| درخت تصمیم          | -        | ۰,۹۲      |
| ماشین بردار پشتیبان | -        | ۰,۷۱      |
| شبکه عصبی           | -        | ۰,۹۲      |
| جنگل تصادفی معمولی  | ۸,۴٪     | ۰,۹۴      |

جدول (۱۶) معیارهای ارزیابی الگوریتم پیشنهادی و درخت تصمیم، ماشین بردار پشتیبان، شبکه عصبی و جنگل تصادفی معمولی مجموعه داده آلمانی با ۲۰ ویژگی

| الگوریتم            | خطای OOB | Precision |
|---------------------|----------|-----------|
| مدل پیشنهادی        | ۲,۸٪     | ۰,۹۶      |
| درخت تصمیم          | -        | ۰,۸۹      |
| ماشین بردار پشتیبان | -        | ۰,۸۳      |
| شبکه عصبی           | -        | ۰,۹۱      |
| جنگل تصادفی معمولی  | ۱۲,۶٪    | ۰,۹۲      |

## نتیجه‌گیری

امروزه استفاده از کارت‌های بانکی در سطوح گسترده‌ای از تعاملات تجاری مطرح است. هرچند این تحولات گامی بزرگ در جهت کارایی و سهولت دسترسی است، معایبی نیز به همراه دارد که مهم‌ترین آن آسیب‌پذیری نسبت به فعالیت‌های متقلبان است؛ به همین دلیل پژوهشگران همواره دنبال ارائه روشی جدید و مؤثر و کارا برای شناسایی به هنگام و یا پیشگیری از وقوع تقلب در کارت‌های بانکی هستند.

در این مقاله برای یافتن بهترین مقادیر برای پارامترهای مورد تنظیم در این الگوریتم آزمایشاتی طراحی و اجرا کردیم. پارامترهای آزاد در این الگوریتم عبارت‌اند از تعداد انشعابات، درجه کوادراتیک، پارامتر  $m$  که دامنه‌ی انتخاب ویژگی شکست را معین می‌کند و معیار تعیین ویژگی شکست. با انجام چهار آزمایش مقادیر این پارامترها برآورد شده و نتایج حاصل از آن‌ها ارائه گردید. سپس الگوریتم پیشنهادی را ارزیابی کردیم. با انجام آزمایش نهایی و ارزیابی الگوریتم دریافتیم که مدل ارائه شده در این پژوهش با قدرت بیش‌تری عمل کرده و کاملاً رضایت‌بخش و نتایج خوبی در پی داشته است. حاصل پژوهش انجام شده مدلی برای شناسایی متقلبین در سیستم‌های بانکداری آنلاین بر مبنای تراکنش‌های کارت‌های اعتباری بوده است که عملکرد آن در دسته‌بندی تراکنش‌ها و متقلبین مناسب به نظر می‌رسد.

## پیشنهاد‌های آینده

پژوهش حاضر به طور خاص بر روی تشخیص متقلبین در سیستم‌های بانکداری آنلاین معطوف شده است. برای ارزیابی روش ارائه شده در این مقاله، می‌توان این مدل را در سایر زمینه‌های مشابه تحقیقاتی مثل شناسایی تقلب در بازار بورس و خرید و فروش سهام که در آن نیز تعاملات مهمی انجام می‌شود و به دلیل طرف‌داران زیادی که سالانه به خود جذب می‌کند می‌تواند بسیار مورد توجه افراد متقلب برای فریب افراد تازه‌وارد، قرار گیرد. می‌توان در پژوهش‌های بعدی این مدل را در زمینه‌های دیگر به کار گرفته و آن را ارزیابی کرد.

## منابع

بنائی، هادی، خوش‌نیت، حسام. (۱۳۹۶). نقش و کاربرد هوش عملیاتی و داده‌کاوی در کشف تقلب برخط. ششمین همایش ملی تجارت و اقتصاد الکترونیک. همایش تخصصی امنیت و اعتماد.

حاتمی‌راد، علی، شهریاری، حمیدرضا. (۱۳۹۷). روش‌ها و راهکارهای شناسایی تقلب در بانک‌داری الکترونیک. فصل‌نامه تازه‌های اقتصاد، سال نهم، شماره ۱۳۴، صص ۲۱۹ تا ۲۲۸.

قلی‌پور سلیمانی، علی، ایمانی، سهیلا. (۱۴۰۰). سیر تکنولوژی در بانکداری. دو ماهنامه مدیریت، شماره ۱۵۹، صص ۲۲ تا ۲۵. وثوق، ملیحه، تقوی‌فرد، محمدتقی و البرزی، محمود. (۱۳۹۸). شناسایی تقلب در کارت‌های بانکی با استفاده از شبکه‌های عصبی مصنوعی. فصل‌نامه علمی-پژوهشی مدیریت فناوری اطلاعات دانشگاه تهران، دوره ۶، شماره ۴، صص ۷۲۱-۷۴۶.

Ata, H. A., & Seyrek, I. H. (۲۰۰۹). THE USE OF DATA MINING TECHNIQUES IN DETECTING FRAUDULENT FINANCIAL STATEMENTS: AN APPLICATION ON MANUFACTURING FIRMS. *Suleyman Demirel University Journal of Faculty of Economics & Administrative Sciences*, 14(۲).

