



فصلنامه علمی پژوهشی  
دانش حسابداری و حسابرسی مدیریت  
دوره ۱۳ / شماره ۵۲ (پیاپی) / زمستان ۱۴۰۳  
صفحه ۱۵ تا ۲۸

## تکنیک‌های داده کاوی و پیش بینی تقلب صورتهای مالی

سید جلال احمدی

دانشجوی دکتری گروه حسابداری، واحد سمنان، دانشگاه آزاد اسلامی، سمنان، ایران  
ahmadijalal.acc@gmail.com

خسرو فغانی ماکرانی

دانشیار گروه حسابداری، واحد سمنان، دانشگاه آزاد اسلامی، سمنان، ایران (نویسنده مسئول)  
kh.makrani@chmail.ir

نقی فاضلی

استادیار گروه حسابداری، واحد سمنان، دانشگاه آزاد اسلامی، سمنان، ایران  
fazeli.nphd@gmail.com

تاریخ دریافت: ۱۴۰۰/۰۲/۲۰ تاریخ پذیرش: ۱۴۰۰/۰۵/۳۱

### چکیده

هدف پژوهش حاضر مقایسه تکنیک‌های داده کاوی شبکه عصبی، درخت تصمیم، نزدیک ترین همسایگی و ماشین بردار پشتیبان در پیش بینی صورتهای مالی متقلبانه و غیر متقلبانه است. روش پژوهش توصیفی - کاربردی و قلمرو زمانی نیز از سال ۱۳۸۷ تا ۱۳۹۶ می باشد. در این پژوهش، نسبت‌های مالی برای دو نمونه متقلب و غیر متقلب و روش‌های داده کاوی مورد تجزیه و تحلیل قرار گرفت. فرضیه های آماری نرمال بودن، همگنی و آزمون هم خطی برای نسبت‌های مالی نمونه های متقلب و غیر متقلب، مورد آزمون قرار گرفت. فرضیه نرمال بودن با استفاده از آزمون کولموگروف اسمیرنوف و آزمون شاپیرو ویلک، انجام پذیرفت. سپس ضریب همبستگی پیرسون در خصوص وجود هم خطی مدل برای نسبت‌های مالی و حذف متغیرهای مستقل همبسته مورد بررسی و آزمون قرار گرفت. در مرحله بعد روشهای داده کاوی برای آزمون آنها در پیش بینی تقلب صورتهای مالی و تمایز صورتهای مالی متقلبانه از غیر متقلبانه به کار برده شده است. به طور کلی، نتایج حاصل نشان می‌دهد که روش‌های داده کاوی در تمایز صورتهای مالی متقلبانه از غیر متقلبانه موثر هستند. بدین ترتیب که روش شبکه عصبی ۶۹/۴ درصد، درخت تصمیم ۶۵.۴ درصد، نزدیکترین همسایگی ۶۴/۴ درصد و ماشین بردار پشتیبان ۷۸ درصد پیش بینی صحیح داشته اند.

**واژه‌های کلیدی:** تقلب، داده کاوی، نسبت‌های مالی.

## ۱- مقدمه

تقلب را به عنوان پنهان سازی عمدی حقایق بااهمیت یا تحریف حقایق برای تحریک دیگران به عمل کردن به زیان خود<sup>۱۰</sup> تعریف کردند. گوئل و گانگولی<sup>۱۱</sup> (۲۰۱۲)، تقلب صورت‌های مالی را به عنوان یک تحریف عمدی یا حذف مبالغ و اطلاعات افشا شده در صورت‌های مالی با هدف فریب دادن استفاده کنندگان صورت‌های مالی تعریف کردند. تقلب یک فعالیت بزهکارانه و فریبکارانه با قصد منفعت مالی و سایر منافع می باشد، همچنین تقلب بعنوان یک عمل بزهکارانه در برگیرنده حقه بازی، حيله گری و رفتار غیر منصفانه بوسیله یک شخص فریبکار و متقلب است ( جاویر و همکاران، ۲۰۱۸). تقلب مالی هم بر عملکرد صنایع مختلف و هم بر زندگی روزمره بشر تاثیرگذار است. تقلب می تواند اعتماد به صنعت و کسب و کار را کاهش داده باعث بی ثباتی در سپرده گذاری مردم و در نتیجه تحمیل هزینه اضافی بر زندگی گردد (سادگالی و همکاران، ۲۰۱۹). اغلب تعاریف تقلب اذعان دارند که تقلب به طور کلی و به طور خاص تقلب صورت‌های مالی، ارائه عمدی نادرست اطلاعات (مثلا ارائه نادرست یا حذف مبالغ، افشا، طبقه بندی) از ترازنامه، صورت سود و زیان، صورت جریان وجه نقد و یا یادداشتهای همراه یک شرکت، برای فریب استفاده کنندگان می باشد.

در این پژوهش معیارهای تقلب بر اساس گزارش انجمن بازرسان رسمی تقلب ۱۲ شامل تفاوت زمانی در شناسایی درآمدها و هزینه ها، ثبت درآمدهای واهی یا کم نمایی درآمدها، پنهان کردن بدهیها و هزینه ها یا بیش نمایی بدهیها و هزینه ها، ارزیابی نادرست داراییها و افشای ناکافی در خصوص رویه های حسابداری و اطلاعات با اهمیت می باشد. بنابراین، گزارشات حسابرسی شرکتها مبنای کار قرار میگیرد. به این صورت که فرض می شود شرکتهایی که گزارش غیر مقبول دارند نسبت به شرکتهایی با گزارش مقبول به احتمال بیشتری مرتکب تقلب می شوند. این روش توسط فرقاندوست حقیقی و برواری (۱۳۸۸)، مهام، کردستانی و ترابی (۱۳۹۰)، اعتمادی و زلّی (۱۳۹۲) و برزگری و همکاران (۱۳۹۵)، پیشنهاد گردیده است. بنابراین شرکتهایی که گزارش غیرمقبول ارائه داده اند مشخص شده و سپس در صورتی که گزارش حسابرسی آنها دربرگیرنده معیارهای تقلب مشخص شده توسط انجمن بازرسان رسمی تقلب باشد به عنوان شرکت متقلب شناسایی می گردند.

صورت‌های مالی برای استفاده کنندگان، اطلاعات مالی مفیدی را برای تخصیص منابع به صورت کارآمد و تصمیم گیری های اقتصادی آگاهانه فراهم میکند (کانیگزگرابر<sup>۱</sup>، ۲۰۱۲). تقلب در حسابداری برای دنیای کسب و کار شایع و پر هزینه است (هایز<sup>۲</sup>، ۲۰۱۴). زمانی که تقلب حسابداری آشکار شود، قیمت سهام شرکت کاهش می یابد، و در نتیجه سهامداران متحمل زیان می‌شوند (ورگرس و بوربا<sup>۳</sup>، ۲۰۱۴). تقلب صورت‌های مالی همچنان در حال افزایش است (کراویتز<sup>۴</sup>، ۲۰۱۲). موارد تقلب مالی شرکت های ایالات متحده آمریکا از ۳۷۶ مورد در سال ۲۰۰۵ به ۷۲۵ مورد در سال ۲۰۱۱ افزایش یافته است (دفتر تحقیقات فدرال<sup>۵</sup>، ۲۰۱۲). موارد تقلبات مدیریت (مانند اترون<sup>۶</sup> و دیگران) به شدت تحت تاثیر اقتصاد جهانی و بازارهای سهام قرار دارند (عباسی و همکاران<sup>۷</sup>، ۲۰۱۲). تشخیص تقلب برای جلوگیری از پیامدهای شدید تقلب صورت‌های مالی و ضرر و زیان به ذینفعان حیاتی است. طرح های تقلب، با استفاده از تکنولوژی پیشرفته پیچیده تر شده و در نتیجه شناسایی آن دشوار شده است (رانکیو<sup>۸</sup>، ۲۰۱۶). علی رغم اهمیت و ضرورت توجه به پدیده صورت های مالی متقلبانه لیست شرکت های متقلب و مصادیق تقلب در صورت های مالی توسط هیچ ارگان یا نهادی در کشور بررسی و ارایه نمی شود. نهادهایی از قبیل سازمان بورس و اوراق بهادار اطلاعات احتمالی مربوط به هر گونه تحریف و به طور خاص تقلب در صورت های مالی را در اختیار عموم و تحلیل گران قرار نمی دهند. موارد تخلف بررسی شده در سازمان بورس اوراق بهادار نیز در صورت صدور رای از طریق محاکم قضایی به طور خصوصی اطلاع رسانی می شود اما برای استفاده عمومی منتشر نمی شود. (کاظمی، سجادی و رحمانی ۱۳۹۵). هدف این پژوهش غیر تجربی کمی، مقایسه تکنیکهای داده کاوی شبکه عصبی، درخت تصمیم، نزدیک ترین همسایگی و ماشین بردار پشتیبان در تمایز صورت‌های مالی متقلبانه از غیر متقلبانه موثر هستند.

## ۲- مبانی نظری و پیشینه پژوهش

## ۲-۱- مبانی نظری

ادبیات حسابداری و حسابرسی هیچ توافقی در خصوص تعریف تقلب صورت های مالی نشان نداده است. وارگر و بوربا<sup>۹</sup> (۲۰۱۴)

<sup>7</sup>Abbasi, Albrecht, Vance, & Hansen, 2012

<sup>8</sup> Ruankaew, 2016

<sup>9</sup> Wuerges and Borba (2014)

<sup>10</sup> to act to their own detriments

<sup>11</sup> Goel and Gangolly (2012)

<sup>12</sup> Association of Certified Fraud Examiners

<sup>1</sup> Königsgruber, 2012

<sup>2</sup> Hays, 2014

<sup>3</sup> Wuerges & Borba, 2014

<sup>4</sup> Kravitz, 2012

<sup>5</sup> Federal Bureau of Investigation [FBI], 2012

<sup>6</sup> Enron

## ۲-۲- نسبت‌های مالی هشداردهنده احتمال تقلب

آگاهی از علائم تقلب عامل مهمی در جلوگیری و تشخیص تقلب است (ورگرس و بوربا، ۲۰۱۴). هرچند جرم تقلب به ندرت دیده می‌شود، علائم یا شاخص‌های تقلب مشاهده می‌شود (آلبرت، ۲۰۱۲). در پژوهش حاضر، نسبت‌های مالی (متغیرهای مستقل) برای پیش‌بینی خطر بالقوه تقلب صورت‌های مالی به دلیل سادگی، انعطاف‌پذیری و محبوبیت آن در میان جامعه مالی مورد استفاده قرار می‌گیرد. در این پژوهش، ۳۵ نسبت مالی با استفاده از تکنیک‌های داده کاوی و اثربخشی آنها در پیش‌بینی احتمال تقلب صورت‌های مالی و تمایز گزارش مالی متقلبانه از غیر متقلبانه آزمون می‌شود. آگاهی از علائم تقلب، عامل مهمی در جلوگیری و تشخیص تقلب است (ورگرس و بوربا، ۲۰۱۴).

## ۲-۳- داده کاوی و کاربردهای آن

در دو دهه قبل توانایی‌های فنی بشر در تولید و جمع‌آوری داده‌ها به سرعت افزایش یافته است. عواملی نظیر استفاده گسترده از بارکد برای تولیدات تجاری، به خدمت گرفتن کامپیوتر در کسب و کار، علوم، خدمات دولتی و پیشرفت در وسائل جمع‌آوری داده، از اسکن کردن متون و تصاویر تا سیستم‌های سنسور از دور ماهواره‌ای، در این تغییرات نقش مهمی دارند (هربرت، ای. ادلشتاین، ۱۹۹۹، ۴). بطور کلی استفاده همگانی از وب و اینترنت به عنوان یک سیستم اطلاع‌رسانی جهانی ما را مواجه با حجم زیادی از داده و اطلاعات می‌کند. این رشد انفجاری در داده‌های ذخیره شده، نیاز مبرم به وجود تکنولوژی‌های جدید و ابزارهای خودکاری را ایجاد کرده که به صورت هوشمند به انسان یاری رسانند تا این حجم زیاد داده را به اطلاعات و دانش تبدیل کند، داده کاوی به عنوان یک راه حل برای این مسائل مطرح می‌باشد. در یک تعریف غیر رسمی داده کاوی فرآیندی است، خودکار برای استخراج الگوهایی که دانش را بازنمایی می‌کنند، که این دانش به صورت ضمنی در پایگاه داده‌های عظیم، انبار داده‌ها و دیگر مخازن بزرگ اطلاعات، ذخیره شده است. داده کاوی بطور همزمان از چندین رشته علمی بهره می‌برد نظیر؛ تکنولوژی پایگاه داده، هوش مصنوعی، یادگیری ماشین، شبکه‌های عصبی، آمار، شناسایی الگو، سیستم

های مبتنی بر دانش، حصول دانش، بازیابی اطلاعات، محاسبات سرعت بالا، و بازنمایی بصری داده. داده کاوی در اواخر دهه ۱۹۸۰ پدیدار گشته، در دهه ۱۹۹۰ گام‌های بلندی در این شاخه از علم برداشته شده و انتظار می‌رود در این قرن به رشد و پیشرفت خود ادامه دهد (دیوید و همکاران، ۲۰۰۱). واژه‌های «داده کاوی» و «کشف دانش در پایگاه داده» اغلب به صورت مترادف یکدیگر مورد استفاده قرار می‌گیرند. کشف دانش در پایگاه داده فرایند شناسایی درست، ساده، مفید، و نهایتاً الگوها و مدل‌های قابل فهم در داده‌ها می‌باشد. داده کاوی، مرحله‌ای از فرایند کشف دانش می‌باشد و شامل الگوریتم‌های مخصوص داده کاوی است، بطوریکه، تحت محدودیت‌های مؤثر محاسباتی قابل قبول، الگوها و یا مدل‌ها را در داده کشف می‌کند (هربرت، ای. ادلشتاین، ۱۹۹۹، ۱۳).

## ۲-۳-۱- شبکه‌های عصبی

شبکه‌های عصبی از پرکاربردترین و عملی‌ترین روش‌های مدل سازی مسائل پیچیده و بزرگ که شامل صدها متغیر هستند می‌باشد. شبکه‌های عصبی می‌توانند برای مسائل کلاس بندی (که خروجی یک کلاس است) یا مسائل رگرسیون (که خروجی یک مقدار عددی است) استفاده شوند. هر شبکه عصبی شامل یک لایه ورودی ۱۴ می‌باشد که هر گره در این لایه معادل یکی از متغیرهای پیش‌بینی می‌باشد. گره‌های موجود در لایه میانی وصل می‌شوند به تعدادی گره در لایه نهان ۱۵. هر گره ورودی به همه گره‌های لایه نهان وصل می‌شود. گره‌های موجود در لایه نهان می‌توانند به گره‌های یک لایه نهان دیگر وصل شوند یا می‌توانند به لایه خروجی ۱۶ وصل شوند. لایه خروجی شامل یک یا چند متغیر خروجی می‌باشد (دیوید و همکاران، ۲۰۰۱، ۱۷).

<sup>10</sup> Data visualization

<sup>11</sup> David Hand, Heikki Mannila, Padhraic Smyth.

<sup>12</sup> Knowledge Discovery in Database

<sup>13</sup> Herbert A. Edelstein

<sup>14</sup> Input Layer

<sup>15</sup> Hidden Layer

<sup>16</sup> Output Layer

<sup>17</sup> David Hand, Heikki Mannila, Padhraic Smyth.

<sup>1</sup> Wuerges & Borba, 2014

<sup>2</sup> Albrech, 2012

<sup>3</sup> Wuerges & Borba, 2014

<sup>4</sup> Herbert A. Edelstein

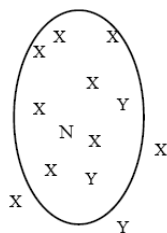
<sup>5</sup> Data warehouses

<sup>6</sup> Knowledge-based system

<sup>7</sup> Knowledge-acquisition

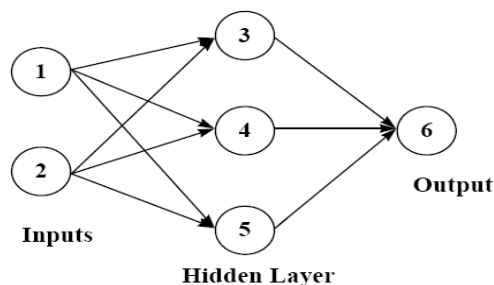
<sup>8</sup> Information retrieval

<sup>9</sup> High-performance computing



شکل (۳): محدوده همسایگی

(بیشتر همسایه ها در دسته X قرار گرفته اند)

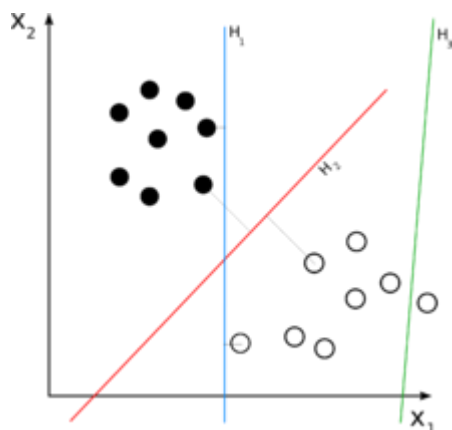


شکل (۱): شبکه عصبی با یک لایه نهان

### ۲-۳-۴- ماشین بردار پشتیبان

در روش ماشین بردار پشتیبان، هر نمونه داده را به عنوان یک نقطه در فضای  $n$ -بعدی روی نمودار پراکندگی داده‌ها ترسیم می‌کنیم. به طوری که مقدار هر ویژگی مربوط به داده‌ها، یکی از مؤلفه‌های مختصات نقطه روی نمودار را مشخص کند. سپس، با ترسیم یک خط راست، داده‌های مختلف و متمایز از یکدیگر را دسته‌بندی می‌نماییم.

