

استفاده از روش نیمه‌نظارتی انتشار برچسب برای رتبه‌بندی اعتباری مشتریان بانک (مطالعه موردی بانک پاسارگاد)

مجید ابتیاع

دانشجوی کارشناسی ارشد دانشکده علوم پایه، گروه ریاضی، دانشگاه آیت‌الله‌العظمی بروجردی
پست الکترونیکی: majid.ebtia@gmail.com

سیدمحمدحسینی*

استادیار دانشکده علوم پایه، گروه ریاضی، دانشگاه آیت‌الله‌العظمی بروجردی
پست الکترونیکی: sm.hoseini@abru.ac.ir

چکیده

کارشناسی شده ندارند، استفاده از روش‌های نیمه‌نظارتی توصیه می‌شود. در روش‌های نیمه‌نظارتی برخلاف روش‌های نظارتی، لزومی به برچسب‌دار بودن تمام داده‌ها نیست و می‌توانند به وسیله مجموعه‌های داده‌ای که بخش زیادی از آن بدون برچسب هستند، مدل‌های قابل اعتمادی ایجاد کنند. روش انتشار برچسب یکی از روش‌های نیمه‌نظارتی موفق و پرکاربرد در این زمینه می‌باشد. در این روش، ویژگی رده داده‌های برچسب‌دار با یک روش تکراری به داده‌های فاقد برچسب منتشر می‌شود. در این پژوهش، ابتدا با کمک روش درخت تصمیم، ویژگی‌های تاثیرگذار بر وضعیت اعتباری مشتریان شناسایی و با کمک روش انتشار برچسب، مدل طبقه‌بندی مشتریان بانک ساخته می‌شود. به منظور نمایش کارایی مدل، از مجموعه داده‌های بانک پاسارگاد استفاده شده است. مدل نهایی بر روی داده‌های بانک پاسارگاد توانست به ۹۹/۷۸ درصد نرخ تشخیص مشتریان بدحساب و دقت کلی ۹۸/۷۷ درصد برسد. همچنین به منظور مقایسه، داده‌های یک

امروزه، بانک‌ها و موسسات مالی مجموعه‌های حجیم از داده‌های مختلف و متنوعی مرتبط با مشتریان خود جمع‌آوری و نگهداری می‌کنند. رتبه‌بندی اعتباری با هدف شناسایی برچسب مشتریان اعتباری اعم از این‌که خوش‌حساب یا بدحساب هستند، سبب کاهش معوقات بانکی و ارائه بهینه و درست تسهیلات به مشتریان می‌شود و به‌عنوان یکی از مهم‌ترین روش‌های تحقیقی و عملیاتی مورد استفاده در بانک‌داری و مهندسی مالی شناخته می‌شود. در این راستا، روش‌های گوناگونی که تحت عنوان روش‌های نظارتی شناخته می‌شوند ابداع و استفاده شده است. در این روش‌ها، مدلی به کمک داده‌های برچسب‌دار موجود، ایجاد و برای پیش‌بینی برچسب داده‌های جدید استفاده می‌شود. بنابراین پیش‌فرض چنین روش‌هایی، وجود پایگاه داده‌ای است که تمام اشیاء داده‌ای آن برچسب داشته باشند. اما از آنجا که تعداد زیادی از داده‌های موجود در بانک‌ها فاقد برچسب هستند، یا برچسب دقیق و

* نویسنده مسئول

بانک خارجی (آلمانی) نیز مورد بررسی قرار گرفت. بر اساس منابع مورد مطالعه در این تحقیق، بهترین دقت در بین مدل‌های نظارتی توسط روش درخت تصمیم مبتنی بر خوشه‌بندی و الگوریتم ژنتیک با دقت ۸۵/۳۷ درصد گزارش شده است. همچنین در بین مدل‌های نیمه‌نظارتی، روش نیمه‌نظارتی گروهی مبتنی بر پرسپترون چندلایه با دقت ۷۵/۴ درصد گزارش شده است، در حالی که روش نیمه‌نظارتی پیشنهادی در این مقاله، دقت ۷۶/۸۵ درصدی را کسب کرد. یافته‌های پژوهش حاکی از عملکرد مطلوب روش نیمه‌نظارتی پیشنهادی است.

واژه‌های کلیدی: روش نیمه‌نظارتی، رتبه‌بندی اعتباری،

بانک، انتشار برچسب

۱- مقدمه

رتبه‌بندی اعتباری شرکتی یک تحلیل از ریسک‌های اعتباری در یک شرکت است که نقش حیاتی در مدیریت ریسک مالی ایفا می‌کند. امتیازات اعتباری در همه جا وجود دارند و ابزاری برای تامین‌کنندگان وام و تنظیم‌کنندگان وام هستند. امتیاز اعتباری به طور گسترده ابزاری در صنعت مالی است که بارها اثر بخشی خود را در به حداقل رساندن ریسک اثبات کرده است. وظیفه اصلی یک امتیاز اعتباری، کمک به ارائه‌دهندگان وام در برآورد ریسک در حین تطبیق‌پذیری و هزینه برای محصولاتشان است. بانک‌ها و نهادهای مالی بخش اصلی نظام مالی هر کشوری را تشکیل می‌دهند و نقش مهمی در بخش‌های مختلف اقتصادی بر عهده دارند. نظام بانکی با فعالیتهای متنوع و گوناگون از قبیل تجهیز منابع، تامین نقدینگی، ابزارهای پرداخت، اعطای تسهیلات و سایر موارد بر عملکرد اقتصاد کشور تاثیر می‌گذارد. مهم‌ترین فعالیت بانک‌ها، جمع‌آوری منابع مالی و سرمایه‌گذاری و تخصیص آن به بخش‌های مختلف اقتصادی است [۱]. در راستای ایفای این نقش، بانک‌ها همواره با ریسک‌های متفاوتی روبه‌رو هستند که یکی از مهم‌ترین آن‌ها ریسک اعتباری می‌باشد. ریسک

اعتباری را به طور خلاصه می‌توان احتمال قصور دریافت کنندگان تسهیلات بانکی نسبت به تعهداتشان طبق شرایط توافق شده نسبت به بانک تعریف کرد. بر این اساس لازم است بانک‌ها برای کنترل و کاهش ریسک اعتباری، قبل از پرداخت تسهیلات به مشتریان وضعیت اعتباری آن‌ها را ارزیابی کنند. به خاطر محدودیت منابع مالی و تسهیلات در اختیار بانک‌ها، ارائه درست و بهینه‌ی تسهیلات و اعتبارات مالی یکی از وظایف بسیار مهم بانک‌ها به شمار می‌رود و بانک‌ها تمایل به اعطای تسهیلات خود به مشتریانی دارند که از ریسک کمتری برخوردار هستند [۲]. ارزیابی درست مشتریان بانک‌ها و موسسات مالی و مدیریت ریسک آن‌ها مباحثی حیاتی هستند. بحران‌های به وجود آمده در نظام بانکی کشورها ناشی از عدم توجه به مدیریت ریسک اعتباری است و حجم عظیم معوقات بانکی، خود گویای نبود مدل‌های کارا و مناسب جهت اندازه‌گیری آن است. در نتیجه بانک‌ها و موسسات مالی برای ارزیابی ریسک اعتباری و کنترل آن نیازمند یک سیستم طبقه‌بندی مشتریان هستند تا با اختصاص رتبه شایسته به هر مشتری، تصمیمات مناسب جهت اعطای تسهیلات گرفته شود [۲]. از گذشته تا کنون پژوهش‌های متفاوتی برای ارائه مدلی مناسب و کارا جهت ارزیابی ریسک اعتباری مشتریان بانک‌ها ارائه شده است که می‌توان به روش‌های آماری، تحلیل ممیزی [۲]، برنامه‌ریزی ریاضی و سیستم‌های خبره اشاره کرد [۲۲]. در عصر حاضر با پیشرفت فناوری‌های نوین اطلاعاتی و کاربردهای وسیع آن در طبقه‌بندی اطلاعات، مدل‌هایی مبتنی بر هوش مصنوعی و یادگیری ماشین مورد توجه بسیاری از پژوهشگران قرار گرفته است. روش‌هایی مانند رگرسیون [۳]، شبکه‌های عصبی [۵]، الگوریتم ژنتیک [۱]، درخت تصمیم [۳]، تحلیل تشخیص خطی، ماشین بردار پشتیبان [۲] و بسیاری از روش‌های دیگر که زیرمجموعه داده‌کاوی، یادگیری ماشین و هوش مصنوعی قرار می‌گیرند، برای اعتبارسنجی مشتریان و تخمین ریسک اعتباری به کار گرفته شده است [۲۲]. روش‌های مذکور را

اصطلاحاً یادگیری نظارتی یا باناظر^۱ می‌گویند. روش‌های یاد شده، عموماً دقت بالایی دارند و روش‌هایی قوی محسوب می‌شوند، در عین حال معایبی دارند که از آن جمله می‌توان به پیچیدگی محاسباتی بالا و زمان آموزش طولانی و حساسیت به متوازن بودن برچسب‌ها اشاره کرد [۱،۲]. اما پیش‌فرض اصلی در روش‌های نظارتی، برچسب‌دار بودن داده‌ها است و بنابراین روش‌های نظارتی را نمی‌توان بر روی داده‌هایی که بخشی یا تمام آن‌ها فاقد برچسب هستند آموزش داد. از طرف دیگر، بخش زیادی از داده‌های موجود در پایگاه‌های داده‌ای و به‌ویژه در بانک‌ها فاقد برچسب هستند. به‌عنوان نمونه، با توجه به فقدان یک سیستم خودکار برچسب‌گذاری در مورد داده‌های مربوط به رتبه‌بندی اعتباری مشتریان، برچسب خوش‌حساب یا بدحساب اکثراً وجود ندارد یا توسط یک کارشناس برای تعداد معدودی از مشتریان انجام شده است. لازم به ذکر است که برچسب‌دار کردن داده‌ها یک فرآیند هزینه‌بر برای بانک‌ها است. علاوه بر این، برچسبی که توسط کارشناس بانک برای هر مشتری در نظر گرفته می‌شود می‌تواند دارای خطای انسانی باشد و گاهی ممکن است به‌طور سلیقه‌ای انجام شود. بنابراین استفاده از روش‌های نظارتی بر روی چنین داده‌هایی، یا به‌خاطر فقدان برچسب بخشی از داده‌های آموزشی ناممکن است، یا در صورت ساخت مدل حتی اگر دقت بالایی داشته باشد، به دلیل خطای موجود در داده‌ها، نتایج مطلوب و موثری ارائه نمی‌کنند.

