

کشف کلاهبرداری بانکی با استفاده از ماشین بردار پشتیبانی

استاد محترم :

جناب آقای امیر شکری

دانشجو:

فاطمه میرزاده

شماره دانشجویی: 40011920009

آذر 1400

خلاصه

با توسعه چشمگیر ارتباطات و محاسبات، کلاهبرداری بانکی در اشکال و مقادیر آن در حال رشد است. در این مقاله، ما اشکال مختلف کلاهبرداری را که بانک‌ها در معرض آنها قرار می‌گیرند و ابزارهای داده‌کاوی در معرض آن قرار می‌گیرند، تجزیه و تحلیل می‌کنیم که به داده‌های تشخیص زودهنگام آن که قبلاً در یک بانک انباشته شده است، منجر می‌شود. ما از روش‌های یادگیری نظارت‌شده از ماشین‌های بردار پشتیبانی با Spark (SVM-S) برای ساختن مدل‌هایی استفاده می‌کنیم که رفتار عادی و غیرعادی مشتری را نشان می‌دهند و سپس از آن برای ارزیابی اعتبار تراکنش‌های جدید استفاده می‌کنیم. نتایج به‌دست‌آمده از پایگاه‌های اطلاعاتی تراکنش‌های کارت اعتباری نشان می‌دهد که این تکنیک‌ها در مبارزه با کلاهبرداری بانکی در داده‌های بزرگ مؤثر هستند. نتایج آزمایش از این مطالعه نشان می‌دهد که SVM-S عملکرد پیش‌بینی بهتری نسبت به شبکه‌های انتشار برگشتی (BPN) دارد. علاوه بر میانگین پیش‌بینی، دقت زمانی که نسبت داده‌ها به 0.8 می‌رسد به حداکثر می‌رسد.

معرفی

سازمان‌های خدمات مالی دارای تعدادی اهداف استراتژیک از جمله جذب و حفظ مشتریان جدید و فعلی از طریق بکارگیری روش‌های مدیریتی مختلف هستند. با توجه به این اهداف، موسسات روزانه حجم زیادی از داده‌های پروفایل، تاریخچه خرید و مرور و داده‌های رسانه‌های اجتماعی را تولید می‌کنند. بانک‌ها به‌عنوان یک مؤسسه مالی، حجم عظیمی از داده‌ها را از منابع مختلف از مشتریان خود تولید می‌کنند و این به نیاز به داده‌های بزرگ طبق ویکی‌پدیا، (2016) کمک کرده است. ماشین بردار پشتیبانی می‌تواند به کاهش ریسک و بهبود کیفیت خدمات ارائه شده به مشتریان برای موفقیت در تجارت کمک کند. تقلب به یک خطر بسیار مهم برای مواجهه با موسسات مالی، اتحادیه‌های اعتباری و به ویژه بانک‌ها تبدیل شده است. مبارزه با کلاهبرداری با تکنیک‌های پیشگیری سنتی مانند پین، رمز عبور و سیستم‌های شناسایی همراه است، اما در سیستم‌های بانکی

مدرن ناکافی شده‌اند. بانک‌ها در چندین فعالیت با کلاهبرداری مواجه شدند اما به نظر می‌رسد استفاده از راه دور از اعتبار آسیب پذیرترین آنها باشد.