در این روش نیز مانند سایر روش‌های یادگیری ماشین، از داده‌های گروه‌های آموزش و آزمایش استفاده می‌شود.



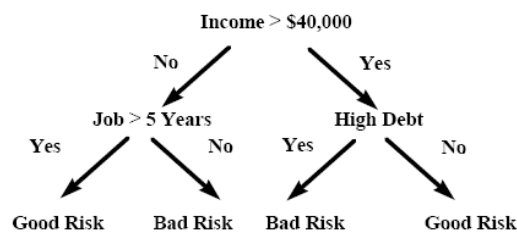
شکل (۴): ماشین بردار پشتیبان

### ۲-۴- پیشینه‌ی پژوهش

اشرف آکل الساید<sup>۳</sup> (۲۰۱۷) در پژوهشی به بررسی اثربخشی ۱۸ نسبت مالی پرداخت. همچنین اثربخشی قانون بنفورد و تکنیک‌های داده کاوی در تشخیص گزارش مالی متقلبانه از غیرمتقلبانه و پیش بینی احتمال تقلب صورتهای مالی (متغیر وابسته) شرکت‌های عمومی ایالات متحده که در سازمان بورس و اوراق بهادار ایالات متحده آمریکا ثبت شده اند را مورد بررسی قرارداد. نتایج نشان داد که درخت تصمیم، رگرسیون لجستیک، قانون بنفورد

### ۲-۳-۲- درختهای تصمیم

درخت‌های تصمیم روشی برای نمایش یک سری از قوانین هستند که منتهی به یک رده یا مقدار می‌شوند. برای مثال، می‌خواهیم متقاضیان وام را به دارندگان ریسک اعتبار خوب و بد تقسیم کنیم. شکل زیر درخت تصمیم را که این مسئله را حل می‌کند نشان می‌دهد و همه مؤلفه‌های اساسی یک درخت تصمیم در آن نشان داده شده است: نود تصمیم، شاخه‌ها و برگ-ها (هربرت.ای. ادلشتاین، ۱۹۹۹).



شکل (۲): درخت تصمیم‌گیری

### ۲-۳-۲- نزدیکترین همسایگی

هنگام تلاش برای حل مسائل جدید، افراد معمولاً به راه‌حل‌های مسائل مشابه که قبلاً حل شده‌اند مراجعه می‌کنند. نزدیکترین همسایگی  $k$  یک تکنیک دسته‌بندی است که از نسخه‌ای از این متد استفاده می‌کند. در این روش تصمیم‌گیری اینکه یک مورد جدید در کدام دسته قرار گیرد با بررسی تعدادی ( $k$ ) از شبیه‌ترین موارد یا همسایه‌ها انجام می‌شود. تعداد موارد برای هر کلاس شمرده می‌شوند، و مورد جدید به دسته‌ای که تعداد بیشتری از همسایه‌ها به آن تعلق دارند نسبت داده می‌شود.

<sup>3</sup> Ashraf akl elsayed

<sup>1</sup> Herbert A. Edelstein

<sup>2</sup> K-nearest neighbor(k-NN)

شرکتهای پذیرفته شده در بورس اوراق بهادار تهران با استفاده از برخی نسبتهای مالی مرتبط تشخیص داده شده است. نمونه آماری تحقیق شامل ۶۸ شرکت در قالب ۳۴ شرکت دارای نشانه های تقلب و ۳۴ شرکت فاقد نشانه های تقلب است. همچنین ۹ نسبت مالی به عنوان پیش بینی کننده های بالقوه برای آزمون انتخاب شده اند. از روش رگرسیون لجستیک جهت تدوین مدل برای شناسایی عوامل مرتبط با تقلب صورت‌های مالی استفاده شده است. این مدل در طبقه بندی صحیح نمونه مورد نظر در این تحقیق از نرخ دقت ۸۳.۸ درصد برخوردار بود.

رهنمای رودپشتی (۱۳۹۱) در پژوهشی با عنوان داده کاوی و کشف تقلب های مالی نشان داد که اولاً تکنیک های داده کاوی، در شناسایی صورت های مالی متقلبان سودمند هستند. ثانیاً داده کاوی، به عنوان کانون هدایت فکر در مدیریت کسب و کارها جهت کشف تقلب می تواند مورد توجه قرار گیرد.

صفر زاده (۱۳۸۹) در پژوهشی با عنوان توانایی نسبت های مالی در کشف تقلب در گزارشگری مالی با استفاد از تحلیل لاجیت در داد های مقطعی به بررسی نقش داده های حسابداری در ایجاد یک الگو برای کشف عوامل مرتبط با تقلب در گزارشگری مالی می پردازد. نتایج تحقیق حکایت از عملکرد مناسب الگو در طبقه بندی شرکت های نمونه داشت به گونه ای که درصد صحت طبقه بندی الگو از ۸۲.۹۸ درصد تجاوز نمود.

### ۳- فرضیه پژوهش

تکنیکهای داده کاوی شبکه عصبی، درخت تصمیم، نزدیک ترین همسایگی و ماشین بردار پشتیبان در تمایز صورت‌های مالی متقلبان از غیر متقلبان موثر هستند.

### ۴- روش شناسی پژوهش

این تحقیق از لحاظ هدف کاربردی و از نوع پژوهشهای تجربی است و با توجه به اینکه از اطلاعات تاریخی استفاده شده است در دسته پژوهشهای شبه آزمایشی قرار می گیرد. جامعه آماری پژوهش حاضر شرکتهای پذیرفته شده در بورس اوراق بهادار تهران و قلمرو زمانی آن هم بین سالهای ۱۳۸۷ الی ۱۳۹۶ است و داده‌های آن نیز طی این دوره مورد اشاره برای تجزیه و تحلیل آماری (توصیفی و استنباطی) با استفاده از بانک های اطلاعاتی و نرم افزارهای موجود گردآوری شده است. متغیر وابسته در این پژوهش، تقلب در صورت‌های مالی است که از ماهیت کیفی برخوردار بوده و دارای مقیاس سنجش اسمی است. در اندازه گیری این متغیر، به شرکتهای متقلب عدد یک و به شرکتهای

و مدل های شبکه عصبی به درستی ۰.۸۹٪، ۰.۹۱٪، ۰.۹۲٪ و ۰.۹۹۲٪ از موارد تقلب و عدم تقلب را پیش بینی کرده اند و بنابراین می تواند به عنوان ابزار پیش بینی تقلب مورد استفاده قرار گیرند.

وسکی (۲۰۱۳)<sup>۱</sup> پژوهشی با عنوان تغییرات قیمت سهم و نسبت قیمت به سود به عنوان شاخص تقلب گزارشات مالی انجام دادند، هدف آزمون فرضیه بازار کارا توسط تعیین محدوده تغییرات در نسبت های قیمت سهم و قیمت به سود قبل از اعلام عمومی تقلب پیش بینی شده یک شرکت بود. نتایج این پژوهش استفاده از معیارهای کمی را برای تشخیص تقلب پیشنهاد می کند. به علاوه محققین می توانند این یافته ها را برای ایجاد یک مدل قوی برای دقت بیشتر تشخیص تقلب به کار ببرند.

ایگو<sup>۲</sup> (۲۰۱۷) در پژوهشی با عنوان نسبتهای نقدینگی و سهم بازار به عنوان شاخص تقلب، بررسی کرد که تغییرات در میزان نقدینگی و سهم بازار میتواند شاخصی از فعالیتهای متقلبان در یک شرکت باشد و اینکه آیا گزارشهای مالی متقلبان منتشر شده میتواند موجب تغییر در قیمت سهام شرکت شود، تجزیه و تحلیل داده ها نشان داد که بین مولفه های نسبت نقدینگی، نسبت سهم بازار و تغییر در قیمت سهام رابطه معنی داری وجود دارد.

خواجوی و ابراهیمی (۱۳۹۶) در پژوهشی با عنوان مدل سازی متغیرهای اثرگذار برای کشف تقلب در صورت‌های مالی با استفاده از تکنیکهای داده کاوی به بررسی این مساله پرداختند که آیا میتوان از طریق شناسایی و انتخاب متغیرهای اثرگذار در کشف تقلب در صورت‌های مالی و با به کارگیری تکنیکهای داده کاوی مدلی برای کشف تقلب در صورت‌های مالی شرکتهای پذیرفته شده در بورس اوراق بهادار تهران ارائه کرد. یافته های پژوهش بیانگر وجود شواهدی دال بر عملکرد مناسب مدل های پیشنهادی و برتری الگوریتم جنگل تصادفی و شبکه بیزین برای پیشبینی تقلب در صورت‌های مالی است. نتایج حاصل از انتخاب ویژگی به روش مبتنی بر همبستگی حاکی از سودمندی متغیرهای نسبت پوشش بهره، نسبت حسابهای دریافتنی به کل داراییها، نسبت موجودی کالا به فروش خالص، نسبت نقدی، لگاریتم طبیعی فروش، نسبت سود خالص به فروش و نسبت جمع داراییهای جاری به کل داراییها برای کشف تقلب بود.