بنابراین آنچه در بالا بیان شد، استفاده از روش‌های یادگیری نیمه‌نظارتی^۲ یکی از رویکردها و راه‌حل‌های مهم و مطرح در این گونه موارد است. در روش‌های نیمه‌نظارتی برخلاف روش‌های نظارتی، لزومی به برچسب‌دار بودن تمام داده‌ها نیست. این روش‌ها، می‌توانند به‌وسیلهٔ مجموعه‌های داده‌ای که بخش زیادی از آن بدون برچسب هستند، مدل‌های قابل اعتمادی ایجاد کنند. دیدگاه‌های نیمه‌نظارتی به چهار گروه شامل مدل‌های مولد، جداسازی

با چگالی کم، مبتنی بر گراف، تغییر نماینده و فرا ابتکاری تقسیم می‌شوند [۲۰]. الگوریتم‌های نیمه‌نظارتی علاوه بر مزیت مهم مذکور، روش‌هایی ساده هستند و فهم آسانی دارند و پایدار هستند [۲۰، ۱۸، ۹].

در این مقاله، مسئلهٔ رتبه‌بندی اعتباری مشتریان موسسات مالی مورد مطالعه قرار گرفته است. از آنجا که بخشی از داده‌های مورد بررسی در این مقاله فاقد برچسب بودند، از روش نیمه‌نظارتی انتشار برچسب به‌عنوان یکی از روش‌های نوین یادگیری ماشین بهره گرفته شده است. پژوهش حاضر به پنج بخش سازماندهی شده است. پس از مقدمه، در بخش دوم مروری جامع بر مبانی نظری و پیشینه پژوهش خواهیم داشت. سپس در بخش سوم به روش‌شناسی، ویژگی‌های داده‌ها و شناخت آن‌ها پرداخته می‌شود. در بخش چهارم برآورد مدل و نتایج گردآوری شده است. نهایتاً در بخش پایانی، نتیجه‌گیری و آرایه‌ی پیشنهادات آورده شده است.

۲- مبانی نظری و مروری بر پیشینه پژوهش

قدمت ارزیابی ریسک اعتباری به زمان ایجاد پول برمی‌گردد و زمانی است که افرادی به گروه‌ها و افراد مختلف پول قرض می‌دادند و توانایی مالی آن‌ها را در نظر می‌گرفتند. اما در آن زمان متغیرها و ویژگی‌های اثرگذار بر وضعیت مالی و اعتباری افراد محدود بود و با توجه به شناخت قبلی، وضعیت اعتباری افراد مشخص می‌شد. اندازه‌گیری و رتبه‌بندی ریسک اعتباری یک موضوع بسیار مهم و واقعی است و نیاز آن امروزه به شدت احساس می‌شود، کسب و کارهای کلان، بانک‌ها، موسسات مالی و ... به شدت نیازمند تجزیه و تحلیل صحیح و دقیق ریسک اعتباری خود می‌باشند. امروزه بخش عظیمی از داده‌ها را داده‌های بدون برچسب تشکیل می‌دهند، به همین علت روش‌های نیمه‌نظارتی به‌عنوان راهکاری کاربردی برای حل این مشکل توصیه می‌شود. روش‌های نیمه‌نظارتی علاوه بر کاهش هزینه‌های محاسبات و کاهش هزینه

1- Supervised learning

2- Semi-supervised learning

برچسب‌گذاری از تعمیم‌پذیری بهتری برخوردار هستند و پایدار می‌باشند [۲۰، ۱۸، ۹]. در این زمینه پژوهش‌های داخلی و خارجی مختلفی با روش‌های متفاوت انجام شده است که به برخی از آن‌ها اشاره می‌شود.

طلوعی و همکارانش (۱۳۸۸) با کمک مدل ترکیبی ماشین بردار پشتیبان مبتنی بر دو راهبرد به امتیازدهی اعتباری متقاضیان کارت‌های اعتباری می‌پردازد. آن‌ها با استفاده از تکنیک رتبه‌F، ویژگی‌های مهم و تاثیرگذار را شناسایی و برای پیدا کردن مقادیر پارامترهای بهینه ماشین بردار پشتیبان از جستجوی شبکه استفاده کردند [۱]. الگوریتم ژنتیک به‌عنوان یکی از تکنیک‌های تکاملی مبتنی بر جمعیت شناخته می‌شود و در بهینه‌سازی کاربرد فراوانی دارد. اقبالی و همکارانش (۱۳۹۶) به بررسی ارزیابی عملکرد توابع شایستگی الگوریتم ژنتیک در رتبه‌بندی مشتریان پرداختند و بین عملکرد این الگوریتم با روش‌هایی چون رگرسیون لجستیک و تحلیل پوششی داده‌ها مقایسه‌ای ارائه می‌کند [۲]. کشاورزحداد و آیتی‌گزار (۱۳۸۶) در پژوهشی، مدل‌های رگرسیون لجستیک و درخت تصمیم برای اعتبارسنجی مشتریان بانکی را مقایسه کردند. آن‌ها نشان دادند که خطای نوع اول و دوم برای مدل درخت تصمیم در مقایسه با رگرسیون لجستیک کمتر است [۳]. پویان‌فر و همکارانش (۱۳۹۲) با به‌کارگیری روش حداقل مربعات ماشین بردار پشتیبان مبتنی بر ژنتیک به تخمین رتبه اعتباری مشتریان بانک‌ها می‌پردازد. نتایج این پژوهش با مدل‌های آماری مانند لاجیت و رویکردهای بهینه‌سازی پارامترهای ماشین بردار پشتیبان مقایسه شد و نتایج این پژوهش حاکی از عملکرد مطلوب این مدل می‌باشد [۴]. مهرآرا و همکارانش (۱۳۸۸) با انتخاب نمونه ۴۰۰ تایی از مشتریان بانک پارسیان، به کمک روش‌های رگرسیون و شبکه‌های عصبی به ارزیابی ریسک اعتباری مشتریان پرداختند. آن‌ها نشان دادند که مدل مبتنی بر شبکه‌های عصبی نسبت به مدل رگرسیونی از دقت بالاتری برخوردار است [۵]. حاجی‌کرد و همکارانش در سال ۱۳۹۵، بر روی

داده‌های مشتریان بانک تجارت با استفاده از مدلی ترکیبی به بررسی ریسک اعتباری پرداختند. در این پژوهش ابتدا از الگوریتم ژنتیک به منظور بهینه‌سازی داده‌های ورودی استفاده شد. سپس مدل، بر اساس روش ماشین بردار پشتیبان به وجود آمد. بهینه‌سازی داده‌های ورودی، منجر به افزایش دقت بر روی داده‌های آزمایشی و تعمیم‌پذیری مدل می‌شود [۶]. شریعت‌پناهی و هاشمی‌برکادهی (۱۳۸۷) با ارائه مدل تحلیل ممیزی، اعتبارسنجی مشتریان بانک صنعت و معدن و همچنین پیش‌بینی نکول شرکت‌های دریافت‌کننده تسهیلات را انجام داد [۷].

معیار فاصله نقش بسیار مهم و کلیدی در بسیاری از الگوریتم‌های یادگیری ماشین دارد، به طوری که انتخاب تابع فاصله مناسب، تاثیر بسزایی بر عملکرد این الگوریتم‌ها دارد. با گذشت زمان و گسترش الگوریتم‌ها، فاصله‌ها جای خود را به هسته‌ها دادند و انتخاب هسته مناسب و پارامترهای آن به یک مسئله بسیار مهم تبدیل شد. زارع بیدکی و همکارانش (۱۳۹۶) الگوریتم یادگیری نیمه‌نظارتی هسته مرکب که از روش‌های یادگیری مبتنی بر فاصله است را تشکیل داد و از آن برای سنجش فاصله داده‌ها در خوشه‌بندی استفاده کرد [۸]. روش‌های یادگیری نیمه‌نظارتی که مبتنی بر گراف هستند بیشتر اوقات بر روی مسئله‌های تک برچسبی اجرا شده‌اند. کردآبادی و همکارانش (۱۳۹۸) در پژوهشی یک مدل نیمه‌نظارتی ترکیبی که مبتنی بر گراف و یادگیری چند برچسبی است تشکیل و نشان دادند این مدل نسبت به سایر مدل‌ها، به‌ویژه زمانی که تعداد نمونه‌های برچسب‌دار کم است عملکرد بهتری دارد [۹]. کریچن (۲۰۱۷) با الگوریتم بیز ساده به ارزیابی ریسک وام ۹۲۴ پرونده اعتبارات اعطایی که به شرکت‌های صنعتی تونسسی توسط یک بانک داده شده بود پرداخت و دقت مناسبی نزدیک به ۶۳ درصد به دست آورد. از نتایج حاصل می‌توان نقش وثیقه را در پیش‌بینی ریسک پیش‌فرض نام برد. در حقیقت این شاخص ظرفیت توضیحی برای پیش‌بینی ریسک اعتباری دارد [۱۱].

هانگ و همکاران (۲۰۰۷) برای امتیازدهی اعتباری مشتریان با یک روش استخراج داده بر اساس ماشین بردار پشتیبان، مدل‌هایی بر روی داده‌های دو بانک آلمانی و استرالیایی ارائه دادند. دقت مدل آن‌ها بر روی داده‌های بانک آلمانی نزدیک به ۷۶ درصد و بر روی داده‌های بانک استرالیایی نزدیک به ۸۵ درصد به دست آمد [۱۲]. در یادگیری ماشین باناظر، فرآیند به دست آوردن داده‌های برچسب‌دار پرهزینه و طولانی است. هان (۲۰۲۰) با کمک الگوریتم یادگیری نیمه‌نظارتی که از داده‌های برچسب‌دار و بدون برچسب بهره می‌گیرد، بر اساس رگرسیون لجستیک افزایشی تعمیم یافته به تشخیص ناهنجاری اعتبار شرکتی پرداخت. در این پژوهش نتایج حاصل شده از مطلوبیت خوبی برخوردار بود و همچنین ویژگی‌های مهم و تاثیرگذار شناسایی شدند [۱۳]. لیویریس (۲۰۱۸) در پژوهشی بر روی سه مجموعه داده بانکی استرالیایی، ژاپنی و آلمانی با کمک روش‌های گروهی نیمه‌نظارتی به ارزیابی ریسک اعتباری پرداخت که نتایج بسیار مناسبی حاصل شد [۱۴]. کیم (۲۰۱۹) با یک روش یادگیری گروهی نیمه‌نظارتی مبتنی بر ماشین بردار پشتیبان، عدم پرداخت بدهی در وام‌های بانکی را پیش‌بینی کرد [۱۵]. فنگ (۲۰۲۱) با کمک یک روش دوفازی و ترکیبی شرکت‌های مالی را رتبه‌بندی کرد. نتیجه حاصل شده از این روش ترکیبی، بهبود عملکرد و دقت است [۱۶].

