

طراحی سامانه تشخیص تقلب با استفاده از ماشین بردار

پشتیبان، انتخاب ویژگی و اعتبارسنجی متقابل

بابک رحمانی^۱، حسین قرایی گرانی^۲

۱- پردیس بین‌المللی کیش، دانشگاه تهران

Babak_rahmani@ut.ac.ir

۲- استادیار پژوهشکده امنیت، مرکز پژوهش‌های مخابرات

gharaee@itrc.ac.ir

چکیده

در سال‌های اخیر پرداخت الکترونیکی، رشد سریعی در میان فعالیت‌های اینترنتی داشته است؛ به طوری که امروزه به دلیل سرعت، کارایی، کاهش هزینه‌ها و سهولت دسترسی، مشتریان زیادی را به خود جذب کرده است. کارت‌های اعتباری یکی از پرکاربردترین ابزارهای پرداخت و مبادلات الکترونیکی هستند. در این پژوهش شناسایی و استخراج ویژگی‌های تراکنش‌های تقلبی در تشخیص تقلب و به دنبال آن طبقه‌بندی صحیح آن‌ها به دو طبقه قانونی و تقلبی با استفاده از ماشین بردار پشتیبان و اعتبارسنجی متقابل انجام شده است. نتایج ارزیابی‌های این روش نسبت به روش‌های انجام‌شده، حاکی از بمبود تشخیص در تقلب است؛ به طوری که خطاهای منفی کاذب، کاهش قابل ملاحظه‌ای به میزان ۷۷٪ داشته و به دنبال آن هزینه‌ها ۸۸٪ کاسته می‌شوند و نرخ تشخیص تقلب ۱۱٪ افزایش پیدا می‌کند.

وازگان کلیدی: تشخیص تقلب، کارت‌های اعتباری، ماشین بردار پشتیبان، اعتبارسنجی متقابل و استخراج ویژگی.

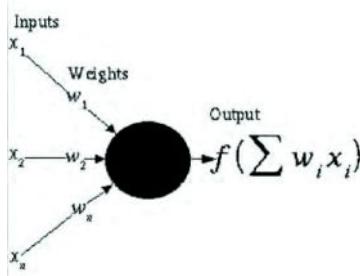
۱- مقدمه

روش‌های ارائه شده برای تشخیص تقلب در این پژوهش در حوزه‌هایی قابل استفاده است که تراکنش‌های مشتریان، قابل ردیابی باشند که بیشتر شامل کارت‌های اعتباری و کارت‌های بدهی می‌شود. خریدهایی را که با کارت‌ها انجام می‌شوند به دو دسته کلی می‌توان تقسیم کرد: دسته نخست خریدهایی است که کارت به صورت فیزیکی مورد استفاده قرار می‌گیرد. تقلب در این نوع خریدها از طریق دزدیدن کارت و یا کپی کارت^۱ صورت می‌پذیرد؛ اما در دسته دوم که کارت به صورت مجازی مورد استفاده قرار می‌گیرد، تنها اطلاعات کارت از جمله شماره کارت، تاریخ انقضا و رمز آن برای انجام خرید مورد نیاز است. این خریدها به طور معمول از طریق اینترنت، رایانامه و یا تلفن انجام می‌پذیرد؛ لذا متقلب با دردست داشتن اطلاعات کارت و بدون اطلاع دارنده کارت به راحتی خرید می‌کند و به طور معمول دارنده کارت خیلی دیر مطلع می‌شود که کارت او مورد سوء استفاده قرار گرفته است. در حال حاضر بیشترین تقلب در حوزه کارت‌ها در این دسته قرار می‌گیرند که امکان کنترل فیزیکی کارت که سطح امنیتی بالاتری را به وجود می‌آورد، وجود ندارد. امروزه ابزارها

کارت‌های بانکی به طور مؤثری بخش اساسی از سامانه پرداخت مدرن به حساب می‌آیند و دامنه وسیعی از خدمات پرداخت الکترونیکی را در اختیار کاربران سامانه قرار می‌دهند. کارت‌های الکترونیکی با وجود اینکه یکی از پیشرفته‌ترین سامانه‌های پرداخت هستند، ولی با این حال همان مشکلات پول نقد، مانند جعلی بودن و یا به سرقت رفتن را دارند. با توجه به گسترش روزافزون فناوری‌های نوین و ارتباطات جهانی، تقلب‌های الکترونیکی نیز در حال افزایش هستند. با افزایش معاملات جعلی، ضرر و زیان قابل توجهی به کسب و کار وارد شده و به همین دلیل، تشخیص تقلب را به یک موضوع مهم تبدیل کرده است. تشخیص تقلب را می‌توان به عنوان یک مشکل در طبقه‌بندی معاملات مشروع از معاملات جعلی دید. تکنیک‌های موجود در تشخیص تقلب به وسیله تعدادی از روش‌ها مانند داده‌کاوی، آمار و هوش مصنوعی اجرا شده‌اند. به طور کلی تشخیص تقلب یک مشکل پیش‌بینی است و هدف آن به حدکث رساندن پیش‌بینی‌های درست و کاهش پیش‌بینی‌های نادرست هزینه‌ها در یک سطح قابل قبول است [۱].

۱-۲ شبکه‌های عصبی^۳

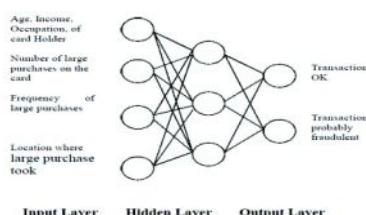
شبکه‌های عصبی یکی از تکنیک‌های هوش مصنوعی است که نشان‌دهنده مدلی از سامانه‌های یادگیری بیولوژیکی است. آنها شبکه‌هایی از پردازنده‌های ساده یا واحد هستند که به هم متصل شده و پردازش مقادیر عددی را انجام می‌دهند. آنها به عنوان یک گراف جهت‌دار با گره‌های بسیار (پردازش عناصر) و کمان (ارتباطات) بین آنها پیکربندی شده‌اند. گره وظیفه محاسبه مجموع وزن از ورودی‌های آن و تولید یک خروجی را بر عهده دارد. این خروجی پس از آن به عنوان یک ورودی برای گره‌های دیگر شبکه در نظر گرفته می‌شود. این فرایند همچنان ادامه می‌یابد تا زمانی که یک یا چند خروجی تولید شده باشد.



شکل ۱- یک گره شبکه عصبی

شبکه‌های عصبی به عنوان یک تکنیک قدرتمند داده‌کاوی در صنعت برای انجام کارهای طبقه‌بندی، دسته‌بندی، تعمیم و پیش‌بینی استفاده می‌شوند. برای مثال می‌توان آن را برای تشخیص قانونی تراکنش‌های جعلی، شناسایی تقلب اینترنتی در یک سایت تجارت الکترونیکی، پیش‌بینی تراکنش‌هایی که ممکن است جعلی باشند، مورد استفاده قرار داد. علاوه‌بر این، شبکه‌های عصبی یکی از نخستین و موفق‌ترین برنامه‌های کاربردی در زمینه تشخیص تقلب در کارت‌های اعتباری هستند^[۸].

شکل ۲ مدل کردن رفتار عادی مشتری را توسط شبکه عصبی نشان می‌دهد.



شکل ۲- لایه‌های شبکه عصبی در کارت‌های اعتباری

و تکنیک‌های پیشرفته‌ای بهمنظور پیش‌گیری از تقلب توسط بانک‌های صادرکننده کارت و فروشنده‌گان، مورد استفاده قرار می‌گیرند. از جمله رایج‌ترین این روش‌ها می‌توان به خدمات بازبینی نشانی^۱ و یا مقدار امنیتی کارت^۲ اشاره کرد. در عین حال که روش‌های پیش‌گیری از تقلب روزبه‌روز پیشرفته‌تر می‌شوند، متقلبان نیز پی‌درپی روش‌های خود را تغییر می‌دهند و روش‌های پیشرفته‌تری برای عدم شناسایی توسط ابزارهای مذکور ابداع می‌کنند؛ به‌گونه‌ای که بروز تقلب اجتناب‌ناپذیر است. زمانی که ابزارهای پیش‌گیری از تقلب نتوانند مانع از بروز تقلب شوند، راه حل بعدی، تشخیص هرچه سریع‌تر تقلب برای جلوگیری از سوء استفاده‌های بعدی است^[۱۰].

در ادامه، ساختار مقاله بدین شرح ارائه می‌شود: در فصل بعدی پژوهش‌هایی را که تاکنون در این حوزه صورت گرفته موردن بررسی قرار می‌دهیم. فصل سوم، روش‌های پیشنهادی این پژوهش برای بهبود تشخیص تقلب شرح داده خواهد شد. فصل چهارم به شرح ارزیابی‌های انجام‌شده برای روش‌های پیشنهادی می‌پردازد و درنهایت فصل پنجم جمع‌بندی، نتیجه‌گیری و چشم‌اندازی برای کارهای آینده ارائه خواهد شد.

۲- کارهای انجام‌گرفته

محققان دو دسته کلی از روش‌های تشخیص را توسعه داده‌اند: نخست سوء استفاده و دوم شناسایی ناهمجاري. در تشخیص سوء استفاده، معاملات جعلی به خوبی شناخته شده هستند و آنها با الگوهایی کدگذاری شده‌اند که سپس مورد استفاده برای مطابقت با معاملات جعلی جدید و شناسایی آنها قرار می‌گیرند. در تشخیص ناهمجاري، رفتار طبیعی از فعالیت‌های کاربر و سامانه برای نخستین بار به عنوان پروفایل‌های عادی خلاصه و پس از آن به عنوان معیارهایی استفاده می‌شوند؛ بنابراین زمان اجرای فعالیت‌ها اگر نتیجه، انحراف قابل توجهی از پروفایل کاربر داشته باشد، آن را به عنوان احتمال وقوع معاملات تقلب در نظر می‌گیریم. در این بخش ما به طور خلاصه به برخی از تکنیک‌های تشخیص تقلب در کارت‌های اعتباری و همچنین توصیف و ذکر مزایا و معایب آنها می‌پردازیم:

^۳Neural Networks

^۱AVS(Address Verification Service)

^۲CVV(Card Verification Value)

نتایج پژوهش بالا جهت مقایسه در بخش ارزیابی نتایج ذکر شده است [۵].

در پژوهش دیگری نیز طبقه‌بندی توسط ماشین بردار پشتیبان بر روی مجموعه داده‌های یک بانک بزرگی انجام شده است و نتایج آن را با نتایج شبکه‌های عصبی ترکیب کرده است، که در بخش ارزیابی نیز جهت مقایسه با نتایج روش ارائه شده ذکر شده است [۱۲].

۳- روش پیشنهادی

۱- تشخیص تقلب با استفاده از ماشین بردار

پشتیبان با حاشیه نامتوازن

تشخیص الگو با استفاده از روش ماشین بردار پشتیبان نسبت به دیگر روش‌ها مزیت‌های بسیاری دارد. مهم‌ترین مزیت این روش پاسخ یکتا به‌ازای مجموعه داده آموزشی است. در حالی که الگوریتم‌های ژنتیک و شبکه عصبی تصادفی بوده و به‌ازای یک مجموعه داده خاص جواب یکتا تولید نمی‌کنند. به همین منظور روش طبقه‌بندی ماشین بردار پشتیبان در این بخش برای طبقه‌بندی داده مورد استفاده قرار گرفته است.

یکی از نکات اساسی در تشخیص تقلب هزینه تشخیص اشتباه نمونه‌های است. مسأله تشخیص تقلب در کارت‌های اعتباری از مسائلی است که تشخیص اشتباه در دسته‌های مختلف هزینه متفاوت هزینه متفاوتی دارد. برای مثال اگر یک نمونه داده تقلب باشد و سامانه باشتباه آن را داده عادی تشخیص دهد برای مشتری و بانک هزینه بالای خواهد داشت و تا حد امکان باید از این اتفاق جلوگیری کرد. همان‌طور که گفته شد به این نوع تشخیص، خطای منفی کاذب یا FN می‌شود. در عوض چنانچه یک نمونه داده تقلب نبوده و سامانه باشتباه تقلب اعلام کند، هزینه زیادی برای بانک و مشتری نخواهد داشت. درنهایت از طریق سامانه به بانک و مشتری اعلام اخطار می‌شود و مشتری با تأیید تراکنش‌های صورت گرفته مشکل را برطرف می‌کند. به عبارت دیگر به دنبال رویکردی هستیم که تعداد FN کمی داشته باشد. البته توجه کنید که کم کردن این تعداد در ازای ازبین‌بردن مقدار جزئی از دقت می‌شود. برای این مسأله رویکرد نوینی بر اساس ماشین بردار پشتیبان ارائه شده است که در ادامه شرح آن ذکر شده است. برای بیان رویکرد ارائه شده در ابتداء روش پایه

۲-۲ سامانه‌های خبره^۱

قوانين تشخیص تقلب را می‌توان از اطلاعات به دست آمده از یک انسان متخصص یا گروهی از خبرگان و یا ذخیره شده در یک سامانه مبتنی بر قانون مانند قوانین IF-THEN ذخیره شود، آن را یک سامانه پایگاه دانش^۲ یا یک سامانه خبره^۳ می‌نامیم. قوانینی که در سامانه خبره استفاده می‌شوند برای انجام عملیات بر روی اطلاعات با استنتاج به‌منظور رسیدن به نتیجه مناسب هستند [۳].

۲-۳ الگوریتم‌های ژنتیک^۴

الگوریتم‌های ژنتیک روش جستجو براساس روش‌های محاسبات تکاملی هستند. با توجه به جامعه‌ای از راه حل‌های فرضی مسأله (افراد)، محاسبات تکاملی با راه حل‌های جدید که به طور بالقوه بهتر نیز هستند، این جامعه را گسترش می‌دهند. در داده‌کاوی، الگوریتم‌های ژنتیک ممکن است برای خوبه‌بندی، پیش‌بینی و حتی قوانین انجمن استفاده شوند. این تکنیک‌ها می‌توانند برای پیداکردن مناسب‌ترین مدل از مجموعه‌ای از مدل‌ها برای نشان دادن داده‌ها مورد استفاده قرار گیرند.

الگوریتم ژنتیک روش را تکرار می‌کند تا زمانی که تعدادی از نسل از پیش مشخص شده را منتقل و بهترین راه حل را پیدا کند. این روش پارامتری است و برای به دست آوردن عملکرد بهتر نیاز به مسائل انجام شده دارد. فهرستی از این پارامترها و تنظیمات مورد نیاز برای تولید تراکنش‌های تقلبی هستند. پارامترهای مورد نیاز برای محاسبه مقادیر بحرانی، تعداد دفعات استفاده از کارت، محل استفاده از کارت، اضافه‌برداشت از کارت، موجودی فعلی بانکی، میانگین روزانه مصرف و برداشت و امثالی از این‌هاست [۴].

۴-۲ ماشین بردار پشتیبان^۵

در پژوهشی طبقه‌بندی تراکنش‌های جعلی و مشرع با استفاده از ماشین بردار پشتیبان به صورت دو رویی انجام گرفته است. هدف این پژوهش بررسی دقت و عملکرد کرنل‌های مختلف ماشین بردار پشتیبان اعم از SVC، C-SVC و SMO در تشخیص تقلب‌های صورت گرفته در کارت‌های اعتباری است.

¹Expert Systems

²Knowledge base

³Knowledge base system

⁴Expert system

⁵Genetic algorithms

⁶Support Vector Machine

مرز طبقه‌بندی باید تا حدی جایه‌جا شود که معیار تعريفشده کمینه شود. این معیار در ۲۰۱۲ (سیده نفیسه اسحاقی) نیز به کار رفته است و برای مقایسه نتایج روش ارائه شده نتایج این روش نیز ذکر و بررسی شده است.

برای تغییر مرز به سمت داده‌های طبقه عادی باید برآورده اولیه‌ای از حاشیه^۱ داشته باشیم. بعد از برآورد نسبت به حاشیه مرز تصمیم‌گیری را به سمت داده‌های عادی نزدیک می‌کنیم. در مرحله یادگیری تنها داده‌های آموزشی در اختیار داریم. از این رو برای ارزیابی مدل بعد از تعیین مرز تصمیم‌گیری نمی‌توان از داده‌های آزمون استفاده کرد. بنابراین برای ارزیابی مرز از روش اعتبارسنجی متقابل^۲ استفاده شده است. همان‌طور که در بخش مرور روش‌ها اشاره شد، ماشین بردار پشتیبان با کمینه‌سازی عبارت زیر عمل می‌کند^[۹]:

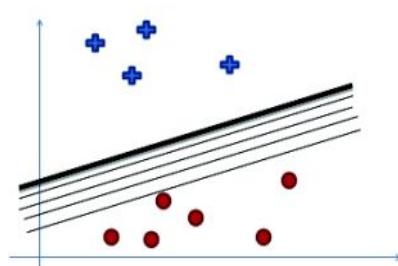
$$\text{Minimize} \quad (\text{in } w, b) \quad \frac{1}{2} \|w\|^2$$

$$y_i(w \cdot x - b) \geq 1 \quad \forall i = 1, \dots, n$$

در این عبارت $\|w\|$ مقدار حاشیه برای هر داده را تعیین می‌کند. همچنین در این عبارت مقدار b مرز تصمیم‌گیری را به سمت طبقه‌های متفاوت منحرف می‌کند.

در پیاده‌سازی‌ها از کتابخانه LIBSVM برای پیاده‌سازی الگوریتم استفاده شده است. برای استفاده از این کتابخانه ابتدا تمامی مجموعه آموزشی به الگوریتم یادگیری ارائه می‌شود. الگوریتم اجراسده و مقدار b تعیین می‌شود، سپس با استفاده از روش اعتبارسنجی متقابل مقادیر مختلف برای b مورد ارزیابی قرار می‌گیرد.

در شکل ۴ داده مقدار اصلی حاشیه با خط سیاه پرنگ و مقدارهای مورد ارزیابی با خطوط سیاه کمرنگ نشان داده شده است. با تعیین هر کدام از خطوط سیاه رنگ به عنوان مرز اعتبارسنجی انجام می‌شود. بعد از اعتبارسنجی و تعیین بهترین مرز با معیار ذکر شده این مرز انتخاب شده و در بخش آزمون برای داده‌های آزمون ارزیابی می‌شود.



شکل ۴- تعیین حاشیه با اعتبارسنجی متقابل

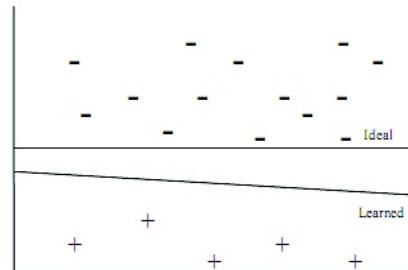
¹Margin

²Cross-Validation

ماشین بردار پشتیبان ذکر شده است. بعد از توضیح روش تغییر ایجادشده برای کمترکردن تعداد FN به تفصیل بیان شده است.

۱-۱-۳ تعیین حاشیه

هدف اصلی ماشین‌های بردار پشتیبان، یافتن خطی است که بیشترین حاشیه را برای دو طبقه ایجاد می‌کند. با استفاده از حاشیه ایجادشده، طبقه‌بندی برای دو طبقه با کمترین ریسک ممکن صورت می‌گیرد^[۵]. در شکل ۳ این خط برای یک مثال از داده‌های آموزشی نشان داده شده است.



شکل ۳- تعیین حاشیه

در شکل ۳ برچسب‌های منفی، داده‌های عادی و برچسب‌های مثبت، داده‌های تقلب را نشان می‌دهد. مرز تصمیم‌گیری ماشین بردار پشتیبان در حالت عادی خطی است که با عنوان Learned در شکل نشان داده شده است. این خط اگر چه از دید کلی و در یک مجموعه داده عمومی مرز مناسب برای تصمیم‌گیری است؛ ولی در مسأله تشخیص تقلب مرز مناسبی نیست.

در مسأله تشخیص تقلب باید مرز تصمیم‌گیری به سمت داده‌های عادی بیشتر منحرف شود. علت این انحراف کم کردن مقدار FN (تقلب از دست رفته) است. در شکل بالا مرز جدید تصمیم‌گیری با عنوان Ideal علامت‌گذاری شده است. مقدار انحراف خط باید توسط معيارهای ارزیابی عملکرد سامانه طبقه‌بندی تعیین شود. همان‌طور که ذکر شد، کم کردن FN موجب کم شدن دقت و همچنین افزایش خطای مثبت کاذب یا FP می‌شود. معياری باید تعریف شود که توازنی بین FN و FP ایجاد کند. طبق نتایج مشاهده شده در مجموعه داده‌های مذکور مقدار FN از FP همیشه بیشتر بوده است در حالی که هدف ما کم بودن FN است. به همین منظور معيار زیر برای برقراری توازن تعریف شده است.

$$m = 100 * FN + 10 * FP + TP$$

$$D = \frac{1}{|M|} \sum_{x_i \in M} I(x_i, c)$$

در رویکرد ارائه شده مقدار بالا باید بیشینه شود. همچنین مجموع اطلاعات بین ویژگی‌ها توسط رابطه زیر محاسبه می‌شود. این مجموع درواقع اطلاعات تکراری موجود در مجموعه ویژگی‌های انتخاب شده را نشان می‌دهد که باید کمینه شود.

$$R = \frac{1}{|M|^2} \sum_{x_i, x_j \in M} I(x_i, x_j)$$

با توجه به دو معیار بالا، معیاری با نام «بیشترین رابطه، کمترین افزونگی»^۲ طبق رابطه زیر تعریف می‌شود:

$$\Phi(D, R) = D - R$$

هدف بیشینه کردن کمیت بالاست. در مقاله [۳۱] اثبات شده است که با بیشینه کردن مقدار بالا ویژگی‌هایی به دست می‌آید که عبارت $I(M, c)$ را بیشینه می‌کند. همچنین اثبات شده است که اضافه کردن تدریجی ویژگی بر طبق معیار تعریف شده یعنی mrmr² معادل با بیشینه کردن رابطه اصلی است. به این ترتیب به جای بررسی تک تک حالت‌ها یعنی 2^n حالت می‌توان در n مرحله ویژگی‌های مناسب را اضافه کرد و به نتایج دلخواه رسید [۲].

۳-۳ اعمال اعتبارسنجی متقابل در ماشین بردار پشتیبان با حاشیه نامتوازن

۳-۳-۱ اعتبارسنجی متقابل

اعتبارسنجی متقابل، که گاهی تخمین گردشی نیز نامیده می‌شود، یک روش ارزیابی است که مشخص می‌کند نتایج یک تحلیل آماری بر روی یک مجموعه داده تا چه اندازه قابل تعمیم و مستقل از داده‌های آموزشی است. این تکنیک به طور ویژه در کاربردهای پیش‌بینی مورد استفاده قرار می‌گیرد تا مشخص شود مدل موردنظر تا چه اندازه در عمل مفید خواهد بود. به طور کلی یک دور از اعتبارسنجی متقابل شامل افزار داده‌ها به دو زیرمجموعه‌ها (داده‌های آموزشی) و اعتبارسنجی تحلیل با آن زیرمجموعه‌ها (داده‌های آموزشی) است (داده‌های اعتبارسنجی یا آزمون). برای کاهش پراکندگی، عمل اعتبارسنجی چندین بار با افزارهای مختلف انجام و از نتایج اعتبارسنجی‌ها میانگین گرفته می‌شود. هنگامی که جمع‌آوری داده‌های بیشتر سخت، پرهزینه و یا غیرممکن باشد. استفاده

۲-۳ انتخاب ویژگی‌ها برای طبقه‌بندی

در این پژوهش برای بهبود طبقه‌بندی از انتخاب ویژگی‌ها استفاده شده است. انتخاب ویژگی از دو لحاظ می‌تواند در طبقه‌بندی مفید باشد. جنبه نخست آن بالابردن سرعت و کم کردن پیچیدگی الگوریتم است. طبقه‌بندی ماشین بردار پشتیبان از نظر پیچیدگی زمانی به تعداد ویژگی‌ها وابسته است و کم کردن این ویژگی‌ها عملکرد سامانه را از لحاظ زمان بهبود می‌دهد.

جنبه دیگر، استفاده از انتخاب ویژگی بالابردن دقت در اثر حذف کردن ویژگی‌های نو فه دارد. درواقع در هر ویژگی مقداری اطلاعات و مقداری نو فه است. برای بیشتر ویژگی‌ها اطلاعات موجود در ویژگی بیش از مقدار نو فه است و دخیل کردن ویژگی در طبقه‌بندی باعث افزایش دقت طبقه‌بندی می‌شود؛ اما اگر مقدار اطلاعات در یک ویژگی کمتر از نو فه باشد، این ویژگی تنها طبقه‌بندی را گمراх می‌کند و حذف آن می‌تواند موجب بهبود دقت شود.

در روش ارائه شده مقدار اطلاعات متقابل برای هر ویژگی و برچسب طبقه‌بندی محاسبه می‌شود. مقدار اطلاعات، راهنمای ما برای انتخاب ویژگی خواهد بود و درنهایت ویژگی‌هایی انتخاب می‌شوند که بیشترین مقدار اطلاعات را داشته باشند. چنانچه ویژگی‌های استفاده شده در مجموعه n نشان دهیم $F = \{x_i | 1 \leq i \leq n\}$ تعداد ویژگی‌ها را نشان می‌دهد). هدف انتخاب مجموعه‌ای از ویژگی‌هایی است که در کنار هم بیشترین مقدار اطلاعات برای طبقه‌بندی را داشته باشد. فرض کنید $M \subseteq F$ مجموعه ویژگی‌های انتخاب شده باشد، آنگاه اطلاعات متقابل این مجموعه ویژگی توسط $I(M, c)$ تعیین می‌شود.

توجه کنید که در این عبارت c بردار برچسب‌های است. برای محاسبه مقدار اطلاعات به ازای هر مجموعه از ویژگی‌ها باید 2^n حالت متفاوت بررسی شود؛ این تعداد درواقع زیرمجموعه‌های مجموعه اصلی ویژگی‌های است. رویکرد جایگزین می‌تواند اضافه کردن ویژگی‌ها به صورت تک تک باشد. در موقع اضافه کردن تک تک ویژگی‌ها باید اطلاعات هر ویژگی و برچسب لحاظ شود. از طرف دیگر در مرحله اضافه کردن ویژگی باید تکراری نبودن اطلاعات بررسی شود. درواقع باید اطلاعات متقابل ویژگی‌ها نسبت به هم نیز بررسی شود. در مقاله رویکرد مناسبی برای برقراری تعادل بین این دو معیار معرفی شده است. مجموع اطلاعات متقابل برای ویژگی‌ها به صورت زیر تعریف می‌شود.

¹Feature Selection

۴- ارزیابی نتایج

یکی از چالش‌های پژوهش‌های انجام‌یافته در زمینه تشخیص تقلب، دسترسی به داده‌های مناسب برای ارزیابی روش‌های پیشنهادی است. بهدلیل محرومانگی داده‌ها در بانک‌ها و مؤسسه‌های مالی صادرکننده کارت اعتباری، بهسختی می‌توان به داده‌های واقعی دست یافته. ما در این پژوهش از داده‌های بانک بزرگی استفاده کردایم که دارای نه مجموعه داده است که ۳/۷۴٪ این داده‌ها تقلب است. این داده‌ها مربوط به بازه زمانی ۱۴ جولای ۲۰۰۴ تا ۱۲ سپتامبر ۲۰۰۴ است. این داده‌ها دارای هفده ویژگی است.

۱-۴ معیارهای ارزیابی

استفاده از تکنیک‌های تشخیص تقلب، درنهایت منجر به تخصیص تراکنش‌ها به یکی از دو طبقه عادی و تقلبی می‌شود؛ لذا روش‌های طبقه‌بندی باید براساس معیارهایی مورد ارزیابی قرار گیرند. معیارهای ارزیابی متعددی برای اعتبارسنجی این روش‌ها وجود دارد که هر کدام از جنبه خاصی کارایی تکنیک‌ها را می‌سنجدن. زمانی که سامانه تراکنش ورودی را برچسب‌گذاری می‌کند یکی از چهار حالت زیر پیش می‌آید:

^۱: تراکنش ورودی جعلی است، و سامانه آن را به درستی تقلب تشخیص داده است.

^۲: تراکنش ورودی عادی است، ولی سامانه آن را به عنوان تقلب تشخیص داده است.

^۳: تراکنش ورودی قانونی است، و سامانه آن را به درستی قانونی تشخیص داده است.

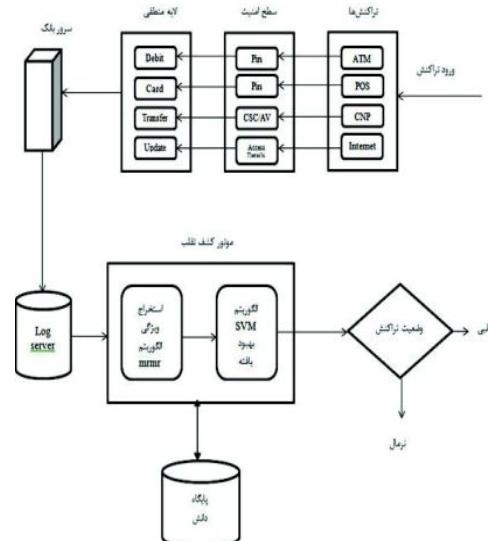
^۴: تراکنش ورودی جعلی است و سامانه آن را به اشتباه قانونی تشخیص داده است.

در حالت بهینه سامانه باید بتواند در عین این که تعداد TP را به حداقل می‌رساند، تعداد FN را در حداقل خود نگه دارد. در این حوزه تراکنش‌هایی که به اشتباه عادی تشخیص داده می‌شوند بسیار خطرناک‌تر از تراکنش‌هایی هستند که به اشتباه تقلبی تشخیص داده می‌شوند. در حالت نخست بانک یا مؤسسه مالی متحمل خسارت مالی قابل توجهی بسته به مبلغ تراکنش خواهد شد؛ اما در حالت دوم فقط هزینه اضافی بررسی تراکنش تحمیل می‌شود؛ لذا ارزش FP و FN با یکدیگر متفاوت است.

از اعتبارسنجی متقابل کمک می‌کند از فرضیات نادرست با داده‌های فعلی که قابل تعمیم نیستند، دوری شود.

K-Fold ۲-۳-۳

در این نوع اعتبارسنجی، داده‌ها به K زیرمجموعه افزار می‌شوند. از این K زیرمجموعه، هر بار یکی برای اعتبارسنجی و K-1 تای دیگر برای آموزش به کار می‌روند. این روال K بار تکرار می‌شود و همه داده‌ها به طور دقیق یکبار برای آموزش و یکبار برای اعتبارسنجی به کار می‌روند. درنهایت میانگین نتیجه این K بار اعتبارسنجی به عنوان یک تخمین نهایی برگزیده می‌شود. البته می‌توان از روش‌های دیگر برای ترکیب نتایج استفاده کرد. به طور معمول از 5-Fold استفاده می‌شود. در روش K-Fold طبقه‌ای سعی می‌شود نسبت داده‌های هر طبقه در هر زیرمجموعه و در مجموعه اصلی یکسان باشد.



شکل ۵- مدل سامانه پیشنهادی

در این پژوهش از روش 5-Fold برای تعیین بهترین مرز استفاده شده است. همان‌طور که گفته شد، حاشیه به پنج قسمت مساوی تقسیم می‌شود و مکان پیشنهادی برای مرز این پنج خط خواهد بود. برای هر کدام از مرزهای پیشنهادی ارزیابی توسط اعتبارسنجی متقابل انجام شده و میانگین آن به‌ازای Fold های متفاوت محاسبه می‌شود. از بین خطوط پیشنهاد شده به عنوان مرز طبقه‌بندی خطی که بهترین کارایی را از نظر Cost داشته باشد، به عنوان مرز انتخاب شده و در بخش آزمون، بر روی داده‌های آزمون ارزیابی می‌شود [۷].

^۱True Positive

^۲False Positive

^۳True Negative

^۴False Negative

$$\text{Accuracy}(A) = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision}(P) = \frac{TP}{TP + FP}$$

$$\text{Recall} = TP r = \frac{TP}{FN + TP}$$

$$F - \text{measure}(F) = \frac{2 * PR}{R + P}$$

$$\text{FP rate} = \frac{FP}{FP + TN}$$

معیارهای مطرح شده اعم از دقت، صحت و نرخ تشخیص تقلب، میزان درصد از پیش‌بینی صحیح را توصیف می‌کند. با این حال هیچ یک از آنها به تنها یک نمی‌تواند کارآیی یک الگوریتم را مشخص کند. دقت بالا تضمین نتیجه بهتر را نمی‌کند، چون هزینه پیش‌بینی‌های نادرست از داده‌های مثبت و منفی به طور معمول متفاوت است. نتیجه دقت بالا ممکن است نرخ تشخیص خوبی نداشته باشد؛ بنابراین باید بین نرخ TP، نرخ FP، نرخ TN و نرخ FN یک تعادل برقرار کرد.

در ادامه به شرح معیارهای مورد استفاده در این پژوهش پرداخته و جهت مقایسه عملکرد روش‌ها از یکتابع هزینه استفاده شده است. در این تابع، بازبینی هر تراکنش صحیح (TP) ۱ دلار و هر تراکنش بهاشتباه تقلب تشخیص داده شده (FP) ۱۰ دلار و هر تراکنش تقلیلی که سامانه نتوانسته آن را تشخیص دهد (FN) ۱۰۰ دلار هزینه در نظر گرفته شده است. بررسی پارامترهای اندازه‌گیری نظری صحت^۱، دقت^۲، نرخ تشخیص^۳ و ... نقش بسیار مهمی در تجزیه و تحلیل سامانه‌ها دارند. متناسب با نوع مشکل و زمینه کاربرد، یک اندازه‌گیری می‌تواند مناسب‌تر از دیگری باشد. به عنوان مثال در علوم رفتاری، ویژگی^۴ و حساسیت^۵ استفاده می‌شود. در علوم پزشکی تجزیه و تحلیل ROC یک استاندارد برای ارزیابی است. در حوزه تشخیص تقلب، میزان تشخیص تقلب، دقت و F-سنجهش به عنوان معیارهای مناسب جهت آزمایش کارآیی سامانه‌ها در نظر گرفته شده است [۱۲].

$$\text{Total cost} = 100 * FN + 10 * FP + 1 * TP$$

جدول ۱ میزان اطلاعات استخراجی در هر ویژگی

	1	2	3	4	5	6	7	8	9
1	0.012 1	0.0127	0.0145	0.0145	0.013	0.0132	0.0129	0.0138	0.0135
2	0.018 6	0.0193	0.0211	0.0206	0.0202	0.0213	0.0197	0.0194	0.02
3	0.003 7	0.0033	0.0033	0.0037	0.0031	0.0038	0.0039	0.0041	0.0039
4	0.015 6	0.0148	0.0153	0.0157	0.0158	0.0159	0.0156	0.0151	0.0146
5	0.019 4	0.0186	0.0181	0.0191	0.0189	0.02	0.0201	0.0192	0.0195
6	0.019 4	0.019	0.0184	0.0196	0.0194	0.02	0.019	0.0186	0.0181
7	0.010 7	0.0117	0.0122	0.0131	0.0121	0.012	0.0122	0.0123	0.013
8	0.001 8	0.0018	0.0018	0.0016	0.0017	0.0017	0.0018	0.0018	0.0015
9	0.000 7	0.0009	0.0008	0.0006	0.001	0.0009	0.0006	0.0008	0.0008
10	0.001	0.0009	0.0011	0.001	0.0011	0.0014	0.0009	0.001	0.0011
11	0.050 2	0.0498	0.0487	0.0474	0.0483	0.0523	0.0509	0.0513	0.0503
12	0.000 2	0.0001	0.0003	0.0002	0.0002	0.0003	0.0003	0.0001	0.0002
13	0.001 9	0.0017	0.002	0.0016	0.0017	0.0017	0.002	0.0017	0.0015
14	0.004 5	0.0043	0.0044	0.0048	0.0038	0.0055	0.0047	0.0043	0.005
15	0.007 8	0.0082	0.0081	0.0085	0.0078	0.0089	0.0083	0.0075	0.0085
16	0.028 8	0.0281	0.0278	0.0283	0.0275	0.0301	0.0296	0.0288	0.0281
17	0.004 6	0.0044	0.0045	0.0044	0.0046	0.0054	0.0041	0.0046	0.0052

¹Accuracy²Precision³Detection Rate⁴F-measure⁵Specificity⁶Sensitivity

زمانی که الگوریتم بالا با شرایط انتخاب ۱۵ ویژگی بر روی مجموعه داده‌ها اجرا می‌شود، ۱۵ ویژگی دارای اهمیت بالا (از نظر اطلاعات) از میان ۱۷ میان ویژگی برای داده‌های تقلیلی انتخاب شده و بعد کل داده‌ها با ۱۵ ویژگی به طبقه‌بندی ماشین بردار پشتیبان داده می‌شوند و درنهایت طبقه‌بندی با دقت بالایی صورت می‌گیرد.

در جدول ۲ پانزده ویژگی لازم برای تشخیص تقلب در هر مجموعه داده به تفکیک مشخص شده که این انتخاب توسط الگوریتم صورت گرفته است.

جدول ۲ ویژگی‌های استخراج شده برای تقلب

۱	۲	۳	۴	۵	۶	۷	۸	۹
11	11	11	11	11	11	11	11	11
4	7	7	7	7	7	7	7	7
5	4	4	4	4	4	4	4	4
7	5	5	5	5	5	5	5	5
16	16	16	16	16	16	16	16	16
17	17	1	1	17	17	1	17	17
1	1	17	17	1	1	13	1	1
13	6	13	6	6	6	17	13	6
6	13	6	13	13	13	6	6	13
14	14	2	14	2	14	14	14	14
2	2	14	2	14	2	2	2	2
15	15	15	15	15	15	15	15	15
8	8	10	10	10	10	8	8	10
3	3	3	3	3	3	3	3	3
10	10	8	8	8	8	12	10	8

جدول ۳ نتایج طبقه‌بندی ماشین بردار پشتیبان برای هر ۹ مجموعه داده بدون انتخاب ویژگی‌ها را نشان می‌دهد. از آنجایی که معیار دقت^۱ برای ارزیابی سامانه تشخیص تقلب و همچنین رویکرد ارائه شده کافی نیست، ستون‌های دیگر در این جدول معیارهای ارزیابی سامانه است که در هر سطر مقدار متناظر با آن برای مجموعه داده خاص ذکر شده است. سطر آخر جدول مقدار میانگین بهزادی هر معیار را نشان می‌دهد.

جدول ۳ نتایج الگوریتم ماشین بردار پشتیبان به تنهایی

set	Accuracy	FPrate	Precision	Recall	cost	FN	FP	TN	TP
1	96.93	0.17	0.592965	0.224762	41628	407	81	12053	118
2	96.97	0.15	0.6	0.2147	41451	406	74	11952	111
3	96.89	0.16	0.586022	0.209213	42079	412	77	11900	109
4	96.72	0.17	0.564356	0.213483	42994	420	88	11805	114
5	96.76	0.17	0.557895	0.200758	43146	422	84	11909	106
6	96.93	0.18	0.522472	0.188641	40943	400	85	11980	93
7	97	0.16	0.6	0.217647	40751	399	74	11851	111
8	97.01	0.15	0.569767	0.194831	41338	405	74	12080	98
9	96.84	0.19	0.539604	0.214145	41039	400	93	11817	109
AVG	96.89	0.17	0.570342	0.208687	41707.67	408	81	11927	108

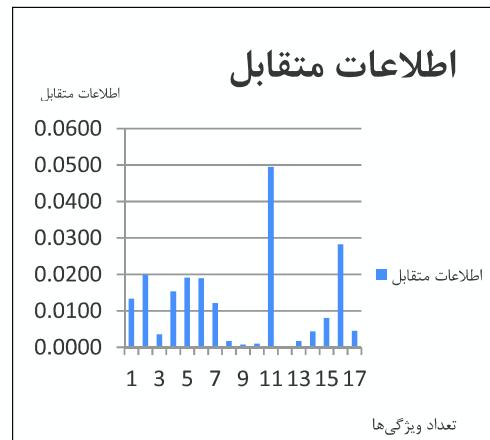
جدول ۴ نتایج طبقه‌بندی را بعد از انتخاب ۱۵ ویژگی و قبل از اعمال اعتبارسنجی متقابل، نشان می‌دهد.

'Accuracy

۲-۴ نتایج انتخاب ویژگی‌ها

در این بخش نتایج کم کردن تعداد ویژگی‌ها بیان شده است. قبل از انتخاب ویژگی‌ها، اطلاعات متقابل هر ویژگی محاسبه می‌شود. اطلاعات ویژگی‌ها برای هر مجموعه داده و تک تک ویژگی‌ها محاسبه شده است. مقادیر به دست آمده از ویژگی‌ها در جدول ۱ ذکر شده است.

در جدول ۱ هر سطر معادل یک ویژگی و هر ستون معادل یک مجموعه داده است. در محاسبه اطلاعات متقابل، تنها داده‌های آموزشی دخیل بوده است. مقادیر این جدول در واقع اطلاعات متقابل هر ویژگی و برچسب طبقه‌بندی (تقلب بودن یا نبودن) را بیان می‌کند. تعداد ویژگی‌ها و مجموعه داده‌ها بسیار زیاد است و این جدول نمی‌تواند به خوبی بیان‌گر ویژگی‌ها و اطلاعات موجود در آن باشد. به همین خاطر میانگین اطلاعات در ویژگی‌ها در نمودار شکل زیر رسم شده است. این نمودار ستون آخر از جدول بالا نشان می‌دهد.



نمودار ۱- اطلاعات متقابل ویژگی‌ها

همان‌طور که در جدول مشاهده می‌شود اطلاعات ویژگی‌ها متفاوت است. تعدادی از ویژگی‌ها حاوی مقدار بسیار کمی اطلاعات هستند؛ البته همان‌طور که گفته شد ویژگی‌ها باید با هم در نظر گرفته شوند. روش استفاده شده برای کم کردن نیز باید در نظر گرفته شود. روش mrmr است که در بخش ۳ توضیح داده شد. یکی از مسائل مهم تعداد ویژگی‌های انتخاب شده است. با توجه به میانگین اطلاعات ویژگی کم کردن تعداد ویژگی از یک ویژگی تا شش ویژگی می‌تواند مناسب باشد. به همین خاطر پاسخ الگوریتم کاهش ویژگی‌ها برای سه مقدار ۱۵، ۱۳ و ۱۱ ذکر شده است.

جدول ۶ نتایج اجرای الگوریتم با استخراج ۱۱ ویژگی

set	Accuracy	Fprate	Precision	Recall	cost	FN	FP	TN	TP
1	96.97	0.19	0.516667	0.195789	39163	382	87	12097	93
2	96.97	0.18	0.497041	0.179872	39234	383	85	11991	84
3	96.85	0.18	0.493902	0.171975	39911	390	83	11944	81
4	96.68	0.16	0.5	0.163223	41369	405	79	11864	79
5	96.84	0.17	0.506173	0.171548	40482	396	80	11963	82
6	97.01	0.2	0.460606	0.171558	37666	367	89	12026	76
7	96.94	0.16	0.54717	0.18913	38107	373	72	11903	87
8	97.02	0.19	0.50289	0.192053	37547	366	86	12118	87
9	96.83	0.19	0.533333	0.20915	37236	363	84	11876	96
AVG	96.9	0.18	0.50642	0.1827	38968.33	381	83	11976	85

۳-۴ نتایج روش ماشین بردار پشتیبان با حاشیه نامتوازن و اعمال اعتبارسنجی متقابل
در ابتدای این بخش نتایج سامانه تشخیص تقلب برای مجموعه داده بزرگی قبل و بعد از اعمال الگوریتم انتخاب ویژگی بیان شده است. همان‌طوری که گفته شد معیار دقت^۱ برای ارزیابی سامانه تشخیص تقلب و همچنین رویکرد ارائه شده کافی نیست. علاوه‌بر این معیار، معیار هزینه نیز تعریف شده است. در جدول ۳ نتایج روش ماشین بردار پشتیبان بدون استفاده از رویکرد ارائه شده برای تغییر مرز طبقه‌بندی نشان داده شده است.

همان‌طور که مشاهده می‌شود مقدار FN در این جدول مقدار بالایی است؛ و حتی در مقایسه با FP بیش از دو یا سه برابر است. این مقدار نشان می‌دهد بسیاری از تراکنش‌های تقلب از دید سامانه پنهان بوده و سامانه این تراکنش‌ها را عادی پنداشته است. در حالی که هزینه این اشتباه برای بانک و مشتری هزینه بالایی است. جدول ۷ نتایج معیارهای ارزیابی را بعد از یافتن بهترین مرز تصمیم‌گیری توسط روش اعتبارسنجی متقابل نشان می‌دهد. سطرها و ستون‌های جدول مشابه جدول ۳ است و مقدارهای آن براساس رویکرد ارائه شده بیان شده است.

جدول ۷ نتایج اجرای الگوریتم با اعتبارسنجی متقابل و انتخاب ۱۵ ویژگی

set	Accuracy	FPrate	Precision	Recall	cost	FN	FP	TN	TP
1	95.98	0.64	0.495895	0.635789	20672	173	307	11877	302
2	95.77	0.69	0.480122	0.672377	19014	153	340	11736	314
3	96.05	0.7	0.496206	0.694268	18047	144	332	11695	327
4	95.82	0.63	0.489431	0.621901	21741	183	314	11629	301
5	95.88	0.67	0.487654	0.661088	19836	162	332	11711	316
6	96.08	0.64	0.48	0.623025	19966	167	299	11816	276
7	96.25	0.67	0.4736	0.643478	19986	164	329	11646	296
8	95.92	0.67	0.466989	0.640177	19900	163	331	11873	290
9	95.68	0.65	0.469522	0.620915	20905	174	322	11638	285
AVG	95.94	0.66	0.482158	0.645891	20007.44	165	323	11736	301

^۱Accuracy

جدول ۴ نتایج اجرای الگوریتم با استخراج ۱۵ ویژگی

set	Accuracy	Fprate	Precision	Recall	cost	FN	FP	TN	TP
1	96.97	0.18	0.606218	0.246316	36677	358	76	12108	117
2	96.97	0.18	0.597884	0.24197	36273	354	76	12000	113
3	96.85	0.18	0.572973	0.225053	37396	365	79	11948	106
4	96.68	0.19	0.556701	0.22314	38568	376	86	11857	108
5	96.84	0.18	0.581633	0.238494	37334	364	82	11961	114
6	97.01	0.2	0.547872	0.232506	34953	340	85	12030	103
7	96.94	0.17	0.58427	0.226087	36444	356	74	11901	104
8	97.02	0.18	0.572222	0.227373	35873	350	77	12127	103
9	96.83	0.21	0.537313	0.235294	36138	351	93	11867	108
AVG	96.9	0.18	0.57301	0.232915	36628.44	357	81	11978	108

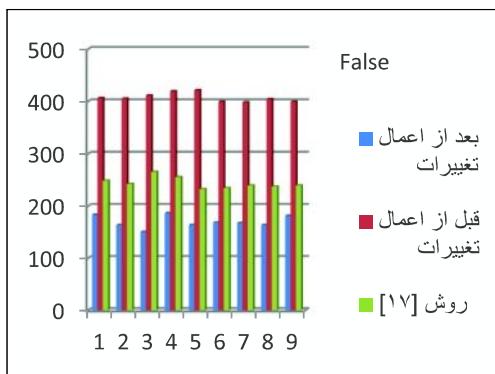
همان‌طور که مشاهده می‌شود، نتیجه دقت و هزینه هر دو نسبت به قبل از حذف ویژگی بهبود کمی یافته است. این بهبود در مورد هزینه مشهودتر است. علاوه‌بر آن، این نتایج نشان می‌دهد که دو ویژگی حذف شده، ویژگی‌هایی ارزشی بوده‌اند و تنها هزینه طبقه‌بندی را افزایش داده‌اند. در جدول پیوست شرح ویژگی‌ها بیان شده است. با توجه به مشاهده ویژگی‌های انتخاب شده توسط الگوریتم mmrm مشاهده کردیم که در تمام آزمایش‌ها بر روی مجموعه داده‌ها ویژگی‌های ۹ و ۱۲ حذف شده است. در نمودار، اطلاعات متقابل نیز مشاهده می‌شود که این ویژگی‌ها اطلاعات کمتری داشته‌اند. جدول ۵ نتایج انتخاب ویژگی با ۱۳ ویژگی را نشان می‌دهد. همان‌طور که دیده می‌شود حذف ویژگی‌های بیشتر، موجب کاهش دقت و بالارفتن هزینه شده است. همچنین با مشاهده ویژگی‌های حذف شده در مجموعه داده‌ها، مشخص شد که در اکثر مجموعه داده‌ها ویژگی‌های ۱۰ و ۸ حذف شده‌اند. با این حال حذف کردن این ویژگی‌ها برای طبقه‌بندی مناسب و مقرن به صرفه نبوده است.

جدول ۵ نتایج اجرای الگوریتم با استخراج ۱۳ ویژگی

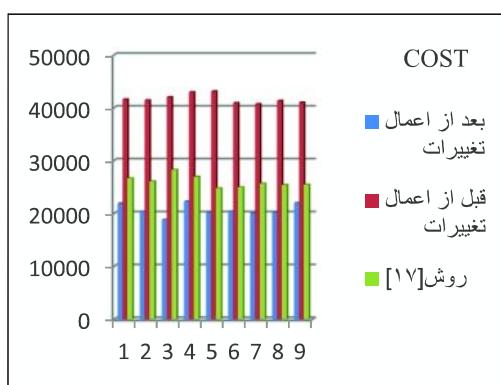
set	Accuracy	FPrate	Precision	Recall	cost	FN	FP	TN	TP
1	96.97	0.17	0.586957	0.227368	37568	367	76	12108	108
2	96.97	0.18	0.597884	0.24197	36273	354	76	12000	113
3	96.85	0.17	0.5625	0.210191	38069	372	77	11950	99
4	96.68	0.18	0.554348	0.210744	39122	382	82	11861	102
5	96.84	0.17	0.579235	0.221757	38076	372	77	11966	106
6	97.01	0.19	0.505814	0.196388	36537	356	85	12030	87
7	96.94	0.18	0.566667	0.221739	36682	358	78	11897	102
8	97.02	0.17	0.6	0.238411	35328	345	72	12132	108
9	96.83	0.19	0.543478	0.217865	36840	359	84	11876	100
AVG	96.9	0.18	0.56632	0.220715	37166.11	363	79	11980	103

در ادامه جدول ۶ عملکرد طبقه‌بندی با ۱۱ ویژگی را نشان می‌دهد.

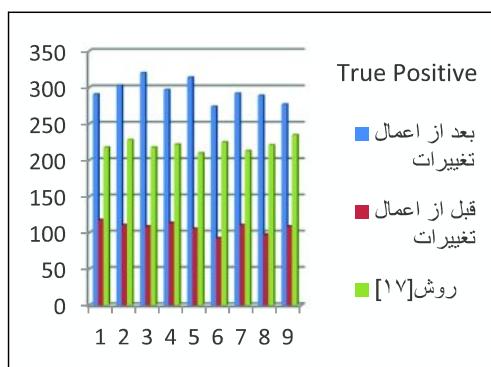
در ادامه چهار نمودار برای مقدارهای FN, TP, Cost Accuracy, Cost نشان داده شده است. در این نمودارها هر سه روش، یعنی ماشین بردار پشتیبان بدون اعمال اعتبارسنجی متقابل، ماشین بردار پشتیبان با اعمال اعتبارسنجی متقابل و روش [۱۲] ذکر شده است.



نمودار ۲- نمودار مقایسه‌ای False Negative برای سه روش



نمودار ۳- نمودار مقایسه‌ای Cost برای سه روش



نمودار ۴- نمودار مقایسه‌ای True Positive برای سه روش

با اعمال رویکرد ارائه شده مقدار به طور تقریبی به نصف مقدار قبلی کاهش پیدا کرده است. عملکرد روش ارائه شده در مقدار Cost یا هزینه نیز مشهود است. مقدار هزینه براساس روش اصلی SVM که در جدول ۳ ذکر شد در حالت میانگین ۴۱۷۰۷ بوده، در حالی که با اجرای رویکرد ارائه شده این مقدار بهشت کاهش یافته و به ۲۰۴۵۱ رسیده است. البته لازم به ذکر است که با تعییر حاشیه، مقدار FP زیاد می‌شود. مقدار FP در واقع داده‌هایی است که تقلیب نبوده و سامانه باشتباه تقلب تشخیص داده است. در دنیای واقعی این اشتباه هزینه بالایی برای مشتری و بانک ندارد و می‌توان گفت نسبت به FN هزینه آن ناچیز است. برای مشاهده کارایی روش ارائه شده نتایج آن با روش [۱۲] نیز مقایسه شده است. البته نتایج مرجع بالا ترکیبی از ماشین بردار پشتیبان با شبکه‌های عصبی است. جدول ۸ نتایج ترکیب دو روش svm با شبکه عصبی توسط [۱۲]

جدول ۸ نتایج ترکیب دو روش svm با شبکه عصبی توسط [۱۲]

set	Accuracy	FPrate	Precision	Recall	cost	FN	FP	TN	TP
1	96.46	0.48	0.525727	0.503212	25555	232	212	11864	235
2	96.42	0.51	0.524793	0.539278	24254	217	230	11797	254
3	96.37	0.5	0.534304	0.530992	25197	227	224	11719	257
4	96.86	0.48	0.551948	0.533473	24625	223	207	11836	255
5	96.99	0.48	0.55711	0.539503	22539	204	190	11925	239
6	96.35	0.46	0.599537	0.563043	22089	201	173	11802	259
7	96.69	0.56	0.492063	0.547461	23308	205	256	11948	248
8	96.8	0.49	0.554566	0.542484	23249	210	200	11760	249
9	97.1	0.48	0.575758	0.56	23126	209	196	11988	266
Avg	96.67	0.49	0.546201	0.539939	23771.33	214.22	209.78	11848.78	251.33

در جدول ۹ مقایسه روش ارائه شده با دیگر روش‌های انجام شده، مقایسه شده است. روش ارائه شده در مقایسه با روش‌های قبلی نیز هزینه کمتر و تعداد FN کمتری دارد. توجه کنید که در روش [۱۲] نیز از تابع هزینه مشابهی استفاده شده است و هدف از ارائه آن طبقه‌بندی بهمنظور کم کردن هزینه در این تابع هدف بوده است. در جدول ردیف یک روش ارائه شده در این پژوهش، روش دوم و سوم مربوط به روش [۱۲] و روش چهارم مربوط به [۵] است.