اعتمادی و زلّی (۱۳۹۲) در پژوهشی با عنوان کاربرد رگرسیون لجستیک در شناسایی گزارشگری مالی متقلبان به بررسی این موضوع پرداختند که داده های صورت‌های مالی حسابرسی شده این توانایی را دارند که هرگونه تقلب صورت‌های مالی را کشف نمایند. در این تحقیق، تقلب صورت‌های مالی در

<sup>2</sup> Ego

<sup>1</sup> WESKE

شده به سود عملیاتی، هزینه استهلاک به اموال، ماشین آلات و تجهیزات، استقلال هیات مدیره، تعداد جلسات هیات مدیره، درصد مالکیت سهام هیات مدیره، درصد مالکیت سهامداران عمده (بالای ۵ درصد) و درصد مالکیت سهامداران نهادی می باشند.

**۵-۱- برازش مدل شبکه عصبی:** داده ها در دو گروه آموزش و آزمایش تقسیم بندی می شوند.

در گروه آموزش حدود ۷۰ درصد مشاهدات که معادل ۵۷۵ مورد می باشد قرار میگیرد، ۳۰ درصد مابقی داده ها یعنی تعداد ۲۴۵ مورد نیز در گروه آزمایش قرار می گیرند. تابع فعالسازی در لایه پنهان Hyperbolic tangent و در لایه خروجی Softmax می باشد. در این بخش پس از آموزش و برازش مدل شبکه عصبی میزان اعتبار نتایج با توجه به فراوانی پیش بینی های صحیح در هر یک از گروه داده های آزمایش بررسی می شود.

#### ۵-۲- برازش مدل درخت تصمیم

روش دیگری که در حوزه ی داده کاوی که مورد استفاده قرار گرفته است روش درخت تصمیم می باشد. هدف بررسی پیش بینی های صحیح شرکت های متقلب و غیرمتقلب می باشد. روش مورد استفاده برای تصمیم گیری ها انتخاب براساس مقدار آماره کای دو (CHAD) بوده است. متغیرهای ورودی شامل متغیرهای مستقل معنادار براساس مدل رگرسیون لجستیک پژوهش احمدی و همکاران (۱۳۹۸) بوده و عبارتند از کل دارایی ها به فروش، سود عملیاتی به فروش، دوره گردش موجودی، چرخه تبدیل وجه نقد می باشند.

#### ۵-۳- برازش مدل نزدیک ترین همسایگی

نزدیک ترین همسایگی<sup>۱</sup> (K) یک روش ناپارامتری است که در داده کاوی، یادگیری ماشین و تشخیص الگو مورد استفاده قرار می گیرد. در این روش نیز مانند سایر روش های یادگیری ماشین، از داده های گروه های آموزش و آزمایش، استفاده می شود. در مرحله پیش پردازش با هدف کسب نتایج بهتر داده ها نرمال سازی می شوند. با نرمال سازی داده ها مقادیر بین «۰» و «۱» قرار می گیرند. به منظور تعیین بهترین مقدار K همسایگی در مدل روی داده های آموزشی با دادن مقادیر مختلف برای K برازش داده شده و ضریب دقت مدل برای هر مقدار K محاسبه می شود.

غیرمتقلب، عددصفر تخصیص داده می شود. متغیرهای مستقل، نسبتهای مالی هستند؛ در ابتدا ۳۵ نسبت مالی انتخاب شدند. که با تحلیل همبستگی و آزمون T منجر به انتخاب نهایی تعداد ۱۹ متغیر مستقل شد؛ که اطلاعات معنادار و غیر همپوشی را ارائه می کنند. سپس آمار توصیفی شامل میانگین، میانه، انحراف او انحراف از میانگین برای هر یک از متغیرهای کمی و به تفکیک نوع شرکت متقلب یا غیرمتقلب بررسی شد. در مرحله بعد برای متغیرهای مستقل، همبستگی بین احتمال تقلب و متغیرهای کمی سنجیده شد. همچنین، با توجه به ضرورت بررسی برخی فرضیات زیربنایی پیش از انجام تحلیل، نرمال بودن (آزمون های کلموگروف اسمیرنوف و شاپیرو ویلک) وجود همخطی و همسانی واریانس ها در مورد متغیرهای مستقل مورد بحث، بررسی شد و تحلیل داده ها با استفاده از نرم افزار EXCEL و SPSS۲۰ صورت گرفته است. در نهایت نیز از آماره آزمون کای دو برای تصمیم گیری در مورد معناداری تفاوت توزیع های فراوانی در هر بخش استفاده شد. این تحقیق شامل شرکتهای پذیرفته شده در بورس اوراق بهادار تهران است که شرایط زیر را داشته باشند.

- از تاریخ ۱۳۸۶/۱/۱ بعد عضو سازمان بورس و اوراق بهادار باشند.
- اطلاعات آنها قابل تهیه و در دسترس باشد.
- جزء شرکت های واسطه مالی، سرمایه گذاری، بانکها، بیمه و بازنشستگی، موسسات اعتباری و شرکتهای چند رشته ای صنعتی نباشند.

#### ۵- مدل و متغیرهای پژوهش

نسبتهای مالی شامل وجه نقد عملیاتی به سود خالص، کل دارایی ها به فروش، سود خالص به حقوق صاحبان سهام، فروش به سود خالص، سود قبل از بهره و مالیات به کل دارایی، سود هر سهم، سود ناخالص به فروش، سود عملیاتی به فروش، سود انباشته به کل داراییها، سود خالص به دارایی های ثابت، کل بدهی به حقوق صاحبان سهام، کل بدهیها به کل داراییها، دارایی های جاری به بدهی های جاری، دوره گردش موجودی، دوره پرداخت بدهی، دوره وصول مطالبات، حسابهای دریافتنی به متوسط فروش در هر روز، چرخه تبدیل وجه نقد، Z آلتمن، بدهی های بلندمدت به کل داراییها، سرمایه در گردش به کل دارایی، بدهی بلندمدت به حقوق صاحبان سهام، بهای تمام شده کالای فروش رفته به فروش، فروش به حسابهای دریافتنی، فروش به موجودی، فروش به کل داراییها، فروش به کل داراییهای ثابت و بهای تمام شده کالای فروش رفته به موجودی کالا، مالیات بر درآمد پرداخت

<sup>۱</sup> k-Nearest Neighbors

#### ۴-۵- ماشین بردار پشتیبان

در روش ماشین بردار پشتیبان یا همان الگوریتم SVM، هر نمونه داده را به عنوان یک نقطه در فضای n-بعدی روی نمودار پراکندگی داده‌ها ترسیم می‌کنیم به طوری که مقدار هر ویژگی مربوط به داده‌ها، یکی از مؤلفه‌های مختصات نقطه روی نمودار را مشخص کند. سپس، با ترسیم یک خط راست، داده‌های مختلف و متمایز از یکدیگر را دسته‌بندی می‌نماییم. در این روش نیز مانند سایر روش‌های یادگیری ماشین، از داده‌های گروه‌های آموزش و آزمایش، استفاده می‌شود. در مرحله پیش پردازش با هدف کسب نتایج بهتر داده‌ها نرمال سازی شده‌اند. با نرمال سازی داده‌ها مقادیر بین «۰» و «۱» قرار می‌گیرند.

پس از آموزش و برازش مدل، میزان اعتبار نتایج با توجه به فراوانی پیش بینی‌های صحیح در گروه داده‌های آزمایش بررسی و گزارش می‌شود.

#### ۶- یافته‌های پژوهش

##### ۶-۱- آمار توصیفی

آمار توصیفی که شامل میانگین، میانه، انحراف استاندارد و انحراف از میانگین برای هر یک از متغیرهای کمی و به تفکیک نوع شرکت متقلب یا غیرمتقلب بوده در جدول شماره (۱)، ارائه شده است. با بررسی این اطلاعات به طور تقریبی تفاوت مقادیر هر یک از این نسبت‌های مالی مربوط به شرکت‌های متقلب را در مقایسه با شرکت‌های غیرمتقلب می‌توان مشاهده نمود.

جدول شماره (۱): توصیفی مربوط به متغیرهای مدل

غیر متقلب				متقلب				شرح متغیرهای مدل پژوهش
انحراف معیار میانگین	انحراف معیار	میانه	میانگین	انحراف معیار میانگین	انحراف معیار	میانه	میانگین	
8.298	168.031	0.786	-4.722	0.726	14.697	1.041	1.890	وجه نقد عملیاتی به سود خالص
0.066	1.334	1.347	1.601	0.031	0.633	1.138	1.261	کل دارایی‌ها به فروش
0.189	3.823	0.195	-0.009	0.074	1.504	0.243	0.338	سود خالص به حقوق صاحبان سهام
135.005	2733.651	6.793	199.890	6.962	140.962	7.396	37.289	فروش به سود خالص
0.008	0.157	0.108	0.115	0.006	0.127	0.147	0.168	سود قبل از بهره و مالیات به کل دارایی
0.000	0.001	0.000	0.000	0.000	0.001	0.000	0.001	سود هر سهم
0.009	0.173	0.173	0.196	0.008	0.156	0.205	0.245	سود ناخالص به فروش
0.014	0.292	0.105	0.104	0.008	0.153	0.142	0.177	سود عملیاتی به فروش
0.015	0.296	0.077	0.044	0.009	0.185	0.115	0.119	سود انباشته به کل داراییها
0.059	1.188	0.215	0.483	1.239	25.079	0.402	2.036	سود خالص به دارایی‌های ثابت
1.172	23.740	1.736	3.934	1.544	31.264	1.572	-0.496	کل بدهی به حقوق صاحبان سهام
0.015	0.306	0.666	0.690	0.017	0.354	0.641	0.654	کل بدهیها به کل داراییها
0.028	0.576	1.109	1.223	0.048	0.982	1.291	1.439	داراییهای جاری به بدهی‌های جاری
5.839	118.225	131.025	149.530	5.491	111.181	133.443	154.672	دوره گردش موجودی
8.983	181.897	65.258	104.140	3.480	70.472	39.828	58.449	دوره پرداخت بدهی
6.996	141.660	105.021	144.872	4.327	87.614	87.470	105.983	دوره وصول مطالبات
7.800	157.942	115.548	154.972	4.758	96.339	89.525	114.186	حسابهای دریافتی به متوسط فروش در هر روز
7.844	158.838	186.333	190.262	6.856	138.822	186.036	202.206	چرخه تبدیل وجه نقد
0.114	2.306	1.534	1.509	0.047	0.955	1.948	2.045	Z آلتمن
0.007	0.149	0.052	0.084	0.015	0.301	0.049	0.111	بدهی‌های بلندمدت به کل داراییها
0.014	0.288	0.066	0.049	0.011	0.215	0.155	0.137	سرمایه در گردش به کل دارایی
1.172	23.740	1.736	3.934	1.544	31.264	1.572	-0.496	بدهی بلندمدت به حقوق صاحبان سهام
0.009	0.173	0.827	0.804	0.008	0.156	0.795	0.755	بهای تمام شده کالای فروش رفته به فروش
4.020	81.391	3.159	15.551	1.981	40.109	4.077	13.323	فروش به حسابهای دریافتی
0.204	4.128	3.677	4.878	0.142	2.881	3.618	4.332	فروش به موجودی
0.030	0.598	0.742	0.879	0.025	0.511	0.879	0.992	فروش به کل داراییها
0.281	5.691	4.338	5.889	2.297	46.503	4.364	9.199	فروش به کل داراییهای ثابت

غیرمتقلب				متقلب				شرح متغیرهای مدل پژوهش
انحراف معیار میانگین	انحراف معیار	میانگین	میانگین	انحراف معیار میانگین	انحراف معیار	میانگین	میانگین	
0.183	3.705	2.700	3.967	0.126	2.561	2.583	3.348	بهای تمام شده کالای فروش رفته به موجودی کالا
0.030	0.597	0.065	0.110	0.135	2.727	0.107	0.275	مالیات بر درآمد پرداخت شده به سود عملیاتی
0.003	0.069	0.100	0.107	0.164	3.328	0.108	0.355	هزینه استهلاک به اموال، ماشین آلات و تجهیزات
0.009	0.173	0.600	0.619	0.007	0.151	0.600	0.645	استقلال هیات مدیره
0.300	6.067	12.000	14.815	0.302	6.117	12.000	15.688	تعداد جلسات هیئت مدیره
0.121	2.454	0.700	0.739	0.011	0.227	0.727	0.664	درصد مالکیت سهام هیات مدیره
0.007	0.146	0.820	0.774	0.007	0.145	0.820	0.772	سهامداران عمده (بالای ۵ درصد)
0.017	0.340	0.305	0.401	0.016	0.331	0.355	0.412	درصد مالکیت سهامداران نهادی

منبع: یافته های پژوهش

۲ باشد. که البته با توجه به این شاخص ها نرمال بودن قابل تایید نیست، بنابراین آماره ی آزمون کلموگروف اسمیرنوف و شاپیرو ویلکز بررسی می شود. در این دو آزمون فرضیه صفر بیانگر نرمال بودن متغیر مورد بحث است. چنانچه در هر یک از متغیرهای گروه شرکت های متقلب یا غیرمتقلب مقدار (p-value) حاصل شده، برای این آزمون ها بیش از سطح خطای آزمون یعنی ۰.۰۵ باشد فرضیه ی نرمال بودن پذیرفته می شود. با توجه به توضیحات ارائه شده نتایج حاصل از جدول شماره (۲) نشان دهنده ی عدم نرمال بودن متغیرهای مستقل کمی است. البته توجه به این نکته ضروری است که با توجه به حجم بالای اطلاعات موجود استفاده از ضریب همبستگی پیرسون جهت محاسبه ی همبستگی میان متغیرها و آزمون های پارامتری تی جهت مقایسه ی میانگین گروه ها، مشکلی ایجاد نمی کند.

## ۶-۲-آمار استنباطی

با توجه به ضرورت بررسی برخی فرضیات زیربنایی پیش از انجام تحلیل، نرمال بودن، وجود همخطی و همسانی واریانس ها در مورد متغیرهای مستقل مورد بحث، بررسی شده است.

## ۶-۲-۱-آزمون بررسی نرمال بودن متغیرهای کمی

به منظور تعیین نرمال بودن متغیرهای کمی مربوط به نسبت های مالی مورد نظر علاوه بر بررسی شاخص های چولگی و کشیدگی، آماره آزمون های کلموگروف اسمیرنوف و شاپیرو ویلکز نیز محاسبه شده است. همچنین نمودارهای Q-Qplot (نمودار چندک های توزیع نرمال) و نمودار جعبه ای نیز ترسیم و بررسی شده است. در ارتباط با شاخص های چولگی و کشیدگی انتظار می رود در صورت نرمال بودن متغیرها مقدار حاصل بین ۲- و

جدول شماره (۲): بررسی نرمال بودن و همسانی واریانس های متغیرهای کمی مدل پژوهش

شرح متغیرهای مدل پژوهش	توصیف متغیرهای پژوهش			آماره <sup>۱</sup> کلموگروف اسمیرنوف <sup>۲</sup>	آماره <sup>۳</sup> شاپیرو ویلکز <sup>۴</sup>	سطح معنی داری	بررسی همسانی واریانس	
	تعداد مشاهدات (N)	چولگی <sup>۵</sup>	کشیدگی <sup>۵</sup>				آزمون لون	سطح معنی داری
وجه نقد عملیاتی به سود خالص	820	774.284	-27.414	.446	.045	.000	3.146	.077
کل دارایی ها به فروش	820	83.158	6.849	.153	.598	.000	17.391	.000
سود خالص به حقوق صاحبان سهام	820	493.864	-18.279	.387	.121	.000	2.290	.131
فروش به سود خالص	820	789.761	27.865	.444	.033	.000	5.015	.025
سود قبل از بهره و مالیات به کل دارایی	820	9.723	-7.19	.091	.903	.000	.468	.494
سود هر سهم	820	10.741	2.260	.177	.785	.000	2.036	.154
سود ناخالص به فروش	820	5.855	-1.167	.090	.940	.000	.085	.771
سود عملیاتی به فروش	820	54.781	-4.716	.166	.712	.000	5.494	.019

<sup>1</sup> Statistic

<sup>2</sup> Kolmogorov-Smirnova

<sup>3</sup> Shapiro-Wilk

<sup>4</sup> Kurtosis

<sup>5</sup> Skewness



شرح متغیرهای مدل پژوهش	توصیف متغیرهای پژوهش			آماره <sup>۱</sup> کلموگروف اسمیرنف <sup>۲</sup>	آماره شاپیرو ویلک <sup>۳</sup>	سطح معنی داری	بررسی همسانی واریانس	
	تعداد مشاهدات (N)	چولگی <sup>۴</sup>	کشیدگی <sup>۵</sup>				آزمون لون	سطح معنی داری
سود انباشته به کل داراییها	820	54.130	-4.676	.160	.745	.000	6.690	.010
سود خالص به دارایی های ثابت	820	772.390	27.488	.432	.035	.000	3.336	.068
کل بدهی به حقوق صاحبان سهام	820	260.427	-7.888	.416	.156	.000	.003	.957
کل بدهیها به کل داراییها	820	49.758	5.283	.163	.654	.000	.091	.762
داراییهای جاری به بدهی های جاری	820	80.029	6.719	.176	.601	.000	2.905	.089
دوره گردش موجودی	820	23.719	3.659	.126	.733	.000	.847	.358
دوره پرداخت بدهی	820	153.111	10.594	.281	.370	.000	19.390	.000
دوره وصول مطالبات	820	29.203	3.787	.147	.737	.000	19.754	.000
حسابهای دریافتی به متوسط فروش در هر روز	820	50.194	4.647	.155	.719	.000	14.653	.000
چرخه تبدیل وجه نقد	820	2.969	.935	.056	.946	.000	3.729	.054
Z آلتمن	820	306.508	-13.795	.184	.476	.000	2.642	.104
بدهی های بلندمدت به کل داراییها	820	139.164	9.832	.327	.346	.000	3.679	.055
سرمایه در گردش به کل دارایی	820	31.198	-3.094	.070	.849	.000	6.438	.011
بدهی بلندمدت به حقوق صاحبان سهام	820	260.427	-7.888	.416	.156	.000	.003	.957
بهای تمام شده کالای فروش رفته به فروش	820	5.855	.167	.090	.940	.000	.085	.771
فروش به حسابهای دریافتی	820	128.735	10.651	.412	.170	.000	2.118	.146
فروش به موجودی	820	24.823	3.615	.154	.724	.000	17.882	.000
فروش به کل داراییها	820	11.997	2.614	.133	.803	.000	.378	.539
فروش به کل داراییهای ثابت	820	387.111	19.425	.412	.088	.000	4.096	.043
بهای تمام شده کالای فروش رفته به موجودی کالا	820	24.197	3.587	.183	.701	.000	17.310	.000
مالیات بر درآمد پرداخت شده به سود عملیاتی	820	724.566	26.063	.411	.068	.000	1.184	.277
هزینه استهلاک به اموال، ماشین آلات و تجهیزات	820	409.671	20.230	.471	.035	.000	6.796	.009
استقلال هیات مدیره	820	.219	-4.77	.245	.858	.000	.653	.419
تعداد جلسات هیئت مدیره	820	9.305	2.805	.292	.606	.000	1.596	.207
درصد مالکیت سهام هیات مدیره	820	785.610	27.728	.434	.064	.000	1.750	.186
سهامداران عمده (بالای ۵ درصد)	820	2.260	-1.290	.142	.909	.000	.169	.682
درصد مالکیت سهامداران نهادی	820	-1.412	.286	.125	.894	.000	1.251	.264

منبع: یافته های پژوهش

#### ۶-۲-۲- هم خطی

در این قسمت ضریب همبستگی میان متغیرهای مستقل نسبت های مالی با یکدیگر محاسبه شده و متغیرهایی که دارای همبستگی معنادار هستند شناسایی شده اند. در ادامه در صورت لزوم به حذف برخی متغیرها پرداخته شده است.

فروش در هر روز «Z آلتمن»، «سرمایه در گردش به کل دارایی»، «فروش به موجودی» و «بهای تمام شده کالای فروش رفته به موجودی کالا»، همسانی واریانس های دو گروه برقرار نیست. در مورد سایر متغیرها مقدار (p-value) حاصل شده، بیش از ۰/۰۵ درصد است، از اینرو، وجود همسان بودن واریانس پذیرفته می شود.

#### ۶-۲-۳- همسانی واریانسها

نتایج حاصل در جدول شماره (۲)، نشان می دهد که در مورد نسبت های مالی «وجه نقد عملیاتی به سود خالص»، «کل دارایی ها به فروش»، «سود خالص به حقوق صاحبان سهام»، «فروش به سود خالص»، «سود هر سهم»، «سود عملیاتی به فروش»، «سود انباشته به کل داراییها»، «دوره پرداخت بدهی»، «دوره وصول مطالبات»، «حسابهای دریافتی به متوسط

#### ۶-۲-۴- محاسبه ضریب همبستگی

نتایج حاصل از محاسبه ضریب همبستگی متغیرهای نسبت های مالی و احتمال تقلب نشان می دهد که متغیرهای «وجه نقد عملیاتی به سود خالص»، «کل دارایی ها به فروش»، «کل بدهیها به کل داراییها»، «دوره گردش موجودی»، «دوره پرداخت بدهی»، «دوره وصول مطالبات»، «چرخه تبدیل وجه

جدول شماره (۴): خلاصه نتایج مقایسه‌ای میزان اهمیت هر یک از متغیرهای وارد شده در مدل شبکه عصبی

متغیرهای مستقل مهم		
ضریب اهمیت نرمال شده	ضریب اهمیت	شرح متغیرهای مدل پژوهش
41.5%	0.055	وجه نقد عملیاتی به سود خالص
57.6%	0.076	کل دارایی‌ها به فروش
38.3%	0.050	سود خالص به حقوق صاحبان سهام
38.1%	0.050	سود هر سهم
32.5%	0.043	سود ناخالص به فروش
57.2%	0.075	سود عملیاتی به فروش
24.3%	0.032	سود انباشته به کل داراییها
25.9%	0.034	کل بدهیها به کل داراییها
61.3%	0.081	داراییهای جاری به بدهی های جاری
21.0%	0.028	دوره گردش موجودی
31.4%	0.041	دوره پرداخت بدهی
85.6%	0.113	دوره وصول مطالبات
29.7%	0.039	چرخه تبدیل وجه نقد
100.0%	0.132	Z آلتمن
43.8%	0.058	سرمایه در گردش به کل دارایی
21.6%	0.028	فروش به حسابهای دریافتی
28.7%	0.038	فروش به موجودی
13.6%	0.018	سهامداران عمده (بالای ۵ درصد)
8.2%	0.011	درصد مالکیت سهامداران نهادی

منبع: یافته‌های تحقیق

۲-۵-۲-۶- برآزش مدل درخت تصمیم

متغیرهای ورودی شامل متغیرهای مستقل معنادار براساس رگرسیون لجستیک بر اساس پژوهش احمدی و همکاران (۱۳۹۸) بوده و عبارتند از کل دارایی‌ها به فروش، سود عملیاتی به فروش، دوره گردش موجودی، چرخه تبدیل وجه نقد. براساس درخت حاصل در Node صفر مطابق انتظار ۵۰ درصد شرکتها متقلب و ۵۰ درصد غیرمتقلب هستند. اولین عامل شناسایی شده در این درخت نسبت مالی دوره پرداخت بدهی است. همچنین در رده ی دوم نسبت دارایی‌های جاری به بدهی‌های جاری و سرمایه در گردش به کل دارایی و همچنین در رده سوم نسبت کل دارایی‌ها به فروش، در درخت قرار گرفته است. به طور کلی نتایج حاصل از آزمون مربوطه در جدول شماره (۵) نشان می‌دهد که پیش بینی‌های صحیح انجام شده توسط این درخت تصمیم حدود ۶۵ درصد است.

نقد»، «بهای تمام شده کالای فروش رفته به فروش» و «فروش به حسابهای دریافتی» همبستگی مثبت با احتمال تقلب دارند و متغیرهای «سود خالص به حقوق صاحبان سهام»، «سود قبل از بهره و مالیات به کل دارایی»، «سود هر سهم»، «سود ناخالص به فروش»، «سود عملیاتی به فروش»، «سود انباشته به کل داراییها»، «دارایی‌های جاری به بدهی های جاری»، «Z آلتمن»، «سرمایه در گردش به کل دارایی»، «فروش به موجودی»، «سهامداران عمده (بالای ۵ درصد)»، «درصد مالکیت سهامداران نهادی»، رابطه‌ی معکوس با احتمال وقوع تقلب را دارا می‌باشند.

۵-۲-۶- نتایج حاصل از آزمون فرضیه

۱-۵-۲-۶- برآزش مدل شبکه عصبی

در گروه آموزش حدود ۷۰ درصد مشاهدات که معادل ۵۷۵ مورد می‌باشد قرار گرفته است، ۳۰ درصد مابقی داده‌ها یعنی تعداد ۲۴۵ مورد نیز در گروه آزمایش قرار دارند. تابع فعالسازی در لایه پنهان Hyperbolic tangent و در لایه خروجی Softmax می‌باشد. پس از آموزش و برآزش مدل شبکه عصبی میزان اعتبار نتایج با توجه به فراوانی پیش بینی‌های صحیح در هر یک از گروه داده‌های آزمایش بررسی شده است. نتایج حاصل در جدول زیر نشان داده شده است. در داده‌های بخش آزمایش شبکه، نرخ دقت آزمون برابر ۶۹.۴ درصد بوده است. این مقدار برای داده‌های گروه آموزش ۶۳ درصد حاصل شده است.

جدول شماره (۳): برآزش مدل شبکه عصبی

طبقه بندی				
نمونه	مشاهده شده	پیش بینی		
		۰	۱	نرخ دقت
آموزش	۰	۱۷۹	۱۱۰	۹۰.۶۱%
	۱	۱۰۳	۱۸۳	۰.۶۴%
	درصد کلی	۰.۴۹%	۰.۵۱%	۰.۶۳%
آزمایش	۰	۸۶	۳۵	۱.۷۱%
	۱	۴۰	۸۴	۷.۶۷%
	درصد کلی	۴.۵۱%	۶.۴۸%	۴.۶۹%

میزان اهمیت هر یک از متغیرهای وارد شده در مدل شبکه عصبی نیز در مقایسه با یکدیگر در جدول شماره (۳) آمده است. بیشترین میزان اهمیت در مدل مربوط به متغیرهای کل دارایی‌ها به فروش، سود عملیاتی به فروش، دارایی‌های جاری به بدهی‌های جاری، دوره وصول مطالبات و Z آلتمن می‌باشد.

66.7%	64	32	1	آزمایش
64.4%	60.9%	68%	درصد کلی	
متغیر وابسته: تقلب				

در این روش اعتبار نتایج با استفاده از گروه داده های آزمایش و آموزش قابل محاسبه می باشد. نتایج حاصل برای گروه آموزش نشان می دهد که نرخ دقت آزمون برابر ۷۸.۴ درصد بوده است. همچنین در گروه آزمایش نیز دقت پیش بینی های صورت گرفته برابر ۶۴.۴ درصد می باشد. بنابراین میتوان گفت این روش با نرخ دقت ۶۴.۴ درصد، توانایی خوبی را در پیش بینی تقلب صورتهای مالی دارد.

#### ۶-۲-۵-۴- ماشین بردار پشتیبان

پس از آموزش و برازش مدل میزان اعتبار نتایج با توجه به فراوانی پیش بینی های صحیح در گروه داده های آزمایش بررسی و در جدول ... گزارش شده است. نتایج حاصل در جدول زیر نشان داده می دهد که به طور کلی در داده های مورد استفاده برای آزمایش ۷۸ درصد مقادیر درست پیش بینی شده است. همچنین در گروه آموزش نیز ۷۳.۲ درصد پیش بینی صحیح انجام شده است. بنابراین میتوان گفت این روش با نرخ دقت ۷۸ درصد، توانایی خوبی را در پیش بینی تقلب صورتهای مالی دارد.

جدول شماره (۷): خلاصه نتایج آزمون ماشین بردار پشتیبان

نمونه	مشاهده	پیش بینی		
		۰	۱	درصد پیش بینی
آموزش	0	۲۳۸	۹۸	%۷۰.۸
	1	۶۴	۲۱۵	%۷۷
	درصد کلی	%۷۸.۸	%۶۸.۶	%۷۳.۲
آزمون	0	۹۰	۲۷	%۷۶.۹
	1	۱۸	۷۰	%۸۱.۴
	درصد کلی	%۸۳.۳	%۷۲.۲	%۷۸

#### ۷- بحث و نتیجه گیری

داده کاوی بطور همزمان از چندین رشته علمی بهره می برد نظیر؛ تکنولوژی پایگاه داده، هوش مصنوعی، یادگیری ماشین، شبکه های عصبی، آمار، شناسایی الگو، سیستم های مبتنی بر دانش<sup>۱</sup>، حصول دانش<sup>۲</sup>، بازیابی اطلاعات<sup>۳</sup>، محاسبات سرعت بالا<sup>۴</sup>

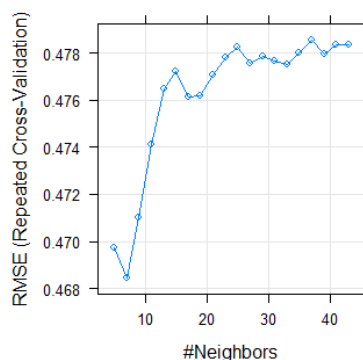
جدول شماره (۵): خلاصه نتایج آزمون درخت تصمیم

دسته بندی			
مشاهده شده	پیش بینی		
	۰	۱	درصد صحیح
۰	۲۷۲	۱۳۸	%۶۶.۳
۱	۱۴۶	۲۶۴	%۶۴.۴
درصد کلی	%۵۱.۰	%۴۹.۰	%۶۵.۴
Growing Method: CHAID			

بنابراین می توان گفت درخت تصمیم با نرخ دقت ۶۵.۴ درصد توانایی خوبی را در پیش بینی تقلب صورتهای مالی دارد.

#### ۶-۲-۵-۳- برازش مدل نزدیک ترین همسایگی

داده های آموزش با دادن مقادیر مختلف برای K برازش داده شده و ضریب دقت مدل برای هر مقدار K محاسبه می شود. نمودار شماره (۲) تغییرات RMSE را براساس مقادیر مختلف K نشان می دهد. همانطور که ملاحظه می شود K=۷ بهترین حالت ممکن را حاصل می کند.



نمودار شماره (۱): نمایش برازش مدل نزدیکترین همسایگی

جدول شماره (۶): خلاصه نتایج آزمون برازش مدل نزدیک

ترین همسایگی

دسته بندی				
نمونه	مشاهده	پیش بینی		درصد صحیح
		0	1	
آموزش	0	۲۳۱	۷۰	23.2%
	1	۶۳	۲۵۱	79.9%
	درصد کلی	78.6%	78.2%	78.4%
	0	68	41	62.4%

<sup>3</sup> Information retrieval

<sup>4</sup> High-performance computing

<sup>1</sup> Knowledge-based system

<sup>2</sup> Knowledge-acquisition

و بازنمایی بصری داده<sup>۱</sup>. واژه های «داده کاوی» و «کشف دانش در پایگاه داده»<sup>۲</sup> اغلب به صورت مترادف یکدیگر مورد استفاده قرار می گیرند. کشف دانش در پایگاه داده فرایند شناسایی درست، ساده، مفید، و نهایتا الگوها و مدل‌های قابل فهم در داده ها می باشد. داده کاوی، مرحله ای از فرایند کشف دانش می باشد و شامل الگوریتم‌های مخصوص داده کاوی است، بطوریکه، تحت محدودیتهای مؤثر محاسباتی قابل قبول، الگوها و یا مدلها را در داده کشف می کند

هدف اصلی پژوهش حاضر بررسی اثربخشی تکنیکهای داده کاوی شبکه عصبی، درخت تصمیم، نزدیک ترین همسایگی و ماشین بردار پشتیبان در پیش بینی تقلب صورت‌های مالی است. به منظور تحقق این هدف از یک نمونه شرکتهای با احتمال تقلب و شرکتهای غیرمتقلب طبق معیارهای تبیین شده قبلی برای دسته بندی شرکتهای؛ استفاده شد. یک مجموعه از ۳۵ نسبت مالی به عنوان پیش بینی کنندگان بالقوه تقلب صورتهای مالی انتخاب شدند. این متغیرها در تحقیقات قبلی با اهمیت ظاهر شدند و از صورتهای مالی انتشار یافته استخراج می شوند. در نهایت ۱۹ متغیر انتخاب شدند. نتایج حاصل از آزمونها که تایید کننده نتایج آزمونهای قبلی است نشان میدهد که صورتهای مالی می تواند جهت پیش بینی تقلب استفاده شود و اعتبار دهندگان، حسابرسان و سایر ذینفعان را یاری رساند. بدین ترتیب که تکنیک شبکه عصبی ۶۹/۴ درصد، درخت تصمیم ۶۵/۴ درصد، نزدیکترین همسایگی ۶۴/۴ درصد و ماشین بردار پشتیبان ۷۸ درصد پیش بینی صحیح داشته اند. بنابراین از نظر کارایی، تکنیک ماشین بردار پشتیبان با ۷۸ درصد صحت طبقه بندی بهترین عملکرد را داشته است.

#### ۷-۱- مقایسه نتایج حاصل از فرضیه ها

(۱) دقت روش نزدیکترین همسایگی (۶۴/۴ درصد) با دقت روش رگرسیون لجستیک (۶۴/۶) تفاوت معناداری ندارد. به بیان دیگر، روش نزدیکترین همسایگی در طبقه بندی صحیح شرکتهای مشابه روش رگرسیون لجستیک است. اما دقت روش ماشین بردار پشتیبان (۷۸ درصد) بالاترین دقت و بالاتر از روش شبکه عصبی (۶۹/۴ درصد)، درخت تصمیم (۶۵/۴ درصد) و دو روش ذکر شده قبلی است. به بیان دیگر عملکرد روش ماشین بردار پشتیبان در طبقه بندی صحیح شرکتهای به شکل معناداری بهتر از عملکرد روشهای نزدیکترین همسایگی، رگرسیون لجستیک، شبکه عصبی و درخت

تصمیم است و روش رگرسیون لجستیک، در شناسایی شرکتهای متقلب به شکل معناداری ناموفقتر از روش ماشین بردار پشتیبان است.

(۲) نرخ خطای نزدیکترین همسایگی (۳۵/۶) با نرخ خطای رگرسیون لجستیک (۳۵/۲) تفاوت معناداری ندارد. به بیان دیگر روش نزدیکترین همسایگی در طبقه بندی اشتباه شرکتهای، مشابه عملکرد الگوریتم شبکه عصبی است. نرخ خطای ماشین بردار پشتیبان (۲۲ درصد) به شکل معناداری کمتر از نرخ خطای روشهای شبکه عصبی (۳۰/۶) و درخت تصمیم (۳۴/۶) و دو روش ذکر شده قبلی است. به بیان دیگر اشتباه در طبقه بندی شرکتهای در ماشین بردار پشتیبان به شکل معناداری کمتر از اشتباه در طبقه بندی شرکتهای با استفاده از روشهای نزدیکترین همسایگی، رگرسیون لجستیک، شبکه عصبی و درخت تصمیم است.

(۳) حساسیت روش نزدیکترین همسایگی (۶۶/۷ درصد) با حساسیت روش شبکه عصبی (۶۷/۶ درصد) تفاوت معناداری ندارد. به بیان دیگر عملکرد روش نزدیکترین همسایگی در شناسایی شرکتهای متقلب، مشابه روش شبکه عصبی است. معیار حساسیت ماشین بردار پشتیبان (۷۹/۵ درصد) به شکل معناداری بیشتر از معیار حساسیت روشهای رگرسیون لجستیک (۶۱/۴ درصد) و درخت تصمیم (۶۴/۴ درصد) و دو روش ذکر شده قبلی است. به بیان دیگر عملکرد ماشین بردار پشتیبان در شناسایی شرکتهای متقلب به شکل معناداری بهتر از عملکرد روشهای نزدیکترین همسایگی، رگرسیون لجستیک، شبکه عصبی و درخت تصمیم است و روش رگرسیون لجستیک، در شناسایی شرکتهای متقلب به شکل معناداری ناموفقتر از روش ماشین بردار پشتیبان است.

(۴) معیار ویژگی ماشین بردار پشتیبان (۷۶/۹) به شکل معناداری بیشتر از ویژگی روشهای رگرسیون لجستیک (۶۸ درصد)، درخت تصمیم (۶۶/۳ درصد)، نزدیکترین همسایگی (۶۲/۴ درصد) و شبکه عصبی (۷۱/۱ درصد) است. به بیان دیگر عملکرد ماشین بردار پشتیبان در شناسایی شرکتهای غیرمتقلب به شکل معناداری بهتر از عملکرد روشهای نزدیکترین همسایگی، رگرسیون لجستیک، شبکه عصبی و درخت تصمیم است. افزون

<sup>۱</sup> Data visualization

<sup>۲</sup> Knowledge Discovery in Database

- مالی، اولین همایش بین المللی اقتصاد سنجی، روشها و کاربردها، سندج، دانشگاه آزاد اسلامی واحد سندج.
- \* فرقاندوست حقیقی، کامبیز؛ برواری، فرید، ۱۳۸۸، بررسی کاربرد روش های تحلیلی در ارزیابی ریسک تحریف صورت های مالی (تقلب مدیریت)، دانش و پژوهش حسابداری، شماره شانزده
- \* خواجهی، شکراله؛ ابراهیمی، مهرداد (۱۳۹۶) مدل سازی متغیرهای اثرگذار برای کشف تقلب در صورتهای مالی با استفاده از تکنیکهای داده کاوی، فصلنامه حسابداری مالی، شماره سی و سه
- \* سجادی، سید حسین؛ کاظمی، توحید (۱۳۹۵) الگوی جامع گزارشگری مالی متقلبانه در ایران به روش نظریه پردازی زمینه بنیان، پژوهشهای تجربی حسابداری، سال ششم، شماره ۲۱، پاییز ۱۳۹۵، صص ۱۸۵-۲۰۴
- \* Abbasi A., Albrecht C., Vance A., & Hansen J. (2012). Metafraud: A meta-learning framework for detecting financial fraud. *MIS Quart Manage Inf Syst MIS Quarterly: Management Information Systems*, 36(4), 1293-1327.
- \* Albrecht, W. S. (2012). *Fraud examination*. Mason, Ohio: South-Western Cengage Learning
- \* ASHRAF AKL ELSAYED . (2017). *Predictability of Financial Statements Fraud-Risk* . Northcentral University.
- \* Association of Certified Fraud Examiners (ACFE). (2016). *Report to the nations on occupational fraud and abuse*. Austin, TX: The association of certified fraud examiners, Inc.
- \* Hays, J. B. (2014). *An investigation of the motivation of management accountants to report fraudulent accounting activity: Applying the theory of planned behavior*. *Dissertation Abstracts International Section A*, 75.
- \* Kravitz, R. H. (2012). Auditors' responsibility for detecting fraud. *CPA Journal*, 82(6), 24-30
- \* Königsgruber, R. (2012). Capital Allocation Effects of Financial Reporting Regulation and Enforcement. *European Accounting Review*, 21(2), 283-296. doi:10.1080/09638180.2011.558294
- \* Goel, S., & Gangolly, J. (2012). Beyond the numbers: Mining the annual reports for hidden cues indicative of financial statement fraud. *Intelligent Systems in Accounting, Finance & Management*, 19(2), 75-89. doi:10.1002/isaf.1326
- \* Mangala, D., & Kumari, P. (2015). Corporate fraud prevention and detection: Revisiting the literature. *Journal of Commerce and Accounting Research*, 4(1).
- \* Ruankaew, T. (2016). Beyond the fraud diamond. *International Journal Of Busines Management & Economic Research*, 7(1), 474-476.
- \* Wuerges, A. E., & Borba, J. A. (2014). Accounting Fraud: an estimation of detection probability. *Revista Brasileira De Gestão De Negócios*, 16(52), 466-483. doi:10.7819/rbgn.v16i52.1555

بر این ویژگی نزدیکترین همسایگی به شکل معناداری کمتر از ویژگی روش ماشین بردار پشتیبان است. به بیان دیگر نزدیکترین همسایگی در شناسایی شرکتهای غیرمتقلب به شکل معناداری ناموفقتر از روش ماشین بردار پشتیبان است.

(۵) این پژوهش میتواند به عنوان زیربنایی برای پژوهشهای آتی استفاده شود و سایر پژوهشگران با استفاده از سایر تکنیک های داده کاوی و یا در نظر قراردادن شرکتهایی مانند شرکتهای سرمایه گذاری، بانکها و بیمه که در نمونه تحقیق قرار نگرفته اند، به پژوهش در این موضوع ادامه دهند. همچنین مطالعه یک صنعت خاص می تواند نتایج ویژه خود را داشته باشد. با ارتقای روشهای مورد استفاده در داده کاوی میتوان داده کاوی را به ابزاری مفیدتر در فرآیند تصمیم گیری مالی تبدیل کرد. کاربرد روشهای داده کاوی بر روی نسبتهای مالی و نیز دیگر اطلاعات، می تواند به حسابرسان در کشف تقلب کمک کند، به طوری که آنان می توانند از نتایج این تحلیلها به عنوان یک علامت اولیه هشداردهنده نسبت به احتمال تقلب صورتهای مالی استفاده کنند. کشف نشانگرهای تقلب در صورتهای مالی، اثری بااهمیت بر تعیین تقلب صورتهای مالی دارد. به طور کل نتایج این پژوهش با نتایج پژوهشهای اشرف آکل الساید (۲۰۱۷)، خواجهی (۱۳۹۵) و رهنمای رودپشتی (۱۳۹۱) مبنی بر سودمندی شیوه های داده کاوی برای پیش بینی تقلب در صورتهای مالی مطابقت دارد.

#### فهرست منابع

- \* رهنمای رود پشته، فریدون (۱۳۹۱)، داده کاوی و کشف تقلب های مالی، فصلنامه علمی و پژوهشی دانش حسابداری و حسابرسی مدیریت، شماره سوم
- \* جمالی، زهرا؛ برزگری خانقاه، جمال؛ عارف منش، زهره؛ انصاری سامانی، حبیب (۱۳۹۵) بررسی رابطه مکانیزمهای حاکمیت شرکتی و کیفیت حسابرسی بر وقوع تقلب در صورتهای مالی شرکت های پذیرفته شده در بورس اوراق بهادار تهران، یزد، دانشگاه یزد
- \* اعتمادی، حسین؛ زلفی، حسن؛ (۱۳۹۲) کاربرد رگرسیون لجستیک در شناسایی گزارشگری مالی متقلبانه، فصلنامه دانش حسابرسی، شماره پنجاه و یک
- \* مهمان، کیهان؛ غلامرضا کردستانی و ابوالفضل ترابی، ۱۳۹۱، ارائه مدلی برای پیش بینی خطر بروز تقلب در گزارشگری



*Accounting Knowledge & Management Auditing*

*Vol. 13/ No. 52/ Winter 2024*

## **Data Mining Techniques and Forecasting Financial Statement Fraud**

**Seyed Jalal Ahmadi**

PhD student in Accounting, Semnan Branch, Islamic Azad University, Semnan, Iran

**Khosro faghani Makrani**

Associate Professor, Department of Accounting, semnan Branch, Islamic Azad University, semnan, Iran.

**Naghi Fazeli**

Assistant Professor , Department of Accounting, Semnan Branch, Islamic Azad University, Semnan, Iran

### **Abstract**

The purpose of this study is to compare neural network, decision tree, nearest neighbor and support vector machine data mining techniques in predicting fraudulent and non-fraudulent financial statements. The research method is descriptive-applied and time domain from 2008 to 2018. In this study, financial ratios for two fraudulent and non-fraudulent samples and data mining methods were analyzed. Statistical hypotheses of normality, homogeneity and linearity test for financial ratios of fraudulent and non-fraudulent samples were tested. The normality hypothesis was tested using Kolmogorov-Smirnov test and Shapiro Wilk test. Then Pearson correlation coefficient for the existence of the model for financial ratios and elimination of correlated independent variables was tested. Next, data mining methods are used to test them in predicting financial statement fraud and distinguishing fraudulent and non-fraudulent financial statements. In general, the results show that data mining methods are effective in differentiating fraudulent and non-fraudulent financial statements. The neural network method had a correct prediction of 69.4%, decision tree 65.4%, nearest neighbor 64.4% and support vector machine 78%.

**Keywords:** Fraud, Data mining, Financial ratios