۲-۱ یادگیری ماشین

یادگیری ماشین زیرمجموعه‌ای از هوش مصنوعی است که می‌توان با کمک آن سیستم یا برنامه‌ای طراحی کرد که الگوها را بر اساس تجربه یاد بگیرد [۱۸]. الگوریتم‌های یادگیری ماشین در یک دسته بندی به سه دسته یادگیری بانظارت، بدون نظارت و نیمه‌نظارتی تقسیم می‌شوند [۲۰]. در یادگیری بدون نظارت، تمامی داده‌ها بدون برچسب هستند و هدف از این روش خوشه‌بندی داده‌ها می‌باشد. منظور از خوشه‌بندی یعنی تقسیم مجموعه داده‌ای به چندین زیرمجموعه به گونه‌ای که داده‌های داخل هر زیرمجموعه

بیشترین شباهت به یکدیگر و کمترین شباهت به داده‌های سایر خوشه‌ها داشته باشد. در یادگیری بانظارت، تمام داده‌ها دارای برچسب هستند و هدف از این روش پیدا کردن مدل یا الگویی است که بتواند ویژگی برچسب هر داده را به طور مطلوب مشخص کند [۱۸]. یکی از روش‌های بسیار مهم در یادگیری نظارتی، درخت تصمیم می‌باشد که برای رگرسیون و دسته‌بندی به کار می‌رود. در این روش یک درخت تصمیم با تقسیمات دودویی تشکیل می‌گردد که در آن برای تقسیم‌بندی و انشعاب بر روی داده‌ها از شاخص جینی استفاده می‌شود.

در یادگیری نیمه‌نظارتی، هدف الگوریتم استفاده همزمان از داده‌های برچسب‌دار و بدون برچسب برای تقویت یادگیری است. الگوریتم‌های یادگیری نیمه‌نظارتی به طور کلی به چهار دسته روش‌های مولد، جداسازی کم چگالی، مبتنی بر گراف و فرا ابتکاری تقسیم می‌شوند [۲۰]. در این مقاله، یک الگوریتم تکراری مبتنی بر تکثیر یا انتشار برچسب که در آن همزمان از داده‌های برچسب‌دار و بدون برچسب استفاده می‌شود بررسی می‌گردد. داده‌های برچسب‌دار برای آموزش الگوریتم‌های نظارت‌شده ضروری هستند هرچند اغلب تعداد کمی از آن‌ها موجود است ولی داده‌های بدون برچسب فراوان است.

مجموعه داده کلی به صورت $X = \{x^1, \dots, x^{n+m}\}$ که در آن $x^i \in R^D$ را در نظر بگیرید. در این مجموعه، هر داده دارای D ویژگی یا متغیر و یک متغیر وابسته یا رده است اما تنها برچسب تعداد n داده از این مجموعه مشخص شده‌اند. بنابراین n داده برچسب‌دار وجود دارد و بقیه‌ی داده‌ها بدون برچسب هستند. بدون کاستن از کلیت، فرض کنید $(x^1, y_1), \dots, (x^n, y_n)$ داده‌های برچسب‌دار این مجموعه باشند. برچسب‌های رده تمام داده‌های برچسب‌دار را با بردار $Y_L = (y_1, \dots, y_n)$ نشان دهید. همچنین فرض کنید C تعداد رده‌های مختلف موجود در داده‌های برچسب‌دار را نشان می‌دهد و از هر رده در مجموعه داده‌ای برچسب‌دار تعدادی موجود است. در نظر

بگیرید $(x^{n+1}, y_{n+1}), \dots, (x^{n+m}, y_{n+m})$ داده‌های بدون برچسب باشند. به عبارت دیگر y_{n+i} به ازای $i = 1, 2, \dots, m$ نامشخص است. از بردار $Y_U = (y_{n+1}, \dots, y_{n+m})$ برای برچسب داده‌های با برچسب نامشخص استفاده کنید. در این مجموعه داده‌ای، به دنبال آن هستیم نقاطی از داده‌ها که نزدیک به هم هستند برچسب‌های مشابه داشته باشند. روش انتشار برچسب به شیوه زیر برچسب داده‌های بدون برچسب را تعیین می‌کند. ابتدا یک گراف کامل وزن دار که در آن راس‌ها همه داده‌های مجموعه شامل داده‌های با برچسب و بدون برچسب هستند ایجاد می‌شود. یال‌های این گراف با ω_{ij} که نشان‌دهنده وزن میان دو راس i, j است وزن دار می‌شوند. در این مقاله از فاصله اقلیدسی برای وزن دار کردن یال‌ها استفاده شده است. نکته مهمی که باید به آن اشاره کرد این است که به جای فاصله اقلیدسی می‌توان هر معیار فاصله دیگری قرار داد. در اینجا با توجه به راحتی و متداول بودن فاصله اقلیدسی از آن استفاده شده است. وزن یال بین دو راس به صورت زیر مقداردهی می‌گردد:

$$\omega_{ij} = \exp\left(-\frac{d_{ij}^2}{\sigma^2}\right) = \exp\left(-\frac{\sum_{d=1}^D (x_d^i - x_d^j)^2}{\sigma^2}\right) \quad (1)$$

که در آن وزن‌ها به کمک پارامتر σ کنترل می‌شوند اما لازم به ذکر است که می‌توان برای هر متغیر یک پارامتر σ_d مخصوص آن به کار برد. در این صورت داریم

$$\bar{\omega}_{ij} = \exp(-\bar{d}_{ij}^2) = \exp\left(-\sum_{d=1}^D \frac{(x_d^i - x_d^j)^2}{\sigma_d^2}\right) \quad (2)$$

در این مقاله، از یک پارامتر σ برای همه ویژگی‌ها استفاده شده است. این شیوه وزن دار کردن یال‌ها به گونه‌ای است که هر چه گره‌ها از لحاظ فاصله اقلیدسی نزدیک‌تر باشند، وزن ω_{ij} بزرگتر است و از منظر هدف الگوریتم انتشار برچسب، دو راسی که وزن بیشتری در نتیجه فاصله اقلیدسی کمتری دارند شباهت بیشتری وجود دارد و اگر یکی از آن‌ها برچسب‌دار و دیگری بدون برچسب باشد همان برچسب به راس بدون برچسب نسبت

داده خواهد شد. در واقع برچسب راس‌های برچسب‌دار از طریق یال‌ها به راس‌های بدون برچسب منتقل و منتشر می‌شود و هرچه وزن یالی بزرگتر باشد انتشار برچسب از طریق آن آسان‌تر خواهد بود. برای انتشار برچسب، ماتریس انتقال T با ابعاد $(n+m) \times (n+m)$ به شکل زیر نیاز خواهد شد

$$T_{ij} = \frac{\omega_{ij}}{\sum_{k=1}^{n+m} \omega_{kj}} \quad (3)$$

که در آن T_{ij} درایه سطر i و ستون j است و احتمال پرش از گره j به i را محاسبه می‌کند. همچنین ماتریس Y با ابعاد $(n+m) \times C$ که ماتریس برچسب نامیده می‌شود و سطر i -ام آن نشان‌دهنده توزیع احتمال برچسب راس x^i در رده‌های مختلف می‌باشد. با توجه به این که برچسب راس‌های برچسب‌دار معلوم است لذا در ماتریس Y سطرها مربوط به آن‌ها در ستون رده‌شان مقدار ۱ و بقیه ستون‌ها صفر است. الگوریتم انتشار برچسب یک روش تکراری است که با مقداردهی اولیه ماتریس Y شروع و در انتهای الگوریتم، ماتریس Y با احتمال تعلق هر داده به هر یک از رده‌ها به عنوان خروجی ارائه می‌شود. ثابت می‌شود که مقدار اولیه سطرهایی از Y که متناظر با راس‌های بدون برچسب است اهمیتی ندارند و با هر مقدار بین صفر و یک، به جواب همگرا می‌شوند (چافیل و همکاران، ۲۰۰۶). مراحل الگوریتم انتشار برچسب به صورت زیر است:

- ۱- برچسب را انتشار بده: TY را محاسبه و جایگزین Y کن.
- ۲- سطرهای Y را نرمال کن.
- ۳- سطرهای مربوط به داده‌های برچسب‌دار را به حالت اول برگردان و به مرحله ۱ برو.

این سه مرحله می‌تواند تا زمان همگرایی Y تکرار شود یا بعد تعداد مشخصی تکرار، تمام شود. همچنین هر معیار دیگری برای توقف الگوریتم می‌تواند استفاده شود. در مرحله ۱، همه راس‌ها برچسب‌های خود را یک واحد منتشر می‌کنند. در مرحله ۲، نرمال شدن سطرهای Y برای تفسیر احتمالاتی برچسب‌ها انجام می‌شود. با توجه

جدول ۱: نتایج پژوهش‌های مشابه در مورد رتبه‌بندی اعتباری مشتریان بانک به کمک روش‌های یادگیری ماشین

منبع اطلاعات	دسته	روش	داده‌ها	مزایا	معایب	دقت
هانگ (۲۰۰۷)	نظارتی	ماشین بردار پشتیبان مبتنی بر جست‌وجوی شبکه و امتیاز F	بانک آلمانی	انتخاب ویژگی‌های مناسب و بهینه‌سازی پارامترها	افزایش هزینه محاسبات و کاهش سرعت الگوریتم	۷۷,۹۲
گتاشیا (۲۰۱۴)	نظارتی	جنگل تصادفی	بانک آلمانی	تفسیر پذیری بالا	پیچیدگی محاسباتی و حساسیت به داده‌ها	۷۷,۲
کشاورز حداد (۱۳۸۶)	نظارتی	درخت طبقه‌بندی و مدل لاجیت	بانک مسکن	بدون پیش‌فرض بر داده‌ها	حساس به تعداد داده‌ها	۸۵,۵
مهرآرا (۱۳۸۸)	نظارتی	شبکه عصبی	مشتریان حقوقی بانک پارسیان	دقت بالا، مقاوم در برابر نوفه و شناسایی روابط غیر خطی	پیچیدگی محاسبات	۸۶
طلوعی (۱۳۸۹)	نظارتی	ماشین بردار پشتیبان مبتنی بر جست‌وجوی شبکه و امتیاز F	مشتریان حقوقی یک بانک داخلی	انتخاب ویژگی‌های مناسب و بهینه‌سازی پارامترها	افزایش هزینه محاسبات و کاهش سرعت الگوریتم	۷۶,۷
پویان‌فر (۱۳۹۲)	نظارتی	ماشین بردار پشتیبان مبتنی بر ژنتیک	بانک آلمانی	افزایش دقت مدل	پیچیدگی محاسبات و افزایش زمان محاسبات	۸۴
محمدیان (۱۳۹۵)	نظارتی	ماشین بردار پشتیبان مبتنی بر ژنتیک	مشتریان حقوقی بانک تجارت	بهینه‌سازی پارامترها	کاهش سرعت الگوریتم	۶۹
اقبال (۱۳۹۶)	نظارتی	مقایسه توابع شایستگی الگوریتم ژنتیک	مشتریان حقوقی یک بانک داخلی	افزایش دقت مدل	هزینه محاسباتی بالایی دارد	۴۴
لیو بریس (۲۰۱۸)	نیمه‌نظارتی	یادگیری نیمه‌نظارتی ترکیبی	بانک آلمانی	افزایش دقت مدل و ساخت دسته‌بند قوی‌تر	افزایش پیچیدگی مدل	۷۵,۴

در جدول ۱، مقایسه‌ای بین پژوهش‌های مشابه گزارش می‌شود.

