

چکیده

در طول سالیان اخیر، حجم مخازن داده‌های الکترونیکی ایجاد شده توسط بانک‌ها، شرکت‌ها و سایر موسسات تجاری، همواره در حال افزایش بوده است. بیت‌های با ارزش اطلاعات در این مخازن داده، جای گرفته‌اند، به گونه‌ای که اندازه و حجم بسیار زیاد این منابع، تجزیه و تحلیل آن‌ها را برای انسان به منظور خلق اطلاعات یا الگوهای مفید و مناسب در جهت فرآیند تصمیم‌گیری، غیر ممکن ساخته است. فنونی نظری داده‌کاوی می‌توانند به طور خودکار، اطلاعات را از میان این مقدادر حجمیم داده استخراج نمایند. داده‌کاوی، تاثیراتی عمیق بر شیوه‌های کسب و کار و مدیریت دانش در سال‌های اخیر داشته است. هوشمندی کسب و کار، مشهورترین کