

کشف تقلبات در کارت های اعتباری با استفاده از شبکه عصبی و درخت تصمیم

رضا روشنی

۱. موسسه آموزش عالی غیر انتفاعی لامعی گرگانی
۲. گروه مهندسی کامپیوتر، دانشگاه ملی مهارت، تهران، ایران

معصومه کیائی

۱. موسسه آموزش عالی غیر انتفاعی لامعی گرگانی

چکیده

استفاده از کارت های اعتباری برای جلوگیری از حوادث ناشی از به همراه داشتن پول نقد رواج پیدا کرده است. اما این موضوع نه تنها باعث عقب نشینی افراد سودجو نشده، بلکه راه های جدید با خطر کمتری را پیش روی آن ها گشوده است. از بررسی آمارهای بانکی و سایر متولیان کارت های اعتباری چنین بر می آید که تقلب در کارت های اعتباری بطور فزاینده ای در حال افزایش است و با توجه به حجم بالای تراکنش های صورت گرفته، نیاز به ابزار مناسبی برای کشف تقلب در داده ها احساس می شود. داده کاوی علم جدیدی است که در این زمینه مطرح شده، بطوریکه با استفاده از تکنیک ها و الگوریتم های آن می توان در این حجم انبوه داده به پردازش پرداخته و موارد مشکوک را جستجو و شناسایی نمود. در این پژوهش به منظور ارائه روشی در زمینه کشف تقلبات صورت گرفته در تراکنش های کارت های اعتباری به بررسی دو تکنیک پر کاربرد داده کاوی یعنی درخت تصمیم و شبکه عصبی بر روی دو مجموعه داده پرداخته می شود و با استفاده از ابزار داده کاوی متلب این الگوریتم ها در دو مرحله بر روی مجموعه داده های انتخابی پیاده سازی خواهند شد و نتایج آن ها مورد بررسی و مقایسه قرار خواهد گرفت. نتایج پیاده سازی نشان می دهد که الگوریتم درخت تصمیم نسبت به شبکه عصبی درصد بیشتری از معیارها را به خود اختصاص داده است.

واژگان کلیدی: کارت اعتباری، کشف تقلب، داده کاوی، درخت تصمیم، شبکه عصبی.

۱. مقدمه

یکی از بزرگترین سازمان‌هایی که با مشتریان به صورت کاملاً مستقیم در تعامل هستند بانک‌ها می‌باشند و با توسعه روز افزون تکنولوژی، و جایگزینی سیستم‌های نوین در تمامی عرصه‌ها، بجای ساختار سنتی و کهنه، نظام بانکداری نیز پیرو این ارتقاء، سعی بر آن دارد که در حوزه فن آوری، مطلوبیت ممکن را برای مشتریان فراهم آورد. یکی از ابزارهای مهم در این عرصه، کارت اعتباری^۱ است که حجم بالای نقل و انتقالات مالی افراد حقیقی و حقوقی را شامل شده است. کارت‌های اعتباری می‌توانند زندگی را آسان‌تر کنند. آن‌ها به مشتریان این امکان را می‌دهند که از اعتبار خود در زمان، مکان و میزان دلخواه، بدون حمل پول نقد استفاده کنند. چالش‌های امنیتی مختلف به دلیل محبوبیت استفاده از کارت‌های اعتباری در حال افزایش است و همین موضوع باعث تشدید کلاهبرداری با هدف به دست آوردن منافع مالی غیر مجاز شده است. تقلب کارت اعتباری به وضعیتی اطلاق می‌شود که متقلب برای رفع نیازهای خود از کارت اعتباری دیگری استفاده می‌کند در حالی که صاحب آن کارت اعتباری از این امر آگاه نیست. روش‌های مختلفی برای کشف تقلب در کارت‌های اعتباری معرفی شده‌اند. با این حال، کلاهبرداران نیز دائماً روش‌های جدیدی برای سرقت اطلاعات ارائه می‌دهند [۱۱]. در حال حاضر اجرای یک سیستم کشف تقلب کارآمد چالش اصلی این موضوع است. روش‌های جدیدی برای کشف تقلب وجود دارد؛ اما از داده‌کاوی^۲ که هدف آن استخراج اطلاعات مورد علاقه از مجموعه‌ی گسترده‌ای از داده‌هاست به طور گسترده‌ای به عنوان ابزار تصمیم‌گیری فعال، استفاده می‌شود.

۱.۱. کارت اعتباری

در فرهنگ آکسفورد، کارت اعتباری چنین تعریف شده است: "کارتی که بانک یا غیر آن در اختیار مشتری قرار می‌دهد و او می‌تواند کالا و خدمات مورد نیاز خود را دریافت کند." بنابراین کارت اعتباری کارت‌ای است که بانک‌ها و موسسات مالی و اعتباری خصوصی در اختیار مشتریان خود قرار می‌دهند تا او بتواند بهای کالا و خدمات مورد نیاز خود را بدون دغدغه پول فیزیکی با استفاده از دستگاه‌های مربوطه نظیر پایانه فروش^۳ (POS)، دستگاه خودپرداز (ATM)^۴ و یا اینترنت پرداخت کند [۱۲].

۱.۲. تقلب کارت‌های اعتباری

معمولاً، تقلب به استفاده غیرمجاز و غیر قانونی از تسهیلات اعتباری یک حساب قانونی اشاره دارد. تقلب می‌تواند به عنوان یک فعالیت بزهکارانه تعریف شود که با نمایش‌های غیرقانونی برای بدست آوردن یک امتیاز ناعادلانه سروکار دارد [۱۳]. امروزه دزدی با کارت اعتباری بسیار امن‌تر، راحت‌تر و پرسودتر از دزدی با اسلحه است. صدمات به تجارت‌ها و بانک‌ها از طریق تقلب کارت‌های اعتباری با افزایش روز افزونی رو به روست. بیان کردن وسعت تقلب کارت اعتباری مشکل است، چون ارقام در طول زمان تغییر می‌کنند (محتملاً رشد می‌کنند). تا حدی که اغلب اوقات شرکت‌ها از منتشر ساختن ارقام تقلب بی‌میل هستند، به دلیل اینکه آن‌ها افکار عمومی را به وحشت می‌اندازند.

۱.۳. انواع تقلب در کارت‌های بانکی و اعتباری [۱] و [۲]

تقلب در کارت‌های اعتباری به طرق مختلفی امکان‌پذیر می‌باشد که در ذیل به آن‌ها اشاره شده است.

¹ Credit Card

² Data Mining

³ Point of Sale

⁴ Automated teller machine

• تقلب ورشکستگی

در این روش فرد با آگاهی از اینکه توان پرداخت بدهی خود را ندارد باز هم با کارت اعتباری خود معاملات انجام می‌دهد. این بدهی‌ها برای بانک از طرف فرد خاطی غیرقابل وصول است و خود بانک متحمل ضرر می‌شود. تنها راه جلوگیری از این تقلب‌ها بررسی وضعیت اعتباری دارنده حساب‌ها از طریق موسسات سنجش اعتباری است.

• کارت گمشده یا دزدی

در این روش کارت فردی مفقود یا دزدیده شده است و تا زمان بلوکه شدن حساب، از کارت برداشت‌هایی صورت می‌گیرد.

• تقلب درخواست

در این نوع تقلب، شخص به سه روش با دادن اطلاعات نادرست، کارت اعتباری دریافت می‌کند:

۱. هویت جعلی : شخصی متقلب با بدست آوردن اطلاعات غیرقانونی از شخص دیگر به افتتاح حساب کاربری با نام آن‌ها دست می‌زند.

۲. تقلب مالی : شخص به هنگام افتتاح حساب درباره وضعیت مالی خود اطلاعات نادرست ارائه می‌دهد.

۳. نرسیدن بسته پستی : بسته پستی حاوی کارت ارسال شده برای فرد در بین راه دزدیده می‌شود.

• کنترل حساب

شخص متقلب در این روش اطلاعات حساب مشتری خاصی را بدست آورده و می‌تواند حساب او را کنترل نماید.

• جعل کارت

در این روش مجرمان با روش‌های مختلفی مثل پاک کردن نوارهای مغناطیسی کارت بوسیله میدان مغناطیسی قوی یا دستکاری اطلاعات و مشخصات کارت و تطبیق آن با اطلاعات یک کارت قانونی یا کپی کردن داده‌های اصلی نوار مغناطیسی یک کارت بصورت الکترونیکی بر روی کارت دیگر یا حتی استفاده از کارت‌های سفید که نیاز به تایید برخی پایانه‌ها ندارد، به جعل کارت می‌پردازند.

• تبانی فروشنده

در این روش فروشنده خاطی اطلاعات مشتریان خود را در اختیار مجرمان قرار می‌دهد.

• تقلب اینترنتی

این نوع تقلب انواع مختلفی دارد:

۱. وب سایت‌های مشابه : شبیه سازی صفحه اینترنتی با فرم خرید از یک فروشگاه معتبر. پس از انجام خرید، اطلاعات کارت خریدار در اختیار مجرم است.

۲. وب سایت تقلبی فروش: سایت‌هایی راه اندازی شده که در آن کالاهای تقلبی با درصد تخفیف‌های بالا به حراج گذاشته می‌شود که ظاهر یک سایت حراج قانونی را دارد.

مولد کارت اعتباری: برنامه‌های کامپیوتری هستند که می‌توانند شماره انواع کارت‌های اعتباری مانند ویزا و مسترکارت را تولید کنند.

۱/۴. معرفی داده‌کاوی

داده‌کاوی عبارت است از کشف ارتباطات یا الگوهای معنی دار در یک پایگاه داده بسیار بزرگ و پیچیده. به عبارتی دیگر داده‌کاوی یک فرآیند خودکار یا نیمه خودکار استخراج دانش در قالب الگوهای پنهان از مجموعه اطلاعات ورودی می‌باشد. در داده‌کاوی ورودی عمدتاً بسیار حجیم و پردازش دستی آن ناممکن است. در نتیجه نتایج حاصل از داده‌کاوی، با روش‌های سنتی پردازش اطلاعات، قابل دستیابی نمی‌باشد [۳].

اکتشاف دانش در پایگاه‌های داده یک فرآیند پیچیده است که شامل سه مرحله کلی آماده سازی داده‌ها، یادگیری مدل و ارزیابی و تفسیر مدل است. این سه مرحله در تمام پروژه‌های داده‌کاوی وجود دارند و اساس داده‌کاوی بر این سه مرحله استوار است.

آماده سازی داده اولین و مهم‌ترین مرحله در فرآیند داده‌کاوی است و هدف آن تأمین ورودی مناسب برای مرحله حیاتی یادگیری مدل است. به طور کلی به مجموعه عملیاتی که منجر به تولید مجموعه‌ای از داده‌های پالایش شده قابل کاوش خواهد شد اصطلاحاً آماده سازی داده می‌گویند و این عملیات عبارت اند از استخراج داده و پیش پردازش داده [۴].

در مرحله یادگیری مدل، نظم حاکم بر داده‌های پیش‌پردازش شده، با توجه به روش کاوش داده‌ای که انتخاب می‌شود، شناسایی شده و مدل تولید شده برای ارزیابی به مرحله بعد یعنی ارزیابی و تفسیر مدل منتقل خواهد شد. میتوان روش‌های مختلف کاوش داده را در دو گروه روش‌های پیش‌بینی (یادگیری باناظر) و روش‌های توصیفی (یادگیری بدون ناظر) طبقه بندی نمود. انواع الگوریتم‌ها و روش‌ها برای یادگیری مدل وجود دارند که در این مقاله ما بر روی الگوریتم شبکه عصبی و درخت تصمیم تمرکز خواهیم کرد و در ادامه توضیحات تکمیلی این دو الگوریتم ارائه می‌شود.

منظور از ارزیابی دانش آن است که می‌بایست میزان صحت دانش تولیدشده مشخص شود تا بتوان به آن اعتماد نمود و به صورت عملی از آن استفاده کرد. تفسیر مدل به معنای آن است که دانش تولید شده را مورد بررسی قرار داده و توجیهی معنایی برای تبیین منطق آن ارائه نماییم [۴].

۱/۴/۱. شبکه‌های عصبی

شبکه عصبی مصنوعی، از مجموعه نورون‌های عصبی تشکیل می‌شود که بصورت موازی با هم برای حل یک مسئله مشخص کار می‌کنند. شبکه‌های عصبی با مثال کار می‌کنند و نمی‌توان آن‌ها را برای انجام وظیفه‌ی خاص برنامه‌ریزی کرد. مثال‌ها باید با دقت انتخاب شوند در غیر این صورت زمان سودمند، تلف می‌شود و یا حتی بدتر از آن شبکه ممکن است نادرست کار کند. امتیاز شبکه عصبی این است که خودش کشف می‌کند که چگونه مسئله را حل کند، عملکرد آن غیرقابل پیشگویی است. مهم‌ترین عواملی که گونه‌ها و کاربردهای شبکه عصبی را از یکدیگر متمایز می‌سازند، نوع نورون به کار گرفته شده، چیدمان یا معماری شبکه، بازه ورودی/خروجی هاست. وضعیت نسبی سلول‌ها در شبکه (تعداد و گروه بندی و نوع اتصالات آن‌ها) را معماری شبکه گویند. در معماری یک شبکه تعداد لایه‌ها و اتصالات بین آن‌ها مهم است. ورودی‌های شبکه با نام "لایه ورودی" و خروجی شبکه با نام "لایه خروجی" و در صورت نیاز، لایه‌های میان این دو "لایه پنهان" نامیده می‌شود [۱۴].

• لایه ورودی

این لایه ورودی‌ها را دریافت می‌کند و بر حسب قدرت ارتباطش با لایه بعد سیگنال ورودی را به لایه بعد می‌فرستد. قدرت ارتباط هر نورون با نورون دیگر را وزن آن نورون می‌گویند.

• لایه میانی

تعداد لایه‌های میانی و تعداد نورون‌های آن دلخواه است. لایه‌های میانی باید با دقت انتخاب شوند تا خروجی مناسب را به ما بدهند.

• لایه خروجی

گروه دیگری از نورون‌ها نیز از طریق خروجی‌های خود، جهان خارج را می‌سازند.

شبکه عصبی مصنوعی مانند یک تابع عمل می‌کند که به تعداد نورون‌های ورودی، ورودی می‌گیرد و به تعداد نورون‌های خروجی، خروجی می‌دهد.

• یادگیری در شبکه عصبی

هر شبکه عصبی سه مرحله آموزش، اعتبار سنجی و اجرا را پشت سر می گذارد. آموزش دیدن شبکه های عصبی در واقع چیزی جز تنظیم وزن های ارتباطی این نورون ها به ازای دریافت مثال های مختلف نیست تا خروجی شبکه به سمت خروجی مطلوب همگرا شود.

یک شبکه عصبی با تعیین رابطه بین ورودی ها و خروجی ها یاد می گیرد. این رابطه با محاسبه اهمیت نسبی ورودی ها و خروجی های سیستم تعیین می شود. با توجه به آزمون و خطا، سیستم نتایجش را با نتایج داده شده توسط متخصص مقایسه می کند، تا زمانی که به سطح مشخصی از دقت برسد. با هر آزمایش وزن های نسبت داده شده به ورودی ها تغییر می کند تا نتایج مطلوب به دست آید.

• مراحل یادگیری

- ۱- مقداردهی اولیه: وزن های اولیه به تصادف در محدوده مورد نظر انتخاب می شوند.
 - ۲- فعال سازی: نورون ها با به کار بردن ورودی ها و خروجی های مطلوب فعال می شوند. خروجی واقعی با تکرار محاسبه می شود.
 - ۳- آموزش وزن: وزن ها با توجه به خروجی مطلوب تصحیح می شوند.
 - ۴- تکرار: برگشتن به مرحله دو و تکرار پروسه تا زمان همگرایی.
- برای محاسبه وزن هر نورون در شبکه های عصبی از فرمول زیر استفاده می شود:

$$net_A = \sum_{i=0}^n w_{iA} x_i \quad \text{رابطه (۱)}$$

w_{iA} : وزن سیگنال ورودی نورون A و خروجی نورون i

x_i : وزن نورون i

n : تعداد نورون های لایه قبل

برای محاسبه میزان خروجی هر نورون نیز از فرمول زیر استفاده می شود:

$$output_A = \frac{1}{1 + e^{-x_A}} \quad \text{رابطه (۲)}$$

تا آخرین نورون برای تمام نورون ها این محاسبات انجام می شود. میزان خروجی بدست آمده برای نورون آخر را با مقدار خروجی واقعی مقایسه کرده و میزان دقت و خطا طبق فرمول زیر محاسبه می شود:

$$SSE = \sum_{records} \sum_{output\ nodes} (Real\ Value - Estimated\ Value)^2 \quad \text{رابطه (۳)}$$

Real Value: مقدار خروجی واقعی

Estimated Value: مقدار خروجی محاسبه شده

حال وزن سیگنال های ورودی خروجی طبق فرمول های زیر از آخرین لایه تا اولین لایه تصحیح می شود:

$$W_{ij, new} = W_{ij, current} + \Delta W_{ij} \quad \text{رابطه (۴)}$$

$$\Delta W_{ij} = \eta \delta_j * output_i \quad \text{رابطه (۵)}$$

$$\delta_j = output_j(1 - output_j)(actual_j - output_j) \quad \text{رابطه (۶)}$$

η : مقداری ثابت که به صورت پیش فرض مشخص می شود.

actual_j: مقدار خروجی واقعی عصب j

برای لایه های میانی از فرمول زیر برای محاسبه δ_j استفاده می شود:

$$\delta_j = output_j(1 - output_j) \sum_{downstream} W_{jk, new} \delta_k \quad \text{رابطه (۷)}$$

۱/۴/۲. درخت‌های تصمیم

درخت تصمیم در حقیقت یک فلوچارت با ساختار درختی است. در این درخت برگ‌ها و گره‌هایی وجود دارند که این گره‌ها یا یک گره تصمیم هستند و یا یک گره برگ هستند. در هر گره تصمیم مقادیر یک متغیر بررسی می‌گردد و توسط شاخه‌ها انشعاب پیدا می‌کنند و در گره برگ نیز مقدار غالب متغیر رده رکوردهای موجود در آن انشعاب نوشته می‌شود. برای تولید درخت تصمیم در حقیقت باید از شاخص‌هایی برای پیدا کردن خلوص انشعاب دهی استفاده نمود و متغیری را که خلوص بیشتری در انشعابات ایجاد می‌کند را در درخت استفاده کرد. در اینجا منظور از خلوص این است که پس از انشعاب در هر شاخه فقط رکوردهایی که دارای یکی از رده‌های موجود در مجموعه داده هستند، حضور داشته باشند. یکی از شاخص‌های انتخاب نقطه‌ای انشعاب ضریب جینی می‌باشد که به صورت زیر محاسبه می‌گردد:

$$I(s) = \frac{|s_1|}{|s|} I(s_1) + \frac{|s_2|}{|s|} I(s_2) \quad \text{رابطه (۸)}$$

$$I(s_1) = 1 - \sum p_i^2 \quad \text{رابطه (۹)}$$

$|s|$: تعداد کل حالات موجود

$|s_1|$: تعداد کل حالات s_1

p_i : فراوانی نسبی برابر تعداد حالات هر کلاس تقسیم بر کل

انشعابات در این درخت یا دودویی هستند و یا چندگانه هستند [۱۴]. که درخت ایجاد شده در متلب برای این پژوهش یک درخت دودویی است.

۱/۵. مروری بر مطالعات انجام شده

در مطالعات انجام شده توسط آقای علی رفیعی و همکارانش با عنوان به کارگیری داده کاوی برای تشخیص تقلب در تراکنش کارت‌های اعتباری روش‌های مختلف داده کاوی مورد بررسی قرار گرفته است. در این مقاله الگوریتم‌های داده کاوی نظیر الگوریتم فازی، رگرسیون، درخت تصمیم، جنگل تصادفی، خوشه بندی، شبکه عصبی، برپایه قوانین، سیستم ایمنی مصنوعی، شبکه بیزین، روش‌های آماری بررسی شده است که برآیند مطالعه نشان می‌دهد استفاده از رهیافت‌های ترکیبی نظیر دو روش دسته بندی موجب بهبود نتایج از نظر معیارهایی نظیر دقت کشف، سرعت کشف، هزینه کشف خواهد شد. به کارگیری مزیت‌های الگوریتم‌ها در کنار یکدیگر می‌تواند علت اصلی این بهبودها باشد [۵].

در مقاله انجام شده توسط خانوم زهره زرین قلمی تقلب در کارت‌های اعتباری از جمله تقلب ورشکستگی، کارت گمشده یا دزدی، تقلب درخواست، کنترل حساب، جعل کارت، تبانی فروشنده و تقلب‌های اینترنتی تشریح و سپس انواع روش‌های داده کاوی برای شناسایی تقلب در کارت‌های اعتباری مانند شبکه‌های عصبی، برپایه قوانین، درخت تصمیم، شبکه‌های بیزین، خوشه‌ای، آماری، الگوریتم‌های ژنتیک، منطق فازی رگرسیون، جنگل تصادفی بررسی شده است. به این نتیجه رسیده است که هیچ کدام از روش‌ها به تنهایی قادر به شناسایی تمامی انواع تقلب‌ها نیستند و هرکدام برای شناسایی نوع خاصی از تقلب مناسب هستند، به همین دلیل در بیشتر مراجع از چندین روش و نیز ترکیب آن‌ها استفاده شده و تنها در شمار بسیار کمی از آن‌ها از یک روش به تنهایی استفاده شده است [۲].

در پژوهشی بهنام حیدری شبیه سازی در زمینه کشف تقلب کارت‌های اعتباری در سیستم‌های بانکداری الکترونیک انجام داده است. این پیاده سازی با استفاده از زبان برنامه نویسی C# و نرم افزارهای داده کاوی Rapidminer و Weka انجام شده است که دریافته استفاده از قوانین وابستگی در الگوریتم‌های دسته بندی مثل KNN می‌توان بهبود قابل توجهی در دقت دسته بندی ایجاد نماید. همچنین ترکیب تکنیک‌های خوشه بندی K-means با تعداد Kهای بهینه و قوانین انجمنی نقش بسیار مهمی در کشف تقلب داشته است. لذا دقت الگوریتم پیشنهادی به منظور کشف تقلب در داده‌های بانکداری الکترونیک، نسبت

به روش معمول بسیار مناسب بوده است. بنابراین در این مقاله به روش جدید دسته بندی تلفیقی با خوشه بندی و قوانین انجمنی مبتنی بر قوانین وابستگی به منظور افزایش دقت کشف تقلب در سیستم های بانکداری الکترونیک ارائه شده است [۶]. هوانگ با استفاده از انتخاب ویژگی در مدل نظارت شده، مدل های آماری خطی و غیر خطی و مدل های یادگیری ماشین از قبیل شبکه عصبی، رگرسیون لاجستیک، درخت تقویت شده و جنگل تصادفی که بر مبنای داده های دارای نویز و داده های تراکنش کارت اعتباری تحقیق شده است به این نتیجه رسید مدل جنگل تصادفی بهترین عملکرد را دارد و می تواند نیمی از تلاش های مربوط به تقلب را تنها در ۳٪ داده های برتر به عنوان مشکوک با امتیاز الگوریتم تقلب تشخیص دهد. با این حال، نتایج خارج از زمان (out-of-time) جنگل تصادفی به دلیل مجموعه داده های کوچک بسیار نزدیک به شبکه عصبی است. از مزایای این تحقیق این است که برای مدل های خطی و غیر خطی استفاده می شوند و چون نظارت شده هستند، بسیار دقیق و قابل اعتمادند. نیاز به داده های بیشتر و همچنین مشکل داده های نامتوازن از معایب کار است [۱۰]. در پژوهش آقای دومان و همکارش تکنیک های پیشرفته داده کاوی به همراه الگوریتم شبکه عصبی برای پوشش بالای تقلب های انجام شده، معرفی شده است. به منظور مقایسه میزان قابلیت اطمینان و درصد صحت تشخیص تقلب الگوریتم های پیشنهادی، دو تکنیک شبکه عصبی و درخت تصمیم انتخاب شده اند و از آنجایی که هر دو روش انتخابی از نوع روش های یادگیری با ناظر می باشند، سعی شده تا این مقایسه از طریق تغییر در درصد یادگیری شان و سپس محاسبه درصد صحت تشخیص بر روی داده های مصنوعی، انجام شود. در ادامه به منظور دست یافتن به نتایج بهتر و قابل اعتمادتر، دو روش را با یکدیگر ادغام نموده و دوباره در میزان یادگیری های مختلف مقایسه صورت می گیرد. نتایج به خوبی نشان می دهد که ترکیب دو مدل در صورتیکه خروجی شبکه عصبی به عنوان منبع برای درخت تصمیم در نظر گرفته شود، در تشخیص تقلب کارت اعتباری کارآمدتر عمل می نماید [۸].

۲. روش پیشنهادی

برای کشف تقلبات صورت گرفته در تراکنش های کارت های اعتباری در این پژوهش به پیاده سازی الگوریتم های شبکه عصبی و درخت تصمیم در نرم افزار متلب پرداخته می شود و این دو الگوریتم را از نظر دقت پیش بینی مقایسه کرده و الگوریتم با خطای کمتر را معرفی می کنیم. جهت آماده سازی داده برای استفاده از الگوریتم های یادگیری ماشین (درخت تصمیم و شبکه عصبی) ستون هایی از مجموعه داده که مقادیر آن ها عددی نیست را با استفاده از دستور `grp2idx` به مقادیر عددی تبدیل می کنیم چون الگوریتم یادگیری با مقادیر عددی کار می کند. سپس با استفاده از دستور `normalize` مقادیر ویژگی ها را نرمالیزه می کنیم تا مقادیر آن ها در یک بازه استاندارد قرار بگیرند.

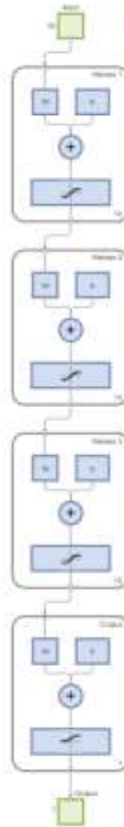
پس از شناخت و درک داده ها و مشخص نمودن فیلد هدف به پیاده سازی دو الگوریتم شبکه های عصبی و درخت تصمیم بر روی دو مجموعه داده معرفی شده می پردازیم. بدین صورت که در نرم افزار ذکر شده ابتدا مجموعه داده های موردنظر را بارگذاری کرده و برای شبکه عصبی در هر دو مجموعه داده ویژگی (Features) از همه ستون ها بجز ستون آخر استخراج می شود و بعنوان ورودی به الگوریتم داده می شود که دارای مجموعه ای از ویژگی های مشترک هستند و آخرین ستون به عنوان تارگت و هدف در نظر گرفته می شود که متغلب بودن یا نبودن را نشان می دهد. پس از این مرحله با استفاده از `CVPartition`، مجموعه داده را به دو بخش مجموعه داده های آموزشی^۵ و مجموعه تست^۶ تقسیم کرده که با استفاده از مجموعه داده آموزشی، الگوریتم آموزش می بیند و به یادگیری می پردازد و مدل نهایی ایجاد می شود و سپس با استفاده از مجموعه تست، مدل ایجاد شده مورد ارزیابی و اعتبارسنجی قرار می گیرد. یک شبکه عصبی با سه لایه پنهان که هر یک شامل ۱۰ نورون

⁵ Training Set

⁶ Test Set

است، با استفاده از patternnet تعریف می شود و الگوریتم شبکه عصبی روی مجموعه داده ها اعمال می شود. قابل ذکر است که این تعداد لایه و نرون به این دلیل انتخاب شده که الگوریتم دقت خوب و بالایی داشته باشد. اگر تعداد نرون و لایه کم باشد، نتیجه با دقت بالایی حاصل نمی شود و از طرفی اگر تعداد نرون و لایه ها زیاد باشد، حجم محاسبات و زمان ران بالا می رود. در شکل ۱ شبکه عصبی ساخته شده با ساختار ذکر شده نشان داده شده است.

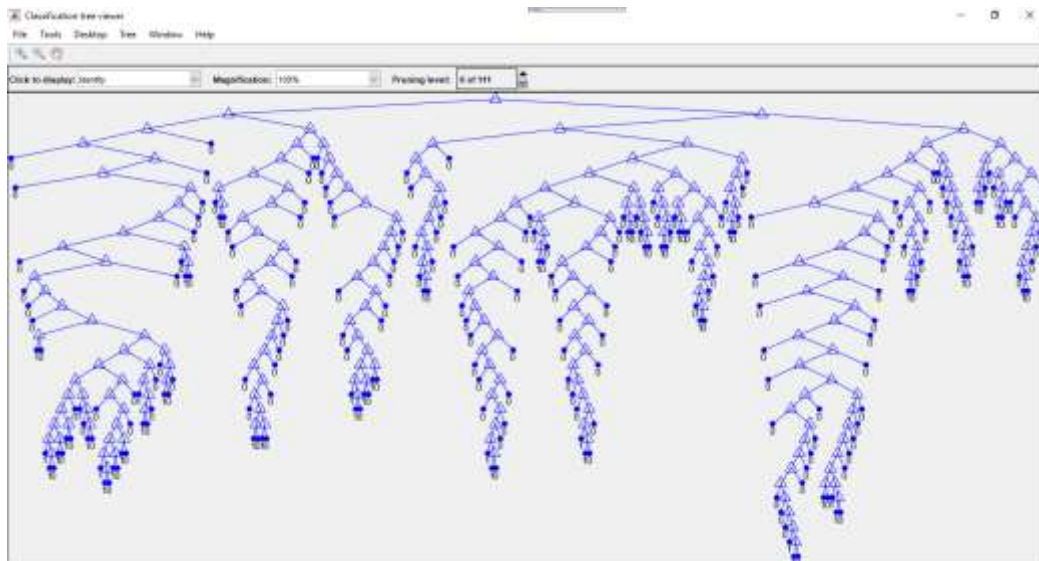
Pattern Recognition Neural Network (view)



شکل ۱ شبکه عصبی ساخته شده توسط الگوریتم

برای درخت تصمیم نیز مانند شبکه عصبی در هر دو مجموعه داده ستون هدف یا تارگت مشخص می شود و داده ها با استفاده از CVPartition مجموعه داده را به دو بخش داده های آموزشی و داده های آزمایش تقسیم بندی می کنیم. قبل از ایجاد درخت تصمیم گیری این نکته را باید در نظر داشته باشیم که ما ابتدا باید نمونه های مثبت و منفی را برای مدیریت عدم تعادل کلاس در مجموعه داده را شناسایی کنیم، نمونه های مثبت برچسب "۱" و نمونه های منفی برچسب "۰" دارند. سپس از کلاس اقلیت (نمونه های مثبت) نمونه برداری بیش از حد انجام دادیم (بیش نمونه برداری یا OverSample کردیم) تا تعداد نمونه های مثبت و منفی را متعادل کنیم و برای این کار از نمونه برداری با جایگزینی استفاده می کنیم. این کار برای جلوگیری از سوگیری درخت تصمیم به سمت کلاس اکثریت (نمونه های منفی) انجام می شود و اطمینان حاصل بشود که مدل به نفع کلاس اکثریت نیست در غیر این صورت می تواند منجر به عملکرد ضعیف در پیش بینی کلاس اقلیت شود. برای شبکه عصبی نمونه برداری بیش از حد یا oversample را اجرا نکردیم چون شبکه عصبی به دلیل پیچیدگی و توانایی یادگیری الگوهای پیچیده تر از داده ها و مکانیسم های متفاوتی مانند backpropagation و معماری های پیچیده تر، ذاتا می تواند عدم تعادل کلاس را بهتر و موثرتر مدیریت کند. پارامترهای استفاده شده در ساخت درخت تصمیم دودویی در متلب به شرح زیر می باشد:

- **MaxNumSplts**: حداکثر تعداد انشعابات در درخت می باشد که در اینجا ۵۰۰۰ در نظر گرفته شده است تا درخت تصمیم بتواند عمیقاً رشد کند و الگوهای پیچیده تری را در داده ها ثبت کند. این زمانی مفید است که داده ها پیچیده هستند و برای دسته بندی دقیق نیاز به مدلی دقیق دارند.
 - **MinLeafSize**: حداقل اندازه هر برگ می باشد که یک در نظر گرفته شده است که درخت می تواند به اندازه لازم برگ های کوچک ایجاد کند و در صورت نیاز الگوهای بسیار خاصی را ثبت کند. این مورد می تواند به بهبود دقت کمک کند، اما خطر بیش از حد پردازش را نیز افزایش می دهد.
 - **SplitCriterion**: معیار تقسیم بندی می باشد که **gdi** در نظر گرفته شده است. **(GDI)** شاخص تنوع جینی معیاری از ناخالصی است که برای تقسیم گره ها استفاده می شود. هدف آن ایجاد گره هایی با کلاس های خالص یا همگن است.
 - **PredictorSelection**: روش انتخاب پیش بینی کننده می باشد که **curvature** در نظر گرفته شده است. این گزینه از آنجا برای انتخاب بهترین پیش بینی کننده ها برای تقسیم استفاده می کند و به انتخاب ویژگی هایی که رابطه غیر خطی با متغیر هدف دارند کمک می کند.
 - **Surrogate**: استفاده از متغیرهای جایگزین می باشد که این گزینه را در حالت **on** قرار دادیم. گزینه **on** یا روشن برای جانشین ها به درخت اجازه می دهد از تقسیم های جایگزین استفاده کند که می تواند با ارائه معیارهای تقسیم جایگزین، داده های از دست رفته را مدیریت کند.
- در شکل ۲ بعنوان نمونه درخت تصمیم ساخته شده توسط الگوریتم برای مجموعه داده **Creditcard** با درصد داده آموزشی ۸۰ درصد و داده آزمایشی ۲۰ درصد نشان داده شده است.



شکل ۲ درخت تصمیم ساخته شده توسط الگوریتم برای دیتاست **Creditcard** با **train=80%**

در این پژوهش در دو مرحله الگوریتم های مورد نظر مورد بررسی قرار گرفتند، به طوریکه در مرحله اول ۸۰ درصد از داده ها به عنوان داده های آموزشی و ۲۰ درصد دیگر به عنوان مجموعه تست تعیین شد و مدل نهایی ایجاد و مورد بررسی قرار خواهد گرفت. در مرحله بعد ۷۰ درصد از داده ها به عنوان داده های آموزشی و ۳۰ درصد دیگر به عنوان مجموعه تست تعیین شد و مدل مورد بررسی قرار می گیرد.

۳. ارزیابی نتایج

برای پیاده سازی الگوریتم های درخت تصمیم و شبکه عصبی بر روی دو مجموعه داده که در ادامه به معرفی آنها می پردازیم، از نرم افزار قدرتمند متلب نسخه ۲۰۲۳ استفاده شده که اجرا شده بر روی سیستم با ویندوز 10 Pro و پردازنده Core i5 و رم 4GB می باشد.

همچنین برای اجرای این دو الگوریتم از جعبه ابزار (ToolBax) آماده متلب و توابع موجود در آن استفاده شده است. برخی از توابع استفاده شده شامل تابع readtable برای خواندن مجموعه داده، تابع cvpartition برای تقسیم بندی داده ها به دو بخش آموزش و تست، تابع patternnet برای ایجاد شبکه عصبی، تابع fitctree برای ایجاد درخت تصمیم، تابع datasample برای oversample یا نمونه برداری بیش از حد، تابع train برای آموزش الگوریتم و تابع confusion برای ایجاد ماتریس اشفستگی می باشد.

از دو مجموعه داده که از سایت کاکل تهیه شده در این پژوهش استفاده شده است. مجموعه داده اول به نام Creditcard که شامل معاملات انجام شده توسط کارت های اعتباری در سپتامبر ۲۰۱۳ می باشد که معاملات رخ داده طی دو روز می باشد و از بین ۲۸۴۸۰۷ معامله ۴۹۲ معامله متقلبه بوده. این مجموعه شامل تنها متغیرهای عددی می باشد و با توجه به مسائل مربوط به محرمانگی اطلاعات ویژگی های اصلی و اطلاعات پس زمینه بیشتر ارائه نشده است. این مجموعه داده دارای ۲۹ ستون ویژگی و یک ستون هدف می باشد که متقلبه بودن یا نبودن تراکنش را مشخص میکند. مجموعه داده دوم به نام Banksim می باشد که این مجموعه داده به صورت مصنوعی تولید شده که شامل پرداخت هایی از مشتریان مختلف است که در دوره های زمانی مختلف و با مبالغ مختلف انجام می شوند. مجموعه داده دارای ۹ ستون ویژگی و یک ستون هدف است.

۳/۱. معرفی معیارهای ارزیابی

برای اعتبارسنجی و مقایسه کارایی الگوریتم ها از معیارهای مختلفی استفاده می شود. از جمله این معیارها میتوان به نرخ فراخوانی^۷، دقت^۸، صحت^۹ و F-measure اشاره نمود که پس از تشکیل ماتریس پراکندگی یا درهم ریختگی متقلب و غیرمتقلب که در جدول ۱ نشان داده شده است به راحتی قابل محاسبه هستند. پارامترهای عنوان شده در این ماتریس به شرح ذیل می باشد [۱۵].

پارامتر (TN)^{۱۰}: بیانگر تعداد رکوردهایی که دسته واقعی آنها متقلب بوده و الگوریتم دسته بندی نیز دسته آنها را به درستی متقلب پیش بینی کرده.

پارامتر (FN)^{۱۱}: بیانگر تعداد رکوردهایی که دسته واقعی آنها غیرمتقلب بوده و الگوریتم دسته بندی آنها را به اشتباه متقلب پیش بینی کرده.

پارامتر (TP)^{۱۲}: بیانگر تعداد رکوردهایی که دسته واقعی آنها غیرمتقلب بوده و الگوریتم دسته بندی آنها را به درستی غیر متقلب پیش بینی کرده.

پارامتر (FP)^{۱۳}: بیانگر تعداد رکوردهایی که دسته واقعی آنها متقلب بوده و الگوریتم دسته بندی آنها را به اشتباه غیرمتقلب پیش بینی کرده.

⁷ Recall

⁸ Precision

⁹ Accuracy

¹⁰ True Negative

¹¹ False Negative

¹² True Positive

¹³ False Positive

جدول ۱ - ماتریس رفتارهای متقلبانه و غیر متقلبانه

| نوع رکورد | رکوردهای پیش بینی شده (Predicated Records) | | |
|-----------|--|-----------|-------|
| | نوع دسته | غیر متقلب | متقلب |
| | غیر متقلب | TP | FN |
| | متقلب | FP | TN |

مهم ترین معیار برای تعیین کارایی یک الگوریتم دسته بندی، معیار صحت می باشد، این معیار دقت کل یک دسته بند را محاسبه می کند. این معیار نشان دهنده این موضوع است که چند درصد از کل مجموعه داده ها به درستی دسته بندی شده است رابطه زیر نحوه محاسبه معیار را نشان می دهد [۱۵].

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN} \quad \text{رابطه (۱۰)}$$

دو مقدار TP و TN مهم ترین مقادیری هستند که باید بیشینه شوند تا کارایی دسته بندی به حداکثر برسد. معیارهای دقت درصدی را نشان می دهند که از میان تمامی دسته ها که توسط دسته بند به آن دسته نسبت داده شده اند، درست دسته بندی شده اند. نحوه محاسبه این معیار در رابطه زیر نشان داده شده است.

$$Precision = \frac{TP}{TP + FP} \quad \text{رابطه (۱۱)}$$

معیار فراخوانی برای یک دسته، که از میان تمامی دسته های متقلب متعلق به آن دسته، به درستی دسته بندی شده است. نحوه محاسبه این معیار به صورت زیر می باشد.

$$Recall = \frac{TP}{TP + FN} \quad \text{رابطه (۱۲)}$$

معیار F-measure از ترکیب معیارهای Precision و Recall بدست می آید و در مواردی استفاده می شود که نتوان اهمیت ویژه ای را برای هر یک از دو معیار Precision و Recall نسبت به یکدیگر قائل شد. رابطه زیر نحوه محاسبه این معیار را نشان می دهد [۱۷].

$$F - measure = \frac{2 * Precision * Recall}{Precision + Recall} \quad \text{رابطه (۱۳)}$$

برای ارزیابی روش های دسته بندی از معیارهای بالا استفاده خواهد شد.

۳/۲. ارزیابی

هدف از دسته بندی، ایجاد مدلی بر مبنای داده های آموزشی عنوان شده که بتواند صفت کلاس را برای داده های مجموعه تست پیشگویی کند. همان طور که گفته شد مجموعه تست برای ارزیابی و اعتبارسنجی مدل ساخته شده مورد استفاده قرار می گیرد به همین دلیل مقادیر گزارش شده برای معیارهای ارزیابی در جداول و نمودارها، بدست آمده از نتایج مجموعه تست مدل ساخته شده می باشد. برای مقایسه کارایی الگوریتم ها این نتایج در جداول و روی نمودارها قرار داده شده است و در ادامه به تشریح آن ها پرداخته می شود.

در جداول ۲، ۳، ۴ و ۵ نتایج بدست آمده از مرحله اول و دوم پیاده سازی قرار داده شده است. که در مرحله اول ۸۰ درصد از داده‌ها به عنوان مجموعه آموزشی و ۲۰ درصد دیگر به عنوان مجموعه تست در نظر گرفته شده بود و در مرحله دوم ۷۰ درصد از داده‌ها به عنوان مجموعه آموزشی و ۳۰ درصد دیگر به عنوان مجموعه تست در نظر گرفته شد. در سطر دوم و سوم جداول به ترتیب، مقادیر معیارهای ارزیابی برای الگوریتم شبکه عصبی و الگوریتم درخت تصمیم قرار داده شده است. در ستون دوم جدول مقادیر معیار Precision که بیان کننده دقت الگوریتم در پیش بینی می‌باشد برای هر دو الگوریتم قرار گرفته است؛ در ستون سوم مقادیر معیار Recall که بیان کننده درصد فراخوانی هر الگوریتم می‌باشد قرار گرفته است؛ در ستون سوم مقادیر معیار Accuracy که بیان کننده درصد صحت پیش بینی انجام شده توسط الگوریتم‌ها می‌باشد قرار گرفته است و در ستون چهارم نیز مقادیر معیار F-measure برای هر دو الگوریتم قرار گرفته است.

جدول ۲ نتایج بدست آمده برای دیتا ست Creditcard با train=80% و Test=20%

| معیار (درصد) | Precision | Recall | Accuracy | F-measure |
|---------------|-----------|---------|----------|-----------|
| الگوریتم | | | | |
| Neural Net | ۹۹.۹۸۵۹ | ۹۹.۹۰۸۶ | ۹۹.۸۹۴۷ | ۹۹.۹۴۷۲ |
| Decision Tree | ۹۹.۹۶۴۷ | ۱۰۰ | ۹۹.۹۸۲۴ | ۹۹.۹۸۲۴ |

جدول Error! No text of specified style in document. نتایج بدست آمده برای دیتا ست Creditcard با train=70% و Test=30%

| معیار (درصد) | Precision | Recall | Accuracy | F-measure |
|---------------|-----------|---------|----------|-----------|
| الگوریتم | | | | |
| Neural Net | ۹۹.۹۶۷۲ | ۹۹.۹۶۰۱ | ۹۹.۹۲۷۴ | ۹۹.۹۶۳۷ |
| Decision Tree | ۹۹.۹۵۳۲ | ۱۰۰ | ۹۹.۹۷۶۶ | ۹۹.۹۷۶۶ |

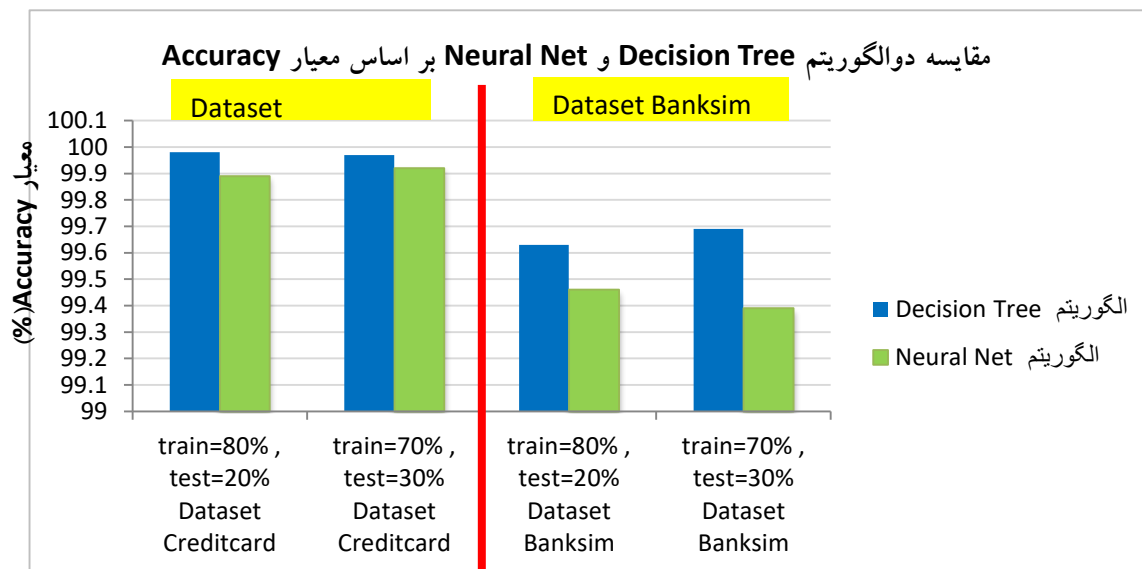
جدول Error! No text of specified style in document. نتایج بدست آمده برای دیتا ست Banksim با train=80% و Test=20%

| معیار (درصد) | Precision | Recall | Accuracy | F-measure |
|---------------|-----------|---------|----------|-----------|
| الگوریتم | | | | |
| Neural Net | ۹۹.۹۲۰۹ | ۹۹.۵۳۶۳ | ۹۹.۴۶۱۹ | ۹۹.۷۲۸۲ |
| Decision Tree | ۹۹.۲۷۸۷ | ۱۰۰ | ۹۹.۶۳۴۵ | ۹۹.۶۳۸۱ |

جدول. Error! No text of specified style in document. نتایج بدست آمده برای دیتا ست Banksim با train=70% و Test=30%

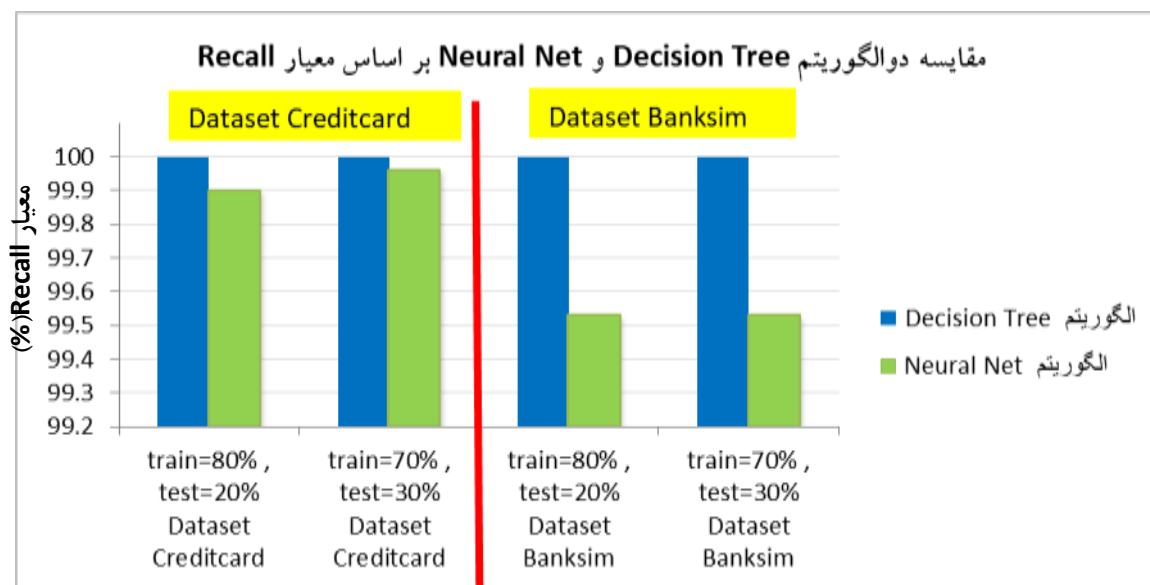
| معیار (درصد) | Precision | Recall | Accuracy | F-measure |
|---------------|-----------|---------|----------|-----------|
| الگوریتم | | | | |
| Neural Net | ۹۹.۸۶۲۷ | ۹۹.۵۳۰۷ | ۹۹.۳۹۹۱ | ۹۹.۶۹۶۴ |
| Decision Tree | ۹۹.۴۰۲۹ | ۱۰۰ | ۹۹.۶۹۷۱ | ۹۹.۷۰۰۶ |

همانطور که مشاهده می کنید طبق نتایج به دست آمده در جداول ۲، ۳، ۴ و ۵ مقادیر Precision برای الگوریتم شبکه عصبی کمی بیشتر از درخت تصمیم بوده است ولی در هر دو الگوریتم این مقادیر درصد مناسبی برای دقت الگوریتم ها می باشد ولی برای مقایسه نمی توان از آن استفاده کرد. مقدار F-measure نیز در هر مرحله مقدار قابل توجهی داشته و در حالت کلی درخت تصمیم در هر مرحله درصد بیشتری اخذ کرده است. معیار Recall در هر مرحله برای الگوریتم درخت تصمیم مقدار برابر ۱۰۰ درصد را اخذ کرده و نسبت به الگوریتم شبکه عصبی درصد بالاتری را به خود اختصاص داده است. و اما معیار Accuracy همانطور که اشاره شد مهم ترین معیار برای الگوریتم های دسته بندی می باشد که این معیار در همه مراحل، مقادیر الگوریتم درخت تصمیم نسبت به الگوریتم شبکه عصبی درصد بیشتری را اخذ نموده و عملکرد بهتری داشته است. مقایسه عملکرد دو الگوریتم Neural Net و Decision Tree بر اساس سه معیار Accuracy، Recall و F-measure روی نمودار در شکل های ۳، ۴ و ۵ نشان داده شده است.

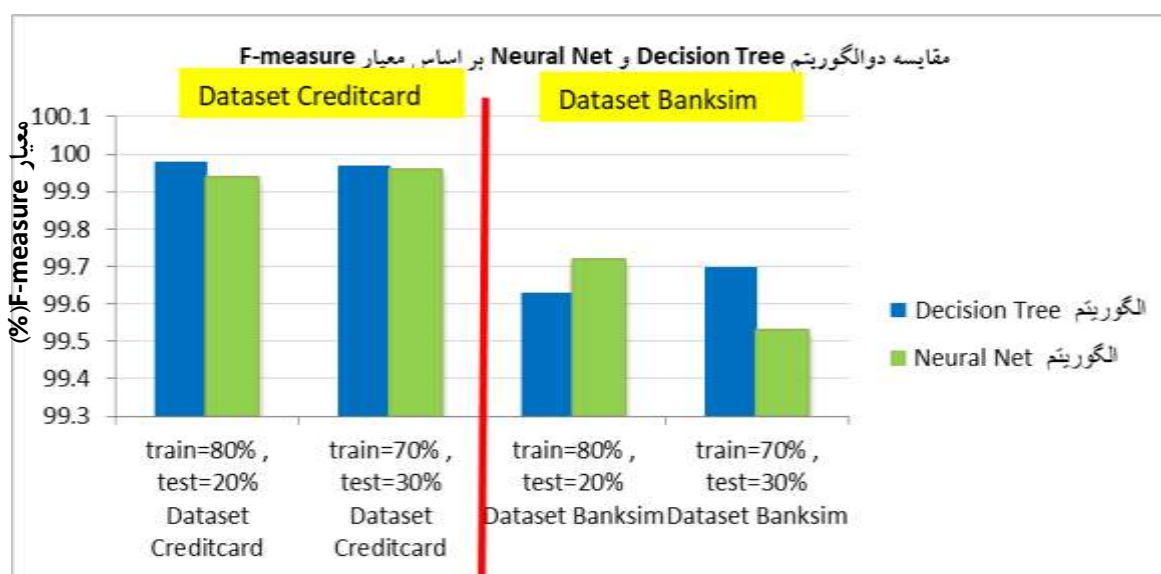


شکل. Error! No text of specified style in document. مقایسه دو الگوریتم Decision Tree و Neural Net بر اساس معیار

Accuracy



شکل Error! No text of specified style in document. مقایسه دو الگوریتم Decision Tree و Neural Net بر اساس معیار Recall



شکل Error! No text of specified style in document. مقایسه دو الگوریتم Decision Tree و Neural Net بر اساس معیار F-

measure

به طور کلی همانطور که در نمودارها نیز مشخص است، در هر مرحله از پیاده سازی با توجه به معیارهای Accuracy، Recall و F-measure، الگوریتم درخت تصمیم عملکرد بهتری نسبت به الگوریتم شبکه عصبی داشته است.

۴. بحث و نتیجه گیری

برای تحلیل تراکنش‌های انجام شده در کارت‌های اعتباری توسط مشتری که از عامل‌های مهم و تاثیرگذار در کشف تقلب است، از الگوریتم‌های دسته بندی درخت تصمیم و شبکه عصبی که به بیان دانش و استخراج در ارتباط با مجموعه‌ای از اطلاعات می‌پردازد، استفاده شده است. طبق گزارشات صورت گرفته و نتایجی که از دو الگوریتم انتخابی طی مراحل پیاده‌سازی بدست آمد؛ الگوریتم درخت تصمیم در هر مرحله عملکرد بهتری نسبت به الگوریتم شبکه‌های عصبی داشت اگر چه طبق معیارهای ارزیابی بدست آمده هر دو الگوریتم درصد قابل توجهی از معیارها را جهت پیش‌بینی و کشف تقلبات به خود اختصاص دادند. در آخر برای جمع‌بندی و بهبود دستاورد پژوهشی مقایسه‌ای با دیگر پژوهش‌های ارائه شده در این زمینه انجام دادیم که طبق اکثر پژوهش‌های انجام شده، پژوهشگران به این نتیجه مهم دست یافتند که هیچ کدام از روش‌ها به تنهایی قادر به شناسایی تمامی انواع تقلب‌ها نیستند و هر کدام برای شناسایی نوع خاصی از تقلب مناسب هستند بنابراین تکنیک‌های داده‌کاوی به تنهایی کارایی بالایی ندارند و بهتر است از رویکرد ترکیبی آن‌ها استفاده شود به همین دلیل در بیشتر مراجع از چندین روش و نیز ترکیب آن‌ها استفاده شده است و تنها در شمار بسیار کمی از آن‌ها از یک روش به تنهایی استفاده شده است.

۵. منابع

- ۱-سیادت، سیدحسین؛ رباب اسمعیلی، هما؛ عزیزی، طوفان؛ "بررسی و مقایسه روش های داده کاوی جهت کشف تقلب در تراکنش های کارت های بانکی"، دومین کنفرانس بین المللی مدیریت در قرن ۲۱، ۱۳۹۴.
- ۲-زرین قلمی، زهره؛ "به کارگیری داده کاوی برای تشخیص تقلب در تراکنش های کارت های اعتباری"، دومین همایش ملی فناوری های نوین در مهندسی برق و کامپیوتر، دانشگاه آزاد اسلامی فسا، ۱۳۹۳.
- ۳-فصیح فر، زهره؛ رخصتی، حمیدرضا؛ وفاخواه، محبوبه؛ "بررسی کاربرد الگوریتم های هوشمند داده کاوی، در شناسایی و جلوگیری از انواع تقلبات بانکی"، سومین کنفرانس بین المللی مهندسی کامپیوتر و سیستم های خبره، ۱۳۹۵.
- ۴-صنّعی آباده محمد، محمودی سینا، طاهرپور محدّثه، "داده کاوی کاربردی"، تهران، نیاز دانش، ۱۳۹۱.
- ۵-رفیعی کشتلی، علی؛ رفیعی کشتلی، میثم؛ دکتر شمسی، محبوبه "به کارگیری داده کاوی برای تشخیص تقلب در تراکنش کارت های اعتباری"؛ اواین همایش ملی؛ ۱۳۹۲.
- ۶-حیدری، بهنام؛ "ارائه یک روش بهینه جهت بهبود کشف تقلب های کارت اعتباری در سیستم بانکداری الکترونیک با استفاده از ترکیب رویکردهای داده کاوی"؛ چهارمین کنفرانس بین المللی مهندسی برق و کامپیوتر؛ ۱۳۹۵.
- ۷-نجف زاده اصل، احسان و کریم پور، جابر و مقدم دیزج هریک، مجید؛ "ارائه‌ی یک روش جدید کشف نفوذ به منظور تشخیص رفتارهای غیرنرمال در شبکه های کامپیوتری با استفاده از تکنیک‌های داده کاوی"، چهاردهمین کنفرانس مهندسی برق ایران، کرمانشاه، ۱۳۹۰.
- 8-Zakaryazad, Ashkan, & Ekrem Duman, "A profit-driven Artificial Neural Network(ANN) with applications to fraud detection and direct marketing" Neurocomputing 175, 2016.
- 9-Zhang, Zhaohui, et al. "A model based on convolutional neural network for online transaction fraud detection", Security and Communication Networks, 2018.
- 10-Huang, "Credit Card Transaction Fraud Using Machine Learning Algorithms, International Conference on Education science and Economic Development, 2020.
- 11-https://en.wikipedia.org/wiki/credit_card_fraud.
- 12-El Madhoun, Nour, and Emmanuel Bertin. "Magic always comes with a price: Utility versus security for bank cards." 2017 1st Cyber Security in Networking Conference (CSNet). IEEE, 2017.
- 13-Miner, G., Nisbet, R., IV, J.E.; "Handbook of Statistical Analysis and Data Mining Applications, 1 edition ed", Academic press, Amsterdam, Boston, 2009.
- 14-Witten, I. H., & Frank, E.; "Data Mining: Practical machine learning tools and techniques", Morgan Kaufmann, 2005.
- 15-Raval, U. Marathe, N. Padiya, P.; "Feature Selection Based Hybrid Anomaly Intrusion Detection System Using K Means and RBF Kernel Function", scientific committee of International Conference on Advance

Computing Technologies and Application (ICACTA-2015), Published by Elsevier B.V. Procedia Coputer Science 45, PP: 428-435, 2015.

Fraud detection in credit cards using neural network and decision tree

Reza Roshani

1-Lamei Gorgani Institute of Higher Education, Gorgan, Iran

2-Department of Mechanical Engineering, National University of Skills (NUS), Tehran, Iran

Masoumeh Kiaee

1-Lamei Gorgani Institute of Higher Education, Gorgan, Iran

Abstract:

Using credit cards to prevent accidents caused by carrying cash has become popular. But this issue not only has not caused the retreat of profiteers, but has opened new ways with less risk for them. From the analysis of bank statistics and other credit card custodians, it appears that fraud in credit cards is increasingly increasing, and due to the high volume of transactions, the need for a suitable tool to detect data fraud is felt. Data mining is a new science that has been proposed in this field, so that by using its techniques and algorithms, it is possible to process this massive amount of data and search and identify suspicious cases. In this research, in order to provide a method in the field of fraud detection in credit card transactions, two widely used data mining techniques, namely, decision tree and neural network, are investigated on two data sets, and these algorithms are analyzed in two stages using MATLAB data mining tool. They will be implemented on selected data sets and their results will be analyzed and compared. The implementation results show that the decision tree algorithm has a higher percentage of the criteria than the neural network.

Keywords: Credit card, fraud detection, data mining, decision tree, neural network.