۳- مدل پژوهش، ویژگی‌ها و شناخت آن‌ها

اعتبارسنجی و ارزیابی ریسک اعتباری شامل طبقه‌بندی متقاضیان تسهیلات اعتباری بر اساس ویژگی‌ها و شرایط متقاضی مثل وضعیت کاری، میزان حساب بانکی، وضعیت مسکن تسهیلات و ... بررسی می‌شود. در این بخش به شرح مدل به کار برده شده در این مقاله، برای طبقه‌بندی جامعه آماری متقاضیان تسهیلات اعتباری و تعیین وضعیت اعتباری آن‌ها پرداخته می‌شود. مدلی که بر اساس داده‌ها ساخته می‌شود را اصطلاحاً یادگیرنده، طبقه‌بند یا دسته‌بند می‌نامند. جامعه آماری در نظر گرفته شده در این پژوهش مشتریانی هستند که از بانک مورد نظر تسهیلات دریافت کرده‌اند و آن‌ها را به بانک بازگشت داده یا نداده‌اند.

به این‌که در مرحله ۱ و ۲ ممکن است توزیع احتمالاتی راس‌های برچسب‌دار به خاطر تاثیر گرفتن از راس‌های مجاور تغییر کرده باشند در مرحله ۲، سطرهای مربوطه به داده‌های برچسب‌دار به حالت ثابت اولیه برگردانده می‌شوند. بنابراین به جای این‌که اجازه داده شود که راس‌های برچسب‌دار به مرور در راس‌های دیگر محو شوند، با دست‌زدن به توزیع برچسب‌های آن‌ها، دوباره آن‌ها اصلاح می‌شوند.

در این مقاله به منظور تنظیم پارامتر σ به صورت زیر عمل می‌شود. ابتدا با کمک الگوریتم کروسکال (کروسکال، ۱۹۵۶)، درخت پوشای مینیمم بر گراف کامل پیدا کنید سپس کوتاه‌ترین یال از این درخت را به گونه‌ای پیدا کنید که بین دو راس با برچسب مختلف است. چنانچه وزن این یال a باشد مقدار $\sigma = a/3$ لحاظ می‌شود. تقسیم به سه به خاطر قانون 3σ در توزیع نرمال است و به گونه‌ای است که وزن این یال نزدیک به میانگین توزیع نرمال استاندارد یعنی صفر باشد به این امید که انتشار محلی برچسب در رده‌های مختلف با همین پراکنندگی توزیع شود [۲۰، ۲۱].

جدول ۲: مقادیر ویژگی‌های داده‌های آلمانی

مقادیر ویژگی										نوع	عنوان ویژگی	نماد	
حساب بیش از DM 200		حساب کمتر از DM 200			حساب خالی			بدون حساب		اسمی	وضعیت حساب	G1	
عددی (برحسب ماه)										عددی	مدت اعتبار	G2	
حساب بحرانی یا اعتبارات در بانک‌های دیگر			تاخیر در بازپرداخت		بازپرداخت به موقع تمام تسهیلات		بازپرداخت به موقع در این بانک		بدون سابقه	اسمی	سابقه اعتبار	G3	
موارد دیگر	بازآموزی	تعمیر	رادیو	ماشین دست دوم	ماشین جدید	میلمان و تجهیزات	مسافرت	تحصیلات	لوازم خانه	کسب و کار	اسمی	هدف تسهیلات	G4
عددی										عددی	مقدار اعتبار	G5	
نامشخص یا بدون حساب پس‌انداز		بیش از DM 1000			DM 1000 تا DM 500		بین 100 تا 500 DM		کمتر از 100 DM	اسمی	وضعیت پس‌انداز	G6	
بیش از ۷ سال			بین ۴ تا ۷ سال		بین ۱ تا ۴ سال		کمتر از ۱ سال		بیکار	اسمی	سابقه کار	G7	
عددی										عددی	تعداد اقساط	G8	
زن متاهل یا مطلقه			زن مجرد		مرد متاهل		مرد مطلقه یا جدا شده		مرد مجرد	اسمی	وضعیت تاهل و جنسیت	G9	
ضامن			متقاضی مشترک				بدون ضامن			اسمی	سایر بدهکاران یا ضامنین	G10	
عددی										عددی	مدت زمان اقامت فعلی	G11	
اموال ناشناخته		املاک و ساختمان			بیمه زندگی		ماشین			اسمی	اموال و دارایی‌ها	G12	
عددی										عددی	سن	G13	
ندارد		فروشگاه‌ها				بانکی				اسمی	سایر تسهیلات	G14	
اجاره		مالک				رایگان				اسمی	مسکن	G15	
عددی										عددی	اعتبارات موجود در این بانک	G16	
مدیر/آزاد/کارمند/فوق تخصص/افسر		کارمند یا متخصص			بی تجربه و مقیم		بیکار/بی تجربه یا غیرمقیم			اسمی	شغل	G17	
عددی										عددی	تعداد عائله	G18	
بلی				خیر						اسمی	مالکیت تلفن	G19	
بلی				خیر						اسمی	کارگر خارجی	G20	
خوب				بد						اسمی	طبقه (رده)	G21	

نمونه آماری این پژوهش شامل یک بانک خارجی (آلمانی) که از پایگاه داده یادگیری ماشین UCI جمع آوری شده است. تعداد نمونه‌های داده‌های بانک آلمانی که در این پژوهش استفاده شده ۱۰۰۰ نمونه است که هر نمونه آن دارای ۲۰ ویژگی است. به طور کلی ویژگی‌ها به دو بخش کیفی و کمی تقسیم می‌شوند. ویژگی‌های این مجموعه، شامل ۲۰ ستون (۱۳ ویژگی اسمی و ۷ ویژگی عددی) است که در جدول ۲ گزارش شده است. در این مجموعه، هر مشتری یک ویژگی هدف یعنی وضعیت اعتباری دارد که برای هر کدام از نمونه‌ها با یکی از دو مقدار خوب (مشتریان خوش حساب) و بد (مشتریان بد حساب) مشخص شده است. در این پژوهش وضعیت اعتباری به صورت دودویی یعنی ۱ برای مشتری خوش حساب و ۰- برای مشتری بد حساب در نظر گرفته شده است.

همچنین یک مجموعه داده‌ای داخلی یعنی بانک پاسارگاد نیز بررسی شده است. در این مجموعه، صورت‌های مالی ۲۲۱۸ مشتری حقوقی وجود دارد که هر مشتری دارای ۱۴ متغیر (ویژگی) ۱ است و در جدول ۳ ارایه شده است. در هر دو مجموعه بانکی ۸۰ درصد از نمونه‌ها برچسب آن‌ها ناشناخته است.

مدل پژوهش

فرآیند مدل پیشنهادی در این پژوهش به دو مرحله تقسیم می‌شود. پس از جمع‌آوری داده‌های مورد نظر و مناسب به پیش‌پردازش آن‌ها پرداخته می‌شود.

● **گام اول:** این مرحله شامل پاک‌سازی، نرمال‌سازی و انتخاب ویژگی می‌باشد. هر دو مجموعه داده بانکی فاقد مقادیر مفقودی و نوفه می‌باشند. با کمک نرم‌افزار یادگیری ماشین پایتون به پاک‌سازی داده‌ها که شامل شناسایی داده‌های پرت است با روش‌های مختلف پرداخته شد و همچنین ویژگی‌های اسمی به عددی تبدیل شدند. برای نرمال‌سازی و مقیاس‌بندی مجموعه داده‌ها روش‌های

متفاوتی وجود دارد. به طور کلی می‌توان هر ویژگی را در محدوده [۱-] یا [۰,۱] مقیاس بندی کرد.

بعد از نرمال‌سازی، ویژگی‌های مهم و تاثیرگذار بر وضعیت اعتباری مشتریان شناسایی و انتخاب می‌شود. در این مقاله، این کار با کمک روش درخت رگرسیون و دسته‌بندی انجام می‌شود. در این روش داده‌ها به شاخه‌های دودویی تقسیم می‌شود و تقسیم داده‌ها و انشعاب روی آن‌ها با کمک معیار جینی صورت می‌گیرد این کار تا زمانی ادامه پیدا می‌کند که یکی از معیارهای توقف در درخت تصمیم مثل هرس کامل درخت یا رسیدن به گره‌های خالص تحقق پیدا کند. با کمک معیار جینی ویژگی‌های مهم و تاثیرگذار بر وضعیت اعتباری مشتریان شناسایی و انتخاب می‌شوند.

● **گام دوم:** داده‌های حاصل از پیش‌پردازش، به الگوریتم انتشار برچسب ارایه می‌شود. برای ساخت مدل ابتدا به روش اعتبارسنجی، پارامترهای مدل انتشار برچسب تنظیم و بهترین آن‌ها انتخاب می‌شود. پارامترهای بهینه الگوریتم به صورت زیر محاسبه می‌شوند. یکی از زیرمجموعه‌های داده‌های آموزشی انتخاب می‌شود. ابتدا یک شبکه از مقادیر ابرپارامترها انتخاب می‌شود و به ازای هر نقطه از این شبکه، ابرپارامترهای روش مقداردهی می‌شوند. سپس از روش اعتبارسنجی متقابل برای ارزیابی عملکرد ابرپارامتر در نقطه مذکور بهره برده خواهد شد.

در روش اعتبارسنجی متقابل، داده‌های آموزشی به k قسمت مساوی تقسیم می‌شوند. از زیرنمونه‌ها، یکی به عنوان داده‌های اعتبارسنجی برای آزمایش مدل حفظ می‌شود و بقیه زیر نمونه‌ها مدل را می‌سازند. بعد از ساخت مدل، زیرنمونه آزمایشی برای میزان اعتبار مدل ساخته شده به کار گرفته می‌شود. به این ترتیب به ازای هر نقطه شبکه، فرآیند اعتبارسنجی متقابل، k بار تکرار می‌شود و در هر بار معیار عملکرد آن نقطه سنجیده می‌شود. در انتها از میانگین معیارهای به دست آمده برای سنجش میزان اعتبار آن نقطه بهره برده می‌شود. در میان نقاط

جدول ۳: مقادیر ویژگی‌های داده‌های پاسارگاد

مقادیر ویژگی		نوع	عنوان ویژگی	نماد						
زن		اسمی	جنسیت	P_1						
مرد										
عددی		عددی	مقدار تسهیلات	P_2						
عددی		عددی	مبلغ وجه التزام	P_3						
سفته تضمینی	افرازنامه	تعهدنامه	چک تضمینی	اسناد عادی	سرمایه‌گذاری سپرده	سند ملکی - ترهینی - مسکونی / اداری	اوراق مشارکت	اسمی	نوع وثیقه	P_4
عددی		عددی	سن	P_5						
سه ماهه		شش ماه	نحوه بازپرداخت	P_6						
دکتری	فوق لیسانس	لیسانس	فوق دیپلم	دیپلم	اسمی	میزان تحصیلات	P_7			
عددی		عددی	میزان سابقه کار	P_8						
اجاره		مالک	وضعیت مسکن	P_9						
ندارد		دارد	سابقه چک برگشتی	P_{10}						
عددی		عددی	حدود درآمد ماهانه	P_{11}						
عددی		عددی	اموال و دارایی فعلی	P_{12}						
عددی		عددی	متوسط سپرده	P_{13}						
بد		خوب	طبقه (رده)	P_{14}						

این صورت معیارهای زیر به منظور عملکرد دسته‌بند به طور متداول استفاده می‌شود.

معیار اول که نرخ تشخیص نمونه‌های بد حساب نامیده می‌شود بسیار مهم و حائز اهمیت می‌باشد چرا که هزینه عدم تشخیص مشتریان بد حساب بسیار بیشتر از هزینه عدم تشخیص مشتریان خوش حساب می‌باشد. هزینه مورد نظر شامل از دست دادن تسهیلات به صورت اصل و فرع به همراه هزینه پیگیری مطالبات معوق می‌باشد. اختلاف این معیار از یک را می‌توان به عنوان ریسک اعتباری مدل ساخته شده در نظر گرفت.

$$TNR = \frac{TN}{TN + FP} \quad (4)$$

معیار دوم که نرخ تشخیص نمونه‌های خوش حساب نامیده می‌شود برای شناسایی مشتریان خوش حساب به کار می‌رود. پیش‌بینی نادرست و اشتباه مشتریان خوش

مختلف شبکه مذکور، نقطه‌ای که بهترین اعتبار را تولید کرده است به عنوان ابرپارامترهای مدل استفاده می‌شود. بعد از انتخاب بهترین پارامترها، الگوریتم انتشار برچسب روی مجموعه داده اعمال و مدل ساخته می‌شود.

معیارهای ارزیابی

برای ارزیابی و مقایسه نتایج حاصل شده و آزمون عملکرد روش‌های ارائه شده از معیارهای زیر بهره گرفته شده است. تعداد نمونه‌های رده بد حساب که به درستی بد حساب تشخیص داده شده‌اند را با TP، تعداد نمونه‌های رده خوش حساب است که به درستی خوش حساب تشخیص داده شده‌اند را با TN، تعداد نمونه‌های رده بد حساب که به اشتباه خوش حساب تشخیص داده شده‌اند را با FP و تعداد نمونه‌های رده خوش حساب که به اشتباه بد حساب تشخیص داده شده‌اند را با FN نشان دهید. در

حساب سبب کاهش حاشیه سود بانکی می‌شود. اختلاف این معیار از یک را می‌توان به‌عنوان ریسک تجاری مدل مورد استفاده در نظر گرفت.

$$TPR = \frac{TP}{TP + FN} \quad (5)$$

معیار سوم دقت کلی مدل است، که نشان‌دهنده توانایی مدل در دسته‌بندی کلی مشتریان چه خوش حساب و چه بدحساب می‌باشد.

$$Acc = \frac{TP + TN}{TP + TN + FP + FN} \quad (6)$$

معیار چهارم نشان‌دهنده پیش‌بینی‌های صحیح انجام شده برای مشتریان بدحساب را نشان می‌دهد.

$$PPV = \frac{TP}{TP + FP} \quad (7)$$

چون معیارهای دوم و چهارم عکس یکدیگر هستند و افزایش یکی باعث کاهش دیگری و بالعکس می‌شود پس از معیار دیگری استفاده می‌شود تا عملکرد کلی مدل ارزیابی شود این معیار برای داده‌های نامتوازن بسیار مناسب است. این معیار ترکیب دو معیار دوم و چهارم (میانگین هندسی) می‌باشد.

$$F_1 = \frac{2TP}{2TP + FP + FN} \quad (8)$$

معیار ششم نشان‌دهنده خطای نوع اول می‌باشد.

$$EI = \frac{FP}{FP + TN} \quad (9)$$

معیار هفتم نشان‌دهنده خطای نوع دوم می‌باشد.

$$EII = \frac{FN}{TP + FN} \quad (10)$$

معیار هشتم که بیان‌گر کیفیت کلاس بندی می‌باشد.

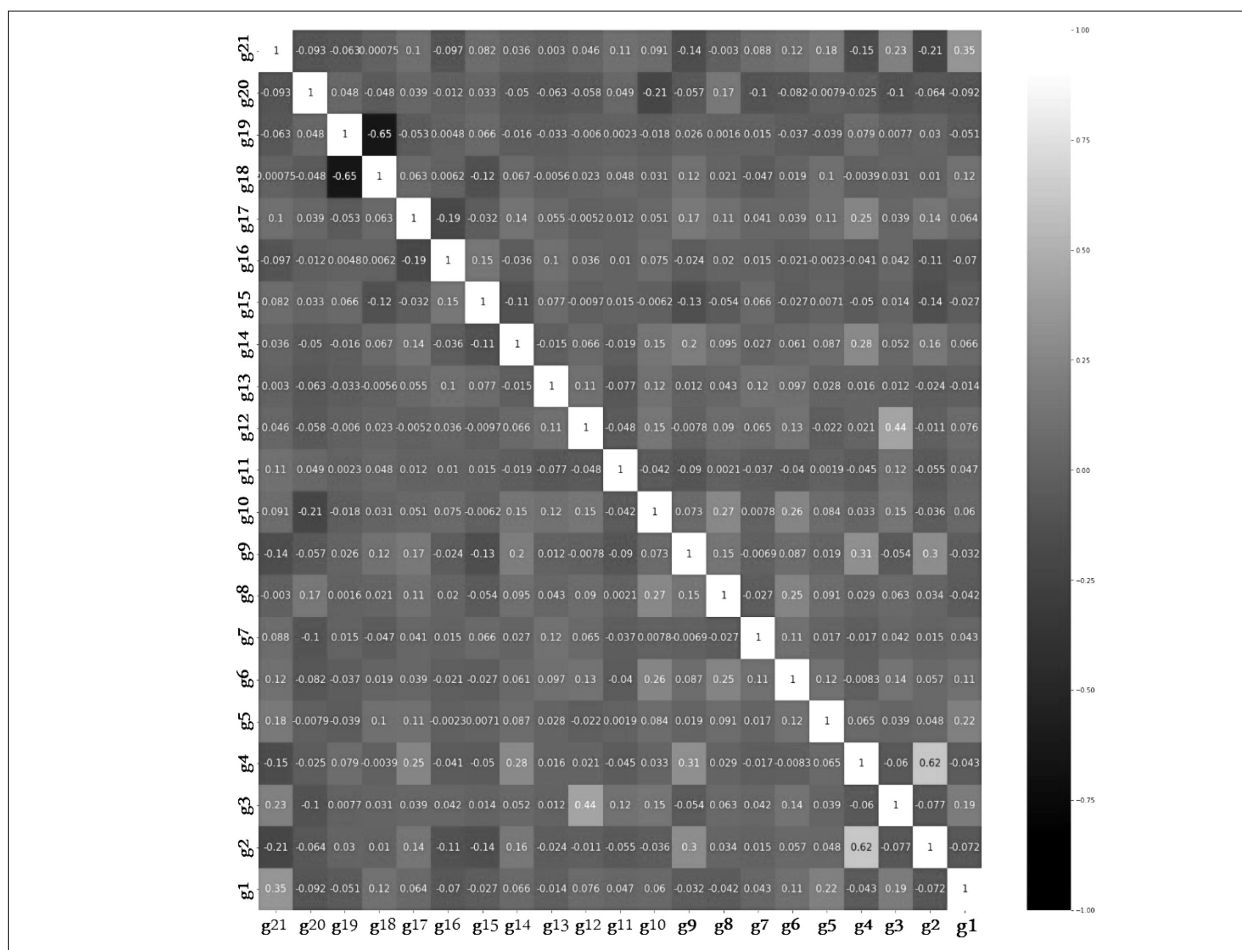
$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (11)$$

مساحت زیر سطح منحنی کارایی نشان‌دهنده قدرت یک مدل یا آزمون است که این مقدار عددی بین ۰ تا ۱ می‌باشد و هرچه به ۱ نزدیک‌تر باشد مدل کارایی بیشتری دارد. معیار نهم مساحت زیر سطح منحنی کارایی است که با AUC نشان داده می‌شود. در نمودار منحنی کارایی، محور عمودی نشان‌دهنده معیار ارزیابی اول و محور افقی نشان‌دهنده معیار ارزیابی پنجم است.

از آنجا که هر یک از معیارهای ارزیابی از اهمیت خاصی برخوردار بوده پس لازم است در ارزیابی عملکرد مدل از هر نه معیار در کنار یکدیگر استفاده شده و مورد بررسی قرار گیرند.

۴- برآورد مدل و نتایج پژوهش

در این قسمت نتایج حاصل از مدل ساخته شده بر روی مجموعه‌های داده‌ای بررسی می‌شود. مطابق بخش ۳-۱ (مدل پژوهش) در گام اول، داده‌ها پاک‌سازی و نرمال‌سازی می‌شوند. سپس عوامل تاثیرگذار بر وضعیت اعتباری مشتریان یک بار با معیار همبستگی و بار دیگر با روش درخت تصمیم شناسایی می‌شود. با روش همبستگی، ارتباط خطی هر ویژگی با وضعیت اعتباری مشتریان بررسی می‌شود. همچنین ارتباط بین ویژگی‌ها نیز مشخص می‌گردد. ماتریس همبستگی یک ماتریس متقارن است و مقادیر آن بین ۱- تا ۱ می‌باشد. در شکل ۱، ماتریس همبستگی متغیرهای مختلف داده‌های بانک آلمانی را مشاهده می‌نمایید. هر یک از سطر یا ستون‌های شکل ۱، نشان‌دهنده ویژگی‌های معرفی شده در جدول ۲ است. در شکل ۱ همبستگی بین ویژگی‌های مختلف بانک آلمانی با یکدیگر نمایش داده شده است. عنوان هر سطر یا ستون بر اساس نماد به‌کار رفته در جدول ۲ می‌باشد. سطر یا ستون‌های G_1 تا G_{20} ویژگی‌های ورودی، و سطر یا ستون G_{21} ویژگی وضعیت اعتباری یا خروجی را نمایش می‌دهد. از شکل ۱، مشخص است که ویژگی G_1 یعنی وضعیت حساب با مقدار ۳۵ درصد، ویژگی G_2 یعنی مدت اعتبار با مقدار ۲۱- درصد و ویژگی G_3 سابقه اعتبار با مقدار ۲۳- درصد بیشترین تاثیر را بر ویژگی رده (برچسب) یعنی وضعیت اعتباری مشتریان دارند. ویژگی‌های سن (G_{13})، مسکن (G_{15})، وضعیت تاهل (G_{19})، تعداد اقساط (G_8) و هدف تسهیلات (G_4) به ترتیب کمترین همبستگی با وضعیت اعتباری مشتریان بانک آلمانی دارند. همچنین ویژگی‌های G_2 و G_4 (مدت اعتبار و هدف



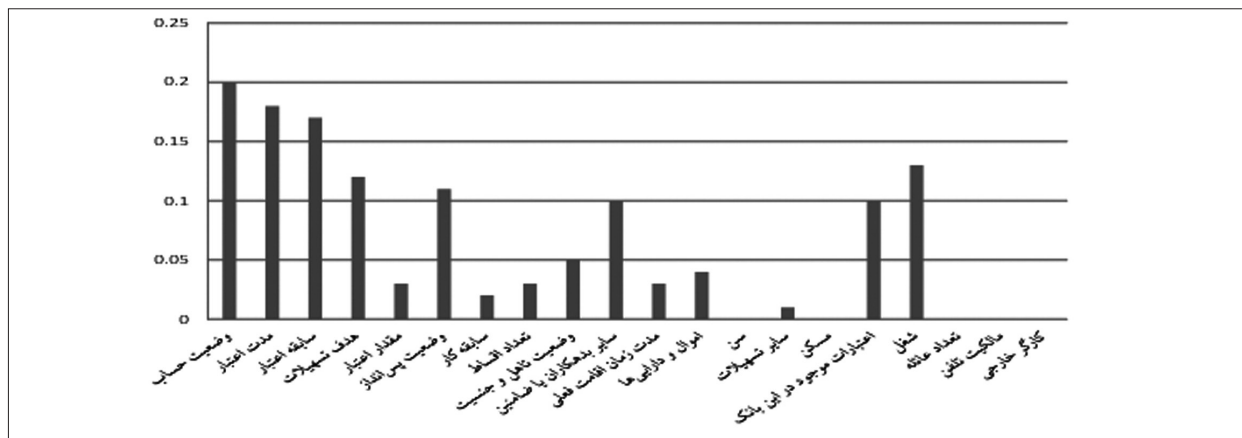
شکل ۱: همبستگی بین ویژگی‌های داده‌های آلمانی

تعداد عائله، مالکیت تلفن و کارگر خارجی تأثیری نزدیک به صفر بر وضعیت اعتباری مشتریان بانک آلمانی دارند و می‌توان آنها را حذف کرد. همچنین تحلیل همبستگی این ویژگی‌های با ویژگی رده نیز این مطلب را تایید می‌کند. در این پژوهش، ویژگی‌هایی که بر اساس الگوریتم درخت تصمیم اهمیتی کمتر از ۰/۰۰۱ دارند حذف شدند.

در شکل ۳، همبستگی بین ویژگی‌های بانک پاسارگاد با یکدیگر نمایش داده شده است. عنوان هر سطر یا ستون بر اساس نماد به کار رفته در جدول ۳ می‌باشد. سطر یا ستون‌های P_1 تا P_{13} ویژگی‌های ورودی و سطر یا ستون P_{14} ویژگی وضعیت اعتباری یا خروجی را نمایش می‌دهد. از شکل ۳، مشخص است که ویژگی P_6 (نحوه بازپرداخت) با مقدار ۸۵ درصد و ستون P_{12} (اموال و دارایی) با مقدار ۹۱- درصد بیشترین همبستگی با ویژگی رده یعنی

تسهیلات) با مقدار ۶۴ درصد بیشترین تأثیر را بر روی یکدیگر دارند که نشان‌دهنده همبستگی بالایی است و در صورت نیاز می‌توان یکی را به دلخواه حذف کرد.

الگوریتم درخت تصمیم به کمک معیار جینی، ارتباط بین هر ویژگی با ویژگی رده را نشان می‌دهد. معیار جینی یک معیار بر پایه آنتروپی است که برای شاخه‌زدن بر روی ویژگی‌ها استفاده می‌شود. هر چقدر مقدار این معیار کمتر باشد ویژگی از اهمیت بیشتری برخوردار است و بر اساس نسبت عکس آن، هر ویژگی یک امتیاز در بازه ۰ تا ۱ کسب می‌کند. شکل ۲، اهمیت ویژگی‌ها بر مبنای تأثیرگذاری بر ویژگی رده یعنی وضعیت اعتباری مشتریان را در مجموعه داده‌های بانک آلمانی نشان می‌دهد. محور افقی نام ویژگی‌ها و محور عمودی اهمیت هر ویژگی را نشان می‌دهد. با توجه به شکل ۲، ویژگی‌های سن، مسکن،



شکل ۲: اهمیت ویژگی‌های داده‌های آلمانی. منبع: محاسبات محقق

جدول ۴، گزارشی از معیارهای مختلف مذکور بر روی داده‌های بانک آلمانی و پاسارگاد ارایه می‌کند.

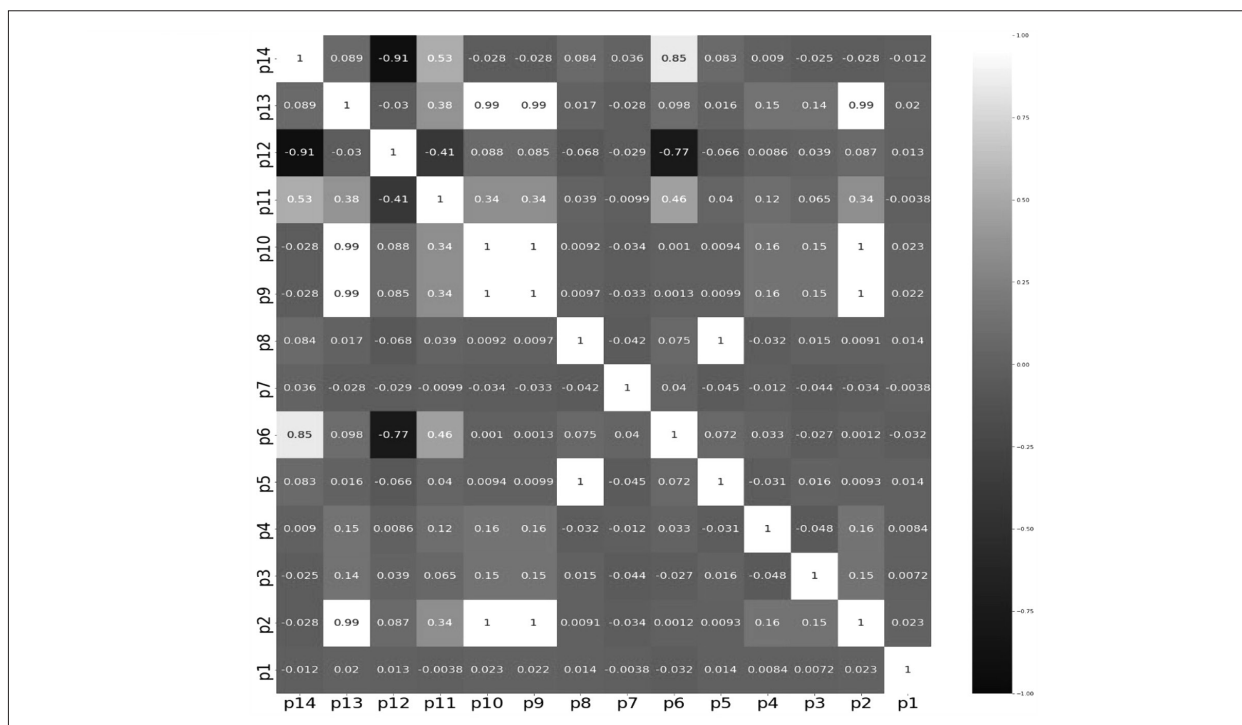
اطلاعات گزارش شده در جدول ۴، معیارهای ارزیابی حاصل از اعمال روش نیمه‌نظارتی پیشنهادی در این مقاله بر روی داده‌های بانک آلمانی و بانک پاسارگاد را نمایش می‌دهد. در مورد بانک پاسارگاد ابتدا روش پیشنهادی بر روی داده‌ها با پنج ویژگی مهم یعنی $P6$ (نحوه بازپرداخت)، $P10$ (سابقه چک برگشتی)، $P11$ (حدود درآمد ماهیانه)، $P12$ (اموال و دارایی فعلی) و $P13$ (متوسط سپرده) اعمال شد و نتایج بسیار خوبی به دست آمد که در جدول ۴ گزارش شده است. علاوه بر این، به منظور این که نشان داده شود پنج ویژگی مهم و منتخب با روش درخت تصمیم تاثیر بسزایی بر روی نتایج دارد روش پیشنهادی بر روی داده‌های بانک پاسارگاد بدون این پنج ویژگی مهم اعمال شد و نتایج آن در جدول ۴ ارایه شد. در ادامه به بررسی مفصل اطلاعات موجود در جدول ۴ پرداخته می‌شود.