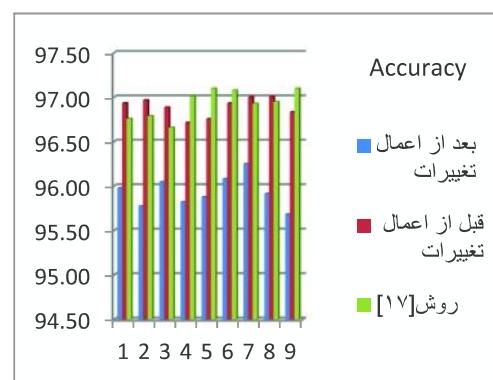
جدول ۹ مقایسه روش ارائه شده با سایر روش‌ها

روش	Accuracy	FPrate	Precision	Recall	cost	FN	FP	TN	TP
1	95.94	0.66	0.482158	0.645891	20007	165	323	11736	301
2	96.93	0.36	0.614503	0.475036	26056	244	139	11919	221
3	96.67	0.49	0.546201	0.539939	23771	214	209	11848	251
4	76.50	*	*	*	*	*	*	*	*

داشت. همچنین به این نتیجه رسیدیم که حذف بیش از حد ویژگی‌ها نیز باعث نتایج نامطلوب خواهد شد. نتایج این پژوهش نشان می‌دهد که اعمال اعتبارستجوی متقابل در الگوریتم ماشین بردار پشتیبان باعث کاهش تعداد FN شده و در نهایت هزینه‌ها را برای یک مؤسسه صادرکننده کارت اعتباری به طور قابل توجهی کاهش می‌دهد.

۶- منابع

- [1]. Adnan M. Al-Khatib;" Electronic Payment Fraud Detection Techniques" ;World of Computer Science and Information Technology Journal (WCSIT)ISSN: 2221-0741Vol. 2, No. 4, 137-141; 2012.
- [2]. HanchuanPeng, Fuhui Long, and Chris Ding" Feature Selection Based on Mutual Information: Criteria of Max-Dependency, Max-Relevance, AndMin-Redundancy", IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 27, No. 8, August 2005.
- [3]. Kevin J. Leonard; "The development of a rule based expert system model for fraud alert consumer credit"; European journal of operational research, vol. 80, p.p. 350-356; 2012.
- [4]. K.RamaKalyani, D.UmaDevi;" Fraud Detection of Credit Card Payment System by Genetic Algorithm ";International Journal of Scientific & Engineering Research Volume 3, Issue 7, July-2012 ISSN 2229-5518.
- [5]. MaryamsadatHejazi, Yashwant Prasad Singh;" Credit Data Fraud Detection using Kernel Methods with Support Vector Machine";Journal of Advanced Computer Science and Technology Research 2 (2012) 35-49ISSN: 2231-8852.
- [6]. Mati as Dimartino, Guzman Hernandez, Marcelo Fiori, Alicia Fernandez," A new framework for optimal classifier design",2013
- [7]. Qi Li *, Raied Salman, Erik Test, Robert Strack, Vojislav Kecman"Parallel multitask cross validation for Support Vector Machine using GPU", © 2012 Elsevier Inc. All rights reserved
- [8]. RaghavendraPatidar, Lokesh Sharma,"Credit Card Fraud Detection Using Neural Network",International Journal of Soft Computing and Engineering(IJSCE) ISSN:2231-2307,Volume-1,Issue-NCAI 2011;June 2011.
- [9]. R.C. Chen, T.S. Chen, C.C. Lin, "A new binary support vector system for increasing detection rate of credit card fraud", International Journal of



نمودار ۴- نمودار مقایسه‌ای Accuracy برای سه روش

همان‌طور که مشاهده می‌شود در تمامی نمودارها به جز دقت، روش ارائه شده بهتر عمل می‌کند. در مورد دقت توجه کنید که تغییر در حدود یک درصد است و همچنین این تغییر با هدف بهترشدن هزینه انجام شده است. به عبارت دیگر معیار دقت در این پژوهش اهمیت کمتری نسبت به هزینه دارد.

۵- نتیجه‌گیری

در این پژوهش به بررسی و مرور پژوهش‌های مرتبط با حوزه تشخیص تقلب در کارت‌های اعتباری پرداخته شد. بیان شد که تشخیص تقلب در این حوزه به معنای فرآیند طبقه‌بندی تراکنش‌ها به دو دسته مشروع و جعلی است. به طور کلی هدف از تشخیص تقلب، بهبود شناسایی و پیش‌بینی‌های درست و حفظ پیش‌بینی‌های نادرست در یک سطح قابل قبول از هزینه است. هیچ یک از سامانه‌های شناسایی تقلب به تنهایی نمی‌تواند تمام شکل‌های تقلب را شناسایی و پوشش دهد. در مسائل شناسایی تقلب باید به نامتوازن‌بودن نسبت تراکنش‌های قانونی به تقلیلی توجه نمود و برای آن چاره‌ای اندیشید؛ چون در کارایی مدل بسیار مؤثر است.

در این پژوهش ابتدا روش طبقه‌بندی با استفاده از ماشین بردار پشتیبان همراه با انتخاب ویژگی و اعتبارستجوی متقابل را ارائه دادیم. در سامانه پیشنهادی از الگوریتم mrmr برای انتخاب ویژگی‌های ضروری داده‌های تقلیلی در مجموعه داده مورد بررسی، استفاده شد. در بررسی‌های انجام شده برای انتخاب ویژگی از هفده ویژگی داده‌های تقلیلی در بهترین وضعیت خود حداقل دو ویژگی از این هفده ویژگی که دارای کمترین اطلاعات نسبت به ویژگی‌های دیگر می‌باشند، قبل حذف است و ما بهترین نتیجه را با پانزده ویژگی خواهیم

Pattern Recognition and Artificial Intelligence,
vol. 20, no. 2, 2006, pp. 227–239.

- [10]. Salvatore J. Stolfo, David W. Fan, Wenke Lee and Andreas L. Prodromidis; "Credit Card Fraud Detection Using Meta-Learning: Issues and Initial Results"; Department of Computer Science- Columbia University; 1997.
- [11]. Xu Wei, Liu Yuan," An Optimized SVM Model for Detection of Fraudulent Online Credit Card Transactions", International Conference on Management of e-Commerce and e-Government, 978-0-7695-4853-1/12 \$26.00 © 2012 IEEE.

[۱۲]. سیده نفیسه اسحاقی، ۱۳۹۲. «رایه طرح جامع جهت تشخیص تقلب در کارت‌های اعتباری و بانکداری الکترونیک»، دانشگاه تهران پر迪س بین الملل کیش



بابک رحمانی اخذ مدرک کارشناسی در رشته کامپیوتر نرم‌افزار در بهمن ماه سال ۱۳۸۹ از دانشگاه علمی کاربردی کیشو مدرک کارشناسی ارشد در رشته مهندسی فناوری اطلاعات گرایش امنیت اطلاعات در بهمن ماه سال ۱۳۹۲ از دانشگاه تهران پر迪س بین الملل کیش. زمینهٔ پژوهشی مورد علاقهٔ تشخیص تقلب در تراکنش‌های مالی، تشخیص نفوذ در شبکه‌های حیاتی و مراکز امنیت اطلاعات (SOC) و رمز نگاری متقارن است.



حسین قرایی تحصیلات خود را در مقطع کارشناسی مهندسی برق الکترونیک در دانشکده مهندسی برق دانشگاه خواجه نصیرالدین طوسی در سال ۱۳۷۷ به اتمام رساند. مقاطع کارشناسی ارشد و دکتری مهندسی برق - الکترونیک در دانشکده مهندسی برق دانشگاه تربیت مدرس را به ترتیب در سال‌های ۱۳۷۹ و ۱۳۸۸ تکمیل کرد. از سال ۱۳۸۱ تاکنون عضو هیئت علمی مرکز تحقیقات مخابرات است و زمینه‌های پژوهشی مورد علاقهٔ ایشان طراحی مدارات دیجیتال، آنالوگ و سیگنال مختلط، پردازش VLSI سیگنال‌های دیجیتال، سامانه‌های تشخیص و پیش‌گیری از نفوذ و مرکز عملیات امنیت و اجزای آن هستند.

افتا
منادی
علمی ترویجی
دوفصلنامه