برنامه‌های کاربردی داده‌های بزرگ با تکنیک‌های داده‌کاوی می‌توانند نقش مهمی در مبارزه با این نوع تقلب‌ها داشته باشند. داده‌کاوی مجموعه‌ای از تکنیک‌ها برای استخراج اطلاعات مهم از مقادیر زیادی داده برای کمک به تصمیم‌گیری است Spark. یک ابزار کلان داده است که به‌ویژه در این زمینه استفاده می‌شود، زیرا به دلیل تکنیک‌های یادگیری ماشینی متعدد و جریان‌سازی در زمان واقعی است. شاغل شدن در موسسات مالی موثر است. جرقه با روش SVM برای این نوع کلاهبرداری‌ها بهترین است. چندین تکنیک برای کشف تقلب توسط بسیاری از محققین پیشنهاد شده و مورد استفاده قرار گرفته است، از جمله تقلب در کارت‌های اعتباری. از جمله این تکنیک‌های داده‌کاوی، شبکه‌های بیزی، زنجیره‌های مارکوف، شبکه‌های عصبی، رگرسیون خطی، هم‌ترازی توالی و غیره هستند. هدف این کار ارائه معماری تشخیص تقلب است که بانک را قادر می‌سازد تا تراکنش‌های متقلبانه را در زمان واقعی با جرقه بر اساس ماشین‌شناسایی کند. ماشین‌بردار پشتیبانی تکنیک یادگیری SVM. که برای تشخیص چهره، شناسایی اثر انگشت، تشخیص صدا و کارهای مشابه بسیار قدرتمند است. هدف ما رفع مشکل کلاهبرداری در بانک‌ها و حل آن از طریق Spark با تکنیک‌های SVM است. ما تحلیلی از مشکل کلاهبرداری بانک‌ها و برای هر نوع طراحی، گونه‌ای از SVM را ارائه می‌کنیم که می‌تواند برای راه‌حل آن و سازگاری‌های لازم استفاده شود.

بقیه مقاله به شرح زیر سازماندهی شده است: ابتدا انواع کشف تقلب و همچنین شاخص‌های مورد استفاده برای کشف آن ارائه می‌شود، سپس انواع راه‌حل‌های کلان داده را که می‌توان مورد استفاده قرار داد مورد بحث قرار می‌دهیم. در بخش سوم استفاده از ماشین‌بردار پشتیبان در جرقه برای رفع نیازهای تشخیص تقلب را مورد بحث قرار می‌دهیم. علاوه بر این، در بخش چهارم اعتبار سنجی راه‌حل‌های پیشنهادی را با آزمایش آنها بر روی پایگاه‌های داده بانک ارائه می‌کند. ما مقاله را با یک نتیجه‌گیری و توصیه به پایان می‌بریم.

II. اشکال کلاهبرداری بانکی و شاخص‌های آنها.

کلاهبرداری در بانک به اشکال مختلف ظاهر می شود. می تواند داخلی باشد که توسط کارمندان خود بانک متعهد شده باشد یا خارجی توسط مشتریان، افراد یا ارگان های خارجی به بانک متعهد شود. ما علاقه مندیم که در این مقاله تقلب خارجی را در نظر بگیریم.

الف- پولشویی

پول شویی نیز یکی از روش های شناخته شده کلاهبرداری است، لوای بین المللی علیه این فعالیت توسط کشورهای مختلف برای کشف و پیگرد قانونی فعالیت های مجرمانه انجام می شود. مبارزه با پولشویی در صنعت مالی مبتنی بر تجزیه و تحلیل و پردازش اظهارات مربوط به تراکنش های مشکوک کشف شده توسط مؤسسات مالی است. به طور کلی، تنها چند تراکنش مشکوک واقعاً عملیات پولشویی هستند، اما تعداد عملیاتی که باید توسط مؤسسات مالی تجزیه و تحلیل شود، به زمان طولانی نیاز دارد. در ادبیات، روش های هوش مصنوعی ممکن است برای بهبود توانایی مؤسسات مالی در پردازش خودکار داده های مشکوک استفاده شود. با این حال، جستجو برای روش های کارآمد برای شناسایی رفتارهای مبادلاتی مشکوک پولشویی، یک زمینه تحقیقاتی بسیار فعال است. تعیین متغیرها و شاخص های پولشویی در حال حاضر آسان است، زیرا چنین فعالیت های غیر رسمی در شرایط پیچیده اقتصادی اجتماعی ایجاد می شوند. موارد زیر چند شاخص سفیدکننده پول مورد استفاده در ادبیات هستند. مبلغ تراکنش در صورتی که از مبلغ از پیش تعیین شده توسط بانک بیشتر باشد، معامله موجه نباشد، مشکوک است. به عنوان مثال، در غنا حداکثر مبلغ خرید با اعتبار شما نباید از 5000 یورو تجاوز کند. موارد زیر بیشتر مورد بررسی قرار می گیرد.