بر اساس نتایج گزارش شده در جدول ۴ با روش انتشار برچسب، نرخ تشخیص نمونه‌های بدحساب یعنی معیار TNR در بین مشتریان بانک آلمانی ۷۱/۸۶ درصد است که نشان‌دهنده طبقه‌بندی قابل قبول مشتریان بدحساب می‌باشد. در مورد بانک پاسارگاد نیز مقدار ۹۹/۷۸ درصد برای این معیار به دست آمده است که بسیار مطلوب می‌باشد. همانگونه که قبلاً اشاره شد این معیار اهمیت بسیار بالایی برای مدیران بانک خواهد داشت زیرا هزینه

وضعیت اعتباری مشتریان را دارند. همچنین ویژگی $P2$ (مقدار تسهیلات) با ویژگی‌های $P9$ (سابقه کار)، $P10$ (چک برگشتی) و $P13$ (متوسط سپرده) همبستگی بالا و نزدیک به ۱۰۰ درصد دارد. این مطلب بدین معناست که می‌توان با داشتن ویژگی $P2$ به طور دقیق مقادیر ویژگی $P9$ و $P10$ را تعیین کرد بنابراین می‌توان سه ویژگی اخیر را از مجموعه داده‌ای حذف کرد.

شکل ۴، اهمیت ویژگی‌ها با استفاده از روش درخت تصمیم را در مجموعه داده‌ای بانک پاسارگاد نشان می‌دهد. محور افقی نام ویژگی‌ها و محور عمودی اهمیت هر ویژگی را نشان می‌دهد. با توجه به شکل ۴، غیر از ویژگی‌های نحوه بازپرداخت، سابقه چک برگشتی، حدود درآمد ماهیانه، اموال و دارایی فعلی و متوسط سپرده بقیه ویژگی‌ها اهمیتی نزدیک به صفر دارند. در این پژوهش در مورد داده‌های بانک پاسارگاد، ویژگی‌های کم اهمیت حذف شدند.

پس از انتخاب ویژگی‌های مهم، روش انتشار برچسب بر روی داده‌های جدید اعمال می‌شود. ابتدا برچسب ۲۰ درصد از داده‌های آموزشی حذف شدند. بنابراین در داده‌های آموزشی ۸۰ درصد برچسب‌دار و ۲۰ درصد بدون برچسب هستند. اکنون داده‌های آموزشی در اختیار الگوریتم انتشار برچسب قرار می‌گیرند و مدل ساخته می‌شود. پس از ساخت مدل نهایی، معیارهای ارزیابی مختلفی که در بخش قبل معرفی شدند محاسبه می‌شود.



شکل ۳: همبستگی بین ویژگی‌های بانک پاسارگاد

بانک آلمانی و پاسارگاد می‌باشد. همان‌طور که پیش از این بیان شد ریسک تجاری به کمک این معیار یا بر اساس معیار خطای نوع دوم (*EII*) قابل محاسبه است. ریسک تجاری روش پیشنهادی بر روی داده‌های بانک آلمانی ۱۸/۱۴ درصد و برای بانک پاسارگاد ۲/۲۲ درصد است که در مورد بانک پاسارگاد بسیار مطلوب است.

معیار ارزیابی *MCC* که به معیار ارزیابی بهره‌وری مدل نیز شناخته می‌شود هر اندازه که به صد نزدیک‌تر باشد نشان‌دهنده توان پیش‌بینی قابل قبول‌تر است. مقادیر درصدی این معیار با توجه به فرمول (۱۱) بین ۱۰۰- تا ۱۰۰+ است. بر این اساس میزان بهره‌وری مدل پیشنهادی بر روی داده‌های بانک آلمانی ۵۳/۹۸ و برای بانک پاسارگاد ۹۷/۵۶ است. بنابراین مدل پیشنهادی بر روی داده‌های بانک پاسارگاد توانایی بیشتری نسبت به داده‌های بانک آلمانی دارد. همچنین با حذف ویژگی‌های مهم در داده‌های بانک پاسارگاد، افت شدید این معیار به ۲۱/۲۵ نتیجه شده است. در مورد کارایی مدل‌ها از معیار *AUC* استفاده می‌شود. همان‌گونه که قبلاً بیان شد

عدم تشخیص مشتریان بدحساب بسیار بیشتر از هزینه عدم تشخیص مشتریان خوش‌حساب می‌باشد. اعطای تسهیلات به مشتری بدحساب منجر به افزایش معوقات بانک خواهد شد. از این رو با توجه به مفهوم این معیار یا بر اساس معیار خطای نوع اول (*EI*)، ریسک اعتباری مدل حاصل از روش انتشار برچسب بر داده‌های بانک آلمانی ۲۸/۱۴ درصد و بر داده‌های بانک پاسارگاد ۰/۲۲ درصد خواهد بود. بنابراین روش پیشنهادی قادر خواهد بود به نحو مطلوبی مشتریان بدحساب بانک آلمانی و با ریسک بسیار پایینی مشتریان بدحساب بانک پاسارگاد را تشخیص دهد. ریسک اعتباری برای داده‌های بانک پاسارگاد بدون در نظر گرفتن پنج ویژگی مهم افزایش چشمگیری می‌یابد (۶۱/۱۱ درصد) و نشان می‌دهد ویژگی‌های مهم به خوبی تشخیص داده شده‌اند.

از سوی دیگر نرخ تشخیص نمونه‌های خوش‌حساب *TPR* مشتریان بانک آلمانی ۸۱/۸۶ درصد و برای بانک پاسارگاد ۹۷/۷۸ درصد به دست آمد که نشان‌دهنده طبقه‌بندی بسیار مطلوب مشتریان خوش‌حساب در هر دو



شکل ۴: اهمیت ویژگی‌های داده‌های آلمانی. منبع: محاسبات محقق

آموزشی برچسب‌دار و بقیه بدون برچسب هستند. با مقایسه کلی نتایج جدول ۴ و نتایج پژوهش هانگ، می‌توان گفت که روش‌های نیمه‌نظارتی از عملکرد مناسبی برخوردار هستند. گتاشیا در سال ۲۰۱۴ به کمک روش نظارتی جنگل تصادفی برای رتبه‌بندی مشتریان بانک آلمانی، دقت کلی ۷۷/۲ درصد را ارایه کرد در حالی که روش نیمه‌نظارتی پیشنهادی در این مقاله، دقت کلی ۷۸ درصد را به دست آورده است. نتایج حاصل شده در مقایسه با جدول ۴ نشان از عملکرد مناسب روش‌های نیمه‌نظارتی در مقایسه با جنگل تصادفی دارد [۱۷]. لیویریس و همکارانش در سال ۲۰۱۸ با روش‌های نظارت شده و نیمه‌نظارت شده برای ارزیابی ریسک اعتباری داده‌های بانک آلمانی، دقت کلی ۷۵/۴ درصد را گزارش کرد [۱۴] در حالی که روش نیمه‌نظارتی پیشنهادی در این مقاله، ۲،۶ درصد بهبود در دقت کلی مدل را نشان می‌دهد.

در میان پژوهش‌های انجام شده بر روی مسئله رتبه‌بندی اعتباری مشتریان بانک، تحقیقاتی متنوعی وجود دارند که بر روی داده‌های بانک آلمانی نتایج گزارشی نکرده‌اند و به‌طور موردی داده‌های یک بانک مشخصی را تحلیل کرده‌اند که در ادامه به برخی از آن‌ها اشاره می‌شود. همچنان از معیار دقت کلی روش‌های ارایه شده برای مقایسه با نتایج حاصل از روش نیمه‌نظارتی پیشنهادی بر روی بانک پاسارگاد استفاده خواهد شد. اطلاعات خلاصه‌ای از آن‌ها در جدول ۱ ارایه شده است. محمدیان و همکارانش در سال ۱۳۹۵ با استفاده از روش

این معیار مساحت زیر منحنی کارایی مدل است و مقدار درصدی آن بین ۰ تا ۱۰۰ درصد خواهد بود هرچه این مقدار به ۱۰۰ درصد نزدیک‌تر باشد مدل کارایی بیشتری خواهد داشت. بر اساس جدول ۴، کارایی مدل نیمه‌نظارتی پیشنهادی بر روی داده‌های آلمانی ۸۱/۵۷ درصد و بر روی داده‌های بانک پاسارگاد ۹۹/۹۵ درصد می‌باشد که نشان از کارایی مطلوب روش پیشنهادی است. همچنین با حذف ویژگی‌های مهم از مجموعه داده‌های بانک پاسارگاد، کارایی مدل به ۶۳/۳۲ کاهش می‌یابد.

توانایی مدل نیمه‌نظارتی پیشنهادی در دقت دسته‌بندی کلی، یعنی *Acc*، برای مشتریان بانک آلمانی ۷۶/۸۵ درصد، برای بانک پاسارگاد ۹۸/۷۷ درصد و برای داده‌های بانک پاسارگاد بدون ویژگی‌های منتخب ۵۹/۶۶ درصد می‌باشد. این معیار که نرخ تشخیص درست همه برچسب‌ها را می‌سنجد به‌عنوان مبنای مقایسه روش پیشنهادی با سایر روش‌های موجود در پژوهش‌های دیگر مورد توجه قرار گرفت. هانگ و همکارانش در سال ۲۰۰۷ با استفاده از روش نظارتی ماشین بردار پشتیبان بر روی داده‌های بانک آلمانی دقت کلی ۷۷.۹۲ درصد را گزارش کرد [۱۲]. در سال ۱۳۹۲ پویان‌فر و همکارانش با انتخاب ویژگی به کمک الگوریتم ژنتیک و ساخت مدل نظارتی مبتنی بر ماشین بردار پشتیبان، دقت کلی نزدیک به ۸۴ درصد دست یافتند [۴]. لازم به یادآوری است که در روش‌های نظارتی، تمامی داده‌های آموزشی باید برچسب‌دار باشند در حالی که در روش‌های نیمه‌نظارتی بخشی از داده‌های

ماشین بردار پشتیبان برای ریسک اعتباری مشتریان بانک تجارت، دقت کلی ۶۹ درصد را گزارش کرد [۶]. طلوعی و همکارانش در سال ۱۳۸۹، چندین روش از جمله ماشین بردار پشتیبان، رگرسیون لجستیک، شبکه بیز و درخت تصمیم به رتبه‌بندی مشتریان بانکی به‌کار بردند که به‌عنوان نمونه روش ماشین بردار پشتیبان مبتنی بر شبکه و امتیاز F، دقت ۷۶/۷ را به‌دست آوردند [۱۰]. کشاورزحداد و همکارانش در سال ۱۳۸۶ بین مدل‌های لاجیت و درخت تصمیم برای اعتبارسنجی مشتریان بانکی مسکن درخت تصمیم با دقت بالاتر ۸۵/۵ درصد را به‌عنوان مدل بهتر نسبت به مدل لاجیت معرفی کرد [۳]. اقبالی و همکارانش در سال ۱۳۹۶ با استفاده از توابع شایستگی الگوریتم ژنتیک در رتبه‌بندی مشتریان حقیقی یک بانک داخلی به دقت کلی ۴۴ درصد دست یافتند [۲]. مهرآرا و همکارانش در سال ۱۳۸۸ با استفاده از شبکه‌های عصبی مصنوعی، رتبه‌بندی مشتریان حقوقی بانک پارسیان را بررسی کرد و دقت کلی نزدیک به ۸۶ درصد را گزارش کرد [۵]. همانطور که از جدول ۴ نتیجه می‌شود روش نیمه‌نظارتی پیشنهادی در این پژوهش بر روی داده‌های بانک پاسارگاد به مراتب دقت کلی بهتری نسبت به روش‌های مذکور دارد و حتی در صورتی که ویژگی‌های مهم از مجموعه داده‌ای این بانک داخلی حذف شود همچنان نسبت به روش‌های ارایه شده در مراجع [۲، ۶] عملکرد قابل توجهی دارد.

۵- نتیجه‌گیری و ارائه پیشنهادات

تخصیص امتیاز اعتباری دقیق به مشتریان بانک‌ها، همواره یک مسئله بسیار مهم برای محققان و بانک‌داران است، زیرا تنها یک درصد افزایش دقت می‌تواند جلوی خسارت‌هایی عظیم به بانک‌ها و موسسات مالی را بگیرد. با توسعه مداوم، گسترده و پویای صنعت اعتبارسنجی، هر روز این صنعت نقش مهم‌تری در اقتصاد ایفا می‌کند و اعتباردهندگان به‌منظور توسعه فرآیند مدیریت اعتباری از روش‌ها، ابزارها و ایده‌های جدید استفاده می‌کنند.

ساخت مدل‌های رتبه‌بندی اعتباری از یک پایگاه داده اعتباری را می‌توان به‌عنوان یک فرآیند داده‌کاوی انجام داد. در این پژوهش با استفاده از روش نوینی مبتنی بر یادگیری نیمه‌نظارتی، مشتریان بانک آلمانی و بانک پاسارگاد رتبه‌بندی شدند و مدلی کارا جهت ارزیابی وضعیت اعتباری مشتریان ارائه شد. نمایش عملکرد و کارایی مناسب و بسیار نزدیک دسته‌بندی‌های نیمه‌نظارتی به روش‌های نظارتی از نتایج اصلی این پژوهش می‌باشد. از آنجا که اعمال روش نیمه‌نظارتی پیشنهادی بر روی داده‌های بانک آلمانی در مقایسه با تحقیقات موجود بر اساس ۹ معیار ارزیابی متنوع، از عملکرد خوب و قابل دفاعی برخوردار است و همچنین بر روی داده‌های بانک داخلی (پاسارگاد) نتایج بسیار مطلوب و موثر به‌دست آورد لذا توصیه می‌شود مدیران بانک‌ها و موسسات مالی یا اعتبارسنجی به موضوع روش‌های نیمه‌نظارتی توجهی خاص بنمایند و از آن برای حل مسائلی از قبیل رتبه‌بندی اعتباری مشتریان خود بهره ببرند. همچنین در اینجا از درخت تصمیم برای انتخاب ویژگی استفاده شد که پیشنهاد می‌شود محققان از روش‌های دیگر انتخاب ویژگی مثل فیلترها، بسته‌بندها، روش‌های توکار و فراابتکاری استفاده کنند. این کار سبب افزایش سرعت الگوریتم می‌شود و همچنین تاثیر بسزایی در شناسایی و درک بهتر عوامل تاثیرگذار بر وضعیت اعتباری خواهد داشت. از آن جایی که داده‌های بدون برچسب ارزان‌تر هستند و همیشه در دسترس می‌باشند و به راحتی به دست می‌آیند پیشنهاد می‌شود به کمک الگوریتم‌های دیگر نیمه‌نظارتی و ترکیب آن‌ها، دسته‌بندی‌های قوی‌تری به وجود آورد. همچنین می‌توان از ترکیب دسته‌بندی‌های نظارت‌شده و نیمه‌نظارتی مدل‌های جدیدی ساخت و به نتایج مناسب، تخمین بهتر و تعمیم‌پذیری بالاتری رسید. به‌کارگیری روش‌های گروهی مانند ارتقا و یا رای‌گیری بین دسته‌بندی‌ها یا ترکیب دسته‌بندی‌ها نیز برای تحقیقات آتی پیشنهاد می‌شود.

- pp. 193–199, 2019.
16. B. Feng and X. Wenfang, “Adversarial Semi-supervised Learning for Corporate Credit Ratings”, *Journal of Software*, Vol. 16, No. 6, pp. 259-266, 2021.
 17. N. Ghatasheh, “Business Analytics using Random Forest Trees for Credit Risk Prediction: A Comparison Study”, *International Journal of Advanced Science and Technology*, Vol. 72, pp. 19-30, 2014.
 18. J. Han, M. Kamber and J. Pei, “Data Mining Concepts and Techniques”, *The Morgan Kaufmann Series in Data Management Systems*, New York, 2012.
 19. V. Papastefanopoulos, S. Karlos and S. Kotsiantis, “Using Semi-Supervised Learning Methods for Credit Score Problem”, *Advances in Smart Systems Research*, Vol. 6, No. 2, pp. 28-40, 2017.
 20. O. Chapelle, B. Schölkopf and A. Zien, “Semi-supervised learning”, *The MIT Press*, London, 2006.
 21. J. B. Kruskal, “On the shortest spanning subtree of a graph and the traveling salesman problem”, In *Proceedings of the American Mathematical Society*, Vol. 7, pp. 48-50, 1956.
 22. Sabzevari H., Soleymani M., Noorbakhsh E. (2006), (n.d.) A Comparison between Statistical and Data Mining Methods for Credit Scoring in Case of Limited Available Data.
۱. عباس طلوعی اشلقی، فرناز مقدوری شریبانی، فرید دانشگر، «امتیازدهی اعتباری متقاضیان کارت‌های اعتباری بانک‌ها با استفاده از تکنیک ماشین بردار پشتیبان»، *کنفرانس شهر الکترونیکی*، دوره ۲، ۱۳۸۸.
 ۲. علی اقبالی، سیدحسین رضوی حاج‌آقا، حنان عموزاد مهدیرچی، «ارزیابی مقایسه‌ای عملکرد توابع شایستگی الگوریتم ژنتیک در رتبه‌بندی مشتریان»، *مدیریت صنعتی (دانش مدیریت)*، دوره ۹، شماره ۲، صفحه ۲۶۴-۲۴۵، ۱۳۹۶.
 ۳. غلامرضا کشاورز حداد، حسین آیتی گازار، «مقایسه کارکرد مدل لاجیت و روش درخت‌های طبقه‌بندی و رگرسیون در فرآیند اعتبارسنجی متقاضیان حقیقی برای استفاده از تسهیلات بانکی»، *پژوهش‌های رشد و توسعه اقتصادی*، دوره ۷، شماره ۴، صفحه ۹۷-۷۱، ۱۳۸۶.
 ۴. احمد پویان‌فر، سعید فلاح‌پور، محمدرضا عزیززی، «رویکرد حداقل مربعات ماشین بردار پشتیبان مبتنی بر الگوریتم ژنتیک جهت تخمین رتبه اعتباری مشتریان بانک‌ها»، *مجله مهندسی مالی و مدیریت اوراق بهادار*، دوره ۴، شماره ۱۷، صفحه ۱۵۸-۱۳۳، ۱۳۹۲.
 ۵. محسن مهرآرا، میثم موسایی، مهسا تصویری، آیت حسن‌زاده، «رتبه‌بندی اعتباری مشتریان حقوقی بانک پارسیان»، *فصلنامه مدل‌سازی اقتصادی*، دوره ۳، شماره ۹، صفحه ۱۵۰-۱۲۱، ۱۳۸۸.
 ۶. امین محمدیان حاجی‌کرد، ملیحه اصغرزاده زعفرانی، مصطفی امام‌دوست، «بررسی ریسک اعتباری مشتریان حقوقی با استفاده از مدل ماشین بردار پشتیبان و مدل هیبریدی الگوریتم ژنتیک-مطالعه موردی بانک تجارت»، *مهندسی مالی و اوراق بهادار*، دوره ۷، شماره ۲۷، صفحه ۱۷-۳۲، ۱۳۹۵.
 ۷. سیدمجید شریعت‌پناهی، سیما هاشمی‌برکادهی، «ارائه مدلی برای اعتبارسنجی مشتریان در بانک صنعت و معدن»، *مطالعات تجربی حسابداری مالی*، دوره ۶، شماره ۲۱، صفحه ۸۲-۶۱، ۱۳۸۷.
 ۸. طاهره زارع‌بیدکی، محمدتقی صادقی، حمیدرضا ابوطالبی، «یادگیری نیمه‌نظارتی کرنل مرکب با استفاده از روش‌های یادگیری معیار فاصله»، *پردازش‌های علم و داده*، دوره ۱۴، شماره ۱، صفحه ۷۱-۵۳، ۱۳۹۶.
 ۹. مجتبی کردآبادی، محرم منصوری‌زاده، حسن ختن‌لو، «روش ترکیبی و نیمه‌نظارتی مبتنی بر گراف برای برچسب زنی خودکار تصاویر»، *مجله ماشین بینایی و پردازش تصاویر*، دوره ۶، شماره ۲، صفحه ۸۸-۷۹، ۱۳۹۸.
 ۱۰. عباس طلوعی اشلقی، هاشم نیکومرام، فرناز مقدوری شریبانی، «طبقه‌بندی متقاضیان تسهیلات اعتباری بانک‌ها با استفاده از تکنیک ماشین بردار پشتیبان»، *مجله پژوهش‌های مدیریت*، شماره ۸۴، ۱۳۸۹.
 11. A. Krichene, “Using a naive Bayesian classifier methodology for loan risk assessment evidence from a Tunisian commercial bank”, *Journal of Economics, Finance and Administrative Science*, Vol. 22, No. 42, pp. 3-24, 2017.
 12. C. Huang, M. Chen and C. Wang, “Credit scoring with a data mining approach based on support vector machines”, *Expert Systems with Applications*, Vol. 33, pp.847-856, 2007.
 13. S. Han, “Semi-supervised learning classification based on generalized additive logistic regression for corporate credit anomaly detection”, *IEEE Access*, Vol 8, pp. 199060-199069, 2020.
 14. I.E. Livieris, N. Kiriakidou, A. Kanavos, V. Tampakas and P. Pintelas, “On ensemble SSL algorithms for credit scoring problem”, *Informatics* Vol. 5, No. 4, 40, pp.1-16, 2018.
 15. A. Kim and S. Cho, “An ensemble semi-supervised learning method for predicting defaults in social lending”. *Engineering Applications of Artificial Intelligence*, Vol. 81,