Bahnsen, A. C., Aouada, D., & Ottersten, B. (۲۰۱۵). Example-dependent cost-sensitive decision trees. *Expert Systems with Applications*, 42(۱۹), ۶۶۰۹-۶۶۱۹.

- Bahnsen, A. C., Aouada, D., Stojanovic, A., & Ottersten, B. (۲۰۱۶). Feature engineering strategies for credit card fraud detection. *Expert Systems with Applications*, 51, ۱۳۴-۱۴۲.
- Bansal, M., & Sharma, D. (۲۰۲۱). A novel multi-view clustering approach via proximity-based factorization targeting structural maintenance and sparsity challenges for text and image categorization. *Information Processing & Management*, 58(۴), ۱۰۲۰۴۶.
- Bhattacharyya, S., Jha, S., Tharakunnel, K., & Westland, J. C. (۲۰۱۱). Data mining for credit card fraud: A comparative study. *Decision support systems*, 50(۳), ۶۰۲-۶۱۳.
- Bose, I., & Mahapatra, R. K. (۲۰۰۱). Business data mining—a machine learning perspective. *Information & management*, 39(۳), ۲۱۱-۲۲۰.
- Breiman, L. (۲۰۱۱). Random forests. *Machine learning*, 45, ۰-۳۲.
- Carta, S., Fenu, G., Recupero, D. R., & Saia, R. (۲۰۱۹). Fraud detection for E-commerce transactions by employing a prudential Multiple Consensus model. *Journal of Information Security and Applications*, 46, ۱۳-۲۲.
- Chandra, V., & Singh, P. (۲۰۱۴). Fuzzy Based High Blood Pressure Diagnosis. *International Journal of Advanced Research in Computer Science & Technology (IJARCST)*, 2(۲), ۲۳۴۷-۸۴۴۶.
- Dal Pozzolo, A., Caelen, O., Le Borgne, Y. A., Waterschoot, S., & Bontempi, G. (۲۰۱۴). Learned lessons in credit card fraud detection from a practitioner perspective. *Expert systems with applications*, 41(۱۰), ۴۹۱۰-۴۹۲۸.
- Dreżewski, R., Sepielak, J., & Filipkowski, W. (۲۰۱۰). The application of social network analysis algorithms in a system supporting money laundering detection. *Information Sciences*, 295, ۱۸-۳۲.
- Eberle, W., & Holder, L. (۲۰۰۷). Anomaly detection in data represented as graphs. *Intelligent Data Analysis*, 11(۶), ۶۶۳-۶۸۹.
- Fang, W., Li, X., Zhou, P., Yan, J., Jiang, D., & Zhou, T. (۲۰۲۱). Deep learning anti-fraud model for internet loan: where we are going. *IEEE Access*, 9, ۹۷۷۷-۹۷۸۴.
- Hirshman, J., Huang, Y., & Macke, S. (۲۰۱۳). Unsupervised approaches to detecting anomalous behavior in the bitcoin transaction network. *Technical report, Stanford University*.
- JYeonkook J. Kim, Bok Baik b, Sungzoon Cho, “Detecting financial misstatements with fraud intention using multi-class cost-sensitive learning”, *Expert Systems With Applications*, Vol. ۶۲, Pages ۳۲-۴۳, (۲۰۱۹).
- Moreira, M. Â. L., Junior, C. D. S. R., de Lima Silva, D. F., de Castro Junior, M. A. P., de Araújo Costa, I. P., Gomes, C. F. S., & dos Santos, M. (۲۰۲۲). Exploratory analysis and implementation of machine learning techniques for predictive assessment of fraud in banking systems. *Procedia Computer Science*, 214, ۱۱۷-۱۲۴.
- Nazeer, I., Prasad, K. D. V., Bahadur, P., Bapat, V., & MJ, K. (۲۰۲۳). Synchronization of AI and Deep Learning for Credit Card Fraud Detection. *International Journal of Intelligent Systems and Applications in Engineering*, 11(۰۵), ۰۲-۰۹.
- Patidar, R., & Sharma, L. (۲۰۱۱). Credit card fraud detection using neural network. *International Journal of Soft Computing and Engineering (IJSCE)*, 1(۳۲-۳۸).

- Phua, C., Lee, V., Smith, K., & Gayler, R. (۲۰۱۰). A comprehensive survey of data mining-based fraud detection research. *arXiv preprint arXiv:1009.6119*.
- Salchenberger, L. M., Cinar, E. M., & Lash, N. A. (۱۹۹۲). Neural networks: A new tool for predicting thrift failures. *Decision Sciences*, 23(۴), ۸۹۹-۹۱۶.
- Shen, A., Tong, R., & Deng, Y. (۲۰۰۷, June). Application of classification models on credit card fraud detection. In *2007 International conference on service systems and service management* (pp. ۱-۴). IEEE.
- Van Vlasselaer, V., Bravo, C., Caelen, O., Eliassi-Rad, T., Akoglu, L., Snoeck, M., & Baesens, B. (۲۰۱۵). APATE: A novel approach for automated credit card transaction fraud detection using network-based extensions. *Decision Support Systems*, 75, ۳۸-۴۸.
- Wang, X., Wang, X., Wilkes, M., Wang, X., Wang, X., & Wilkes, M. (۲۰۲۱). A k-nearest neighbour spectral clustering-based outlier detection technique. *New Developments in Unsupervised Outlier Detection: Algorithms and Applications*, ۱۴۷-۱۷۲.