1. منابع انتقال

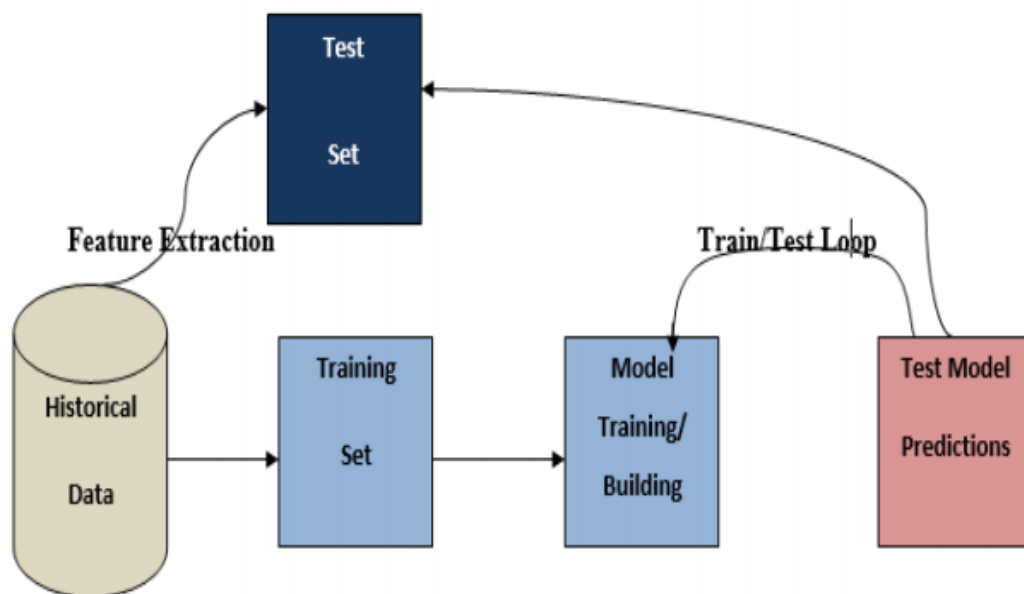
2. تاریخ معامله

3. تغییر آدرس

4. زمان معامله، معامله شبانه با مبلغ زیاد مشکوک است

ب. تقلب بر اساس کارت های اعتباری

کلاهبرداری مبتنی بر کارت اعتباری همچنان در حال افزایش است. بانک ها و شرکت های مالی سالانه مبالغ هنگفتی را از طریق کلاهبرداری با استفاده از کارت اعتباری از دست می دهند. تشخیص تقلب در کارت اعتباری اغلب بر اساس تعدادی از شاخص های پیش بینی است که عموماً از اطلاعات تراکنش های بازبایی شده از پایگاه داده تاریخی به دست می آیند. ما شاخص هایی مانند استفاده مکرر از کارت، مانده مانده پرداخت نشده هر چرخه، حداکثر تعداد روزهای تاخیر، دفعات خرید، تراکنش روزانه، بیشترین تعداد دفعات در پایگاه داده تاریخی و غیره را بررسی می کنیم. این ویژگی ها برای هر تراکنش استخراج می شوند و برای کشف الگوهای متقلبانه ثبت می شوند مدل تشخیص تقلب در شکل 1 نشان داده شده است.



شکل 1: مدل تشخیص تقلب

مدل داده پیشنهادی بر اساس داده های تاریخی موجود در انبار بانک، با استفاده از الگوریتم پیشنهادی که از SVM-S تشکیل شده است برای بررسی اینکه آیا نتیجه تقلبی است یا خیر، ساخته شده است. مجموعه ای از داده های مشابه برای پیش بینی نتیجه برای اثربخشی مدل استفاده خواهد شد. این مدل برای ارزیابی یک تراکنش جدید استفاده می شود. تراکنش پذیرفته شده توسط مدل اجرا می شود و سپس برای بهبود مدل به پایگاه داده اضافه می شود. تراکنش های رد شده

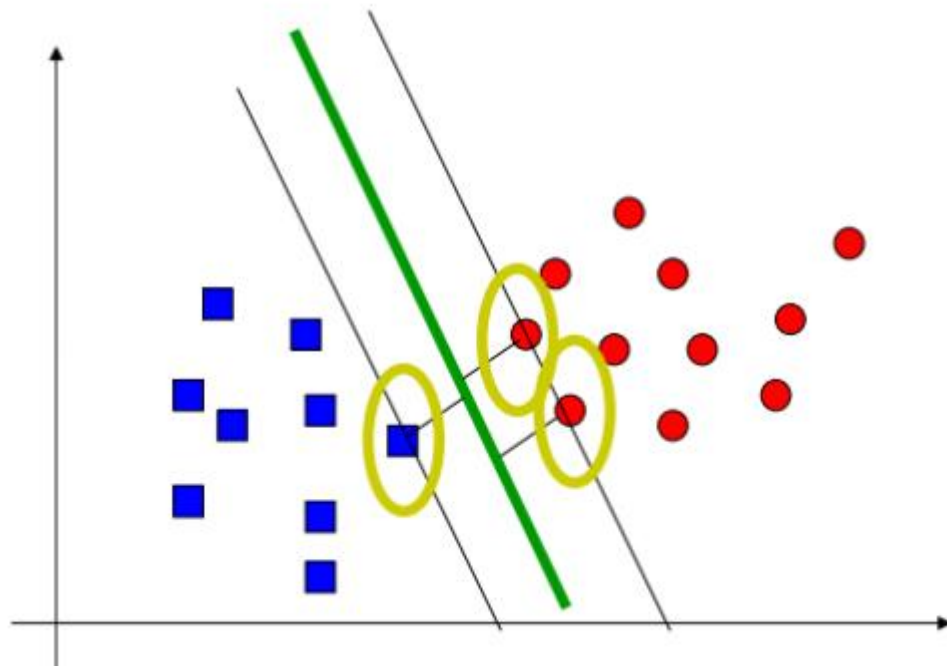
توسط مدل اجرا نمی شوند بلکه به عنوان مشکوک علامت گذاری می شوند. اگر تراکنش عادی باشد، همانطور که در بالا گفته شد اجرا و اضافه می شوند.

۱۱۱. داده های بزرگ برای کشف کلاهبرداری بانکی

به دلیل تغییر سریع و توسعه تکنیک های مورد استفاده توسط کلاهبرداران، ابزارهای داده کاوی دیگر نمی توانند رفتارهای غیرعادی را تجزیه و تحلیل کنند. داده های بزرگ، در این زمینه، با تکنیک های یادگیری ماشینی برای تشخیص تقلب در پایگاه داده ارائه می شوند که بهترین کار برای مبارزه با کلاهبرداری بانکی است. در ادبیات، از دو شکل یادگیری خودکار استفاده می شود: نظارت شده و بدون نظارت. روشهای یادگیری نظارت شده به عنوان قوانین ارتباطی شبکه های بیزی استفاده شده است. این روش ها دانش قبلی از ماهیت معاملات، تقلبی یا واقعی را فرض می کنند. یادگیری در این مورد شامل ساخت مدلی است که فضا را با توجه به نمونه های موجود به دو قسمت تقسیم می کند و سپس نمونه های جدید را بر اساس عضویت آنها در یکی از این دو کلاس طبقه بندی می کند. روش های بدون نظارت مانند شبکه های عصبی نیازی به طبقه بندی قبلی از نمونه های آموزشی ندارد. بیشتر بر اساس تشخیص تراکنش های عجیب و غریب است. برنامه های کاربردی داده های بزرگ، مانند Cassandra، Hadoop، Spark و غیره با الگوریتم های موثر برای مدیریت ساختار، داده های بدون ساختار و نیمه ساختار یافته ارائه می شوند.

الف. ماشین بردار پشتیبانی (SVM)

ماشین بردار پشتیبانی (SVM) یک الگوریتم یادگیری ماشینی نظارت شده است که می تواند برای چالش های طبقه بندی و رگرسیون استفاده شود. با این حال، بیشتر در مسائل طبقه بندی استفاده می شود. در این الگوریتم، ما هر آیتم داده را به عنوان یک نقطه در یک فضای بعدی رسم می کنیم (که در آن n تعداد ویژگی هایی است که شما دارید (با مقدار هر ویژگی مقدار یک مختصات خاص. شکل 2، فوق صفحه را نشان می دهد که این دو کلاس را متمایز می کند.



شکل 2: هایپر صفحه

ب. مورد خطی قابل تفکیک

در حالت قابل جداسازی خطی، یک یا چند ابر صفحه وجود دارد که ممکن است دو کلاس ارائه شده توسط داده های آموزشی را با 100٪ از هم جدا کنند

ج. مورد غیر خطی قابل تفکیک

در حالت غیر خطی قابل تفکیک، نمی توان یک ابر صفحه خطی پیدا کرد که تمام مثال های مثبت و منفی را از هم جدا کند. برای حل این مورد، تکنیک به حداکثر رساندن حاشیه ممکن است با اجازه دادن به برخی از نقاط داده در سمت اشتباه حاشیه، یعنی اجازه دادن درجه ای از خطا در جداسازی، تسهیل شود. متغیرهای ϵ Slack برای نشان دادن درجه خطا برای هر نقطه داده ورودی معرفی می شوند. ماشین بردار پشتیبانی برای طبقه بندی ویژگی های کارت های اعتباری توسط

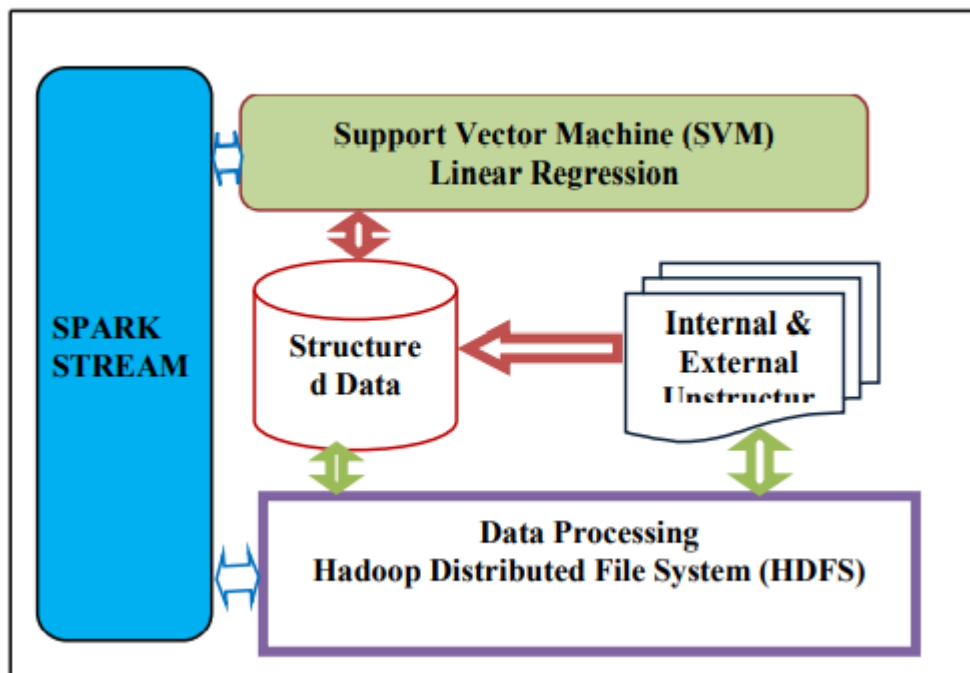
هر مشتری استفاده می شود SVM. باینری مشکل جداسازی دو کلاس را که با n مثال از هر کدام از ویژگی های m نشان داده شده است، حل می کند. مشکل زیر را در نظر بگیرید:

$$\{(x_1, y_1), \dots, (x_n, y_n)\}, x_i \in \mathcal{R}^m, y_i \in \{-1, +1\}$$

جایی که x_i در حال یادگیری مثال ها و y_i کلاس های مربوطه خود هستند. هدف روش SVM یافتن تابع خطی (f معادله 1) به نام hyperplane است که امکان جداسازی دو کلاس را فراهم می کند:

$$F(x) = (x \cdot w) + b$$

جایی که x نمونه ای برای طبقه بندی است، w یک بردار و b یک بایاس است. بنابراین ما باید وسیع ترین حاشیه را بین دو کلاس پیدا کنیم که معادل به حداقل رساندن است



شکل 3: معماری SVM-S برای تشخیص تقلب در داده های بزرگ

شکل 3، معماری SVM-S را برای چارچوب تشخیص تقلب نشان می دهد. برنامه spark همراه با پایگاه داده نیست. آن را همیشه در پایگاه داده خارجی برای عملیات آن سزاوار است HDFS. که در مدیریت کلان داده محبوب است، به عنوان پایگاه داده استفاده می شود. با تکنیک های یادگیری ماشین، ما SVM را برای آموزش مجموعه داده برای پیش بینی ها مستقر می کنیم.

2. Back Propagation Network: استفاده از BPN برای آموزش داده ها نیاز به تنظیم برخی پارامترها دارد. از جمله مهمترین پارامترها می توان به اعداد لایه پنهان N_i ، گره پنهان N_j و دوره های آموزشی N_k ، نرخ یادگیری R_j و نرخ تکانه R_m اشاره کرد. علاوه بر این، تنظیم مقادیر پارامتر به جای علم به عنوان یک هنر باقی می ماند. مشکلات پیچیده را می توان با افزایش لایه های پنهان به طور تدریجی بهتر مدل سازی کرد، اما این بهبود عموماً با هزینه های مرتبط از نظر زمان آموزش و تطبیق داده ها همراه است.

3. تست و نتیجه

آ. داده های مورد استفاده: به دلیل ماهیت محرمانه داده ها، دستیابی به داده های واقعی که رفتار مشتری بانک را توصیف می کند بسیار دشوار بود. با این حال، پایگاه های داده استاندارد وجود دارد که در ادبیات برای آزمایش روش های تشخیص تقلب استفاده می شود. برای آزمایش پیشنهاد هیبریداسیون، از پایگاه های داده آلمانی و استرالیایی کارت های اعتباری استفاده کردیم.

مجموعه داده برای پیاده سازی به حافظه برای پردازش فراخوانی شد. مجموعه داده پسوند فایل CSV بود. این داده های ایجاد شده به عنوان داده های تاریخی برای آموزش استفاده شد، یعنی با استفاده از الگوریتم نظارت شده. نام فایل برای مجموعه داده اول "creditcard.csv" و برای مجموعه داده دوم "creditcard.csv1" نامیده می شد. انتزاع اولیه Spark مجموعه ای توزیع شده از آیتم ها به نام مجموعه داده های توزیع شده انعطاف پذیر (RDD) است. در آن صورت، فایل

ها برای عمل و تبدیل به RDD تبدیل شدند. مجموعه داده با استفاده از SVM طبقه‌بندی شد و با رگرسیون خطی و رگرسیون منطقی برای تشخیص ناهنجاری با استفاده از ویژگی‌های کارت اعتباری آموزش دید. ترکیب این سه روش در مدل ما باعث افزایش دقت تشخیص می‌شود. جدول 1 پارامترهای تنظیم شده برای مجموعه داده را نشان می‌دهد.

Data Set	Classifier Model	Parameters	Definition	Source
1 st Data Set	SVM	C=10 RBF $\epsilon = 0.1$		Djeffal et al. 2014
2 nd Data Set	BPN	C=10 RBF $\epsilon = 0.1$		Soltani et al. 2014

جدول 1: پارامترهای ورودی برای SVM

مدل SVM-S ساخته شده برای دو مجموعه داده از یک تابع هسته استفاده می‌کند. دو پارامتر مرتبط با هسته RBF وجود دارد C: و ϵ . پارامتر تنظیم (C) مبادله بین حداکثر کردن حاشیه و به حداقل رساندن مدت خطای آموزشی را کنترل می‌کند. افزایش مقدار باعث بهبود دقت طبقه‌بندی (کاهش خطای رگرسیون) برای آموزش داده می‌شود. برای یک SVM-S، مقدار ϵ در تابع ضرر غیر حساس به ϵ نیز باید انتخاب شود. ϵ روی صاف بودن پاسخ SVM تأثیر می‌گذارد و بر تعداد بردارهای پشتیبانی تأثیر می‌گذارد، بنابراین هم پیچیدگی و هم قابلیت تعمیم شبکه به مقدار آن بستگی دارد.

Samples		# of Records	
		Training Set	Test Set
Set -1F-To-1N (SVM-S)	Normal	100	86
	Fraud	25	20
Set -1F-To-1N (BPN)	Normal	150	125
	Fraud	25	12

جدول 2: اندازه مجموعه آموزش و آزمون برای نمونه ها

دو مجموعه داده مختلف برای آزمایش هر مدل الگوریتم، با شرایط مختلف استفاده شد. اولین مجموعه داده یک تراکنش عادی برای هر تراکنش متقلبانه دارد. در حالی که مجموعه دوم، دارای چهار نرمال برای هر یک تقلبی است. برای هر داده نمونه، 70 درصد از داده ها، 70 درصد تراکنش عادی و 70 درصد جعلی، به عنوان مجموعه آموزشی برای مدل در نظر گرفته می شود. در حالی که 30 درصد از داده ها به عنوان مجموعه آموزشی برای ارزیابی عملکرد مدل مستقر شده در نظر گرفته می شود. اندازه مجموعه آموزش و آزمون در جدول 2 در بالا آورده شده است.

ب نتیجه: عملکرد معماری SVM-S در جدول 3 ارائه شده است، ستون سمت چپ روشی که برای ساخت مدل های معماری استفاده شده است، ستونی با نام "train" نشان دهنده دقت پیش بینی معماری پیشنهادی بر روی مجموعه داده های آموزشی است. در نمونه های داده شده، ستون هایی با برچسب "تست" دقت پیش بینی مدل معماری را در مجموعه داده های آزمایشی نمونه های داده شده نشان می دهد. ستونی که به عنوان "Build Time" نامگذاری شده است، زمان سپری شده برای ساخت مدل معماری داده شده را در نمونه داده شده نشان می دهد و ستون های "Frauds" تعداد تراکنش های تقلبی را در مجموعه داده های اختصاص داده شده به عنوان تقلب (TruePositive) توسط معماری بر روی نمونه های داده شده نشان می دهد. .

مقایسه نتایج پیش‌بینی شده بین SVM-S و BPN، جدول 3، مقایسه نتایج پیش‌بینی شده بین SVM-S و BPN را از آزمایش نشان می‌دهد.

Model/ Data Set	Set-1F-To-1N				Set-1F-To-4N			
	Tra in	Tes t	Bu ild Ti me	Fra ud	Trai n	Tes t	Bu ild Ti me	Fra ud
D1: SVM	98.7 8%	84.3 7%	<2 0m	20	96.3 4%	82.5 4%	<2 0	20
D2: BPN	99.8 6%	85.5 8%	<2 5m	13	97.3 46%	84.6 7%	<2 5	10

جدول 3: دقت عملکرد نسبت به مجموعه های آموزشی و آزمایشی

از جدول 3، مشخص است که معماری SVM-S با عملکرد SVM بسیار کارآمد است. با این حال، دقت میزان تخصیص واقعی را بدون توجه به اینکه یک تقلب واقعی است یا انتساب عادی واقعی، نشان می‌دهد. با این وجود، تعداد تراکنش‌های متقلبانه که توسط مدل به عنوان جعلی در نظر گرفته شده بود، برای مدل دقیق بود. این واقعیت که نتیجه مورد انتظار تعداد تراکنش‌های تقلبی تزریق شده با هر مجموعه داده، نتیجه SVM-S برابر است با کار در مقایسه با نتیجه BPN. این به وضوح نشان می‌دهد که SVM-S از BPN قابل اعتمادتر و دقیق‌تر است. با پیچیدگی زمانی، SVM-S از زمان کمی برای پیش‌بینی ناهنجاری‌ها در مقایسه با الگوریتم‌های دیگر مانند BPN استفاده می‌کند.

۷. نتیجه

در این زمینه دو مورد کلاهبرداری در بانک‌ها را بررسی کردیم: کلاهبرداری از کارت اعتباری و پولشویی. عملکرد سیستم پیشنهادی بر روی معیارهای دفتر کل، داده‌های پرداختی، که شبیه به

پایگاه داده بانک ایجاد شده است، آزمایش شد. دقت به دست آمده برای روش SVM تک کلاس حدود 80 درصد بود که نشان دهنده پیشرفت قابل توجهی در مقایسه با مرجع آثار مشابه است. برای این روش، بهبود جزئی در پایگاه های داده امتیازدهی اعتباری به دلیل دشواری دستیابی به پایگاه های داده واقعی بود. نتایج را می توان با مطالعه تأثیر پارامترهای مختلف مورد استفاده در معماری SVM-S بهبود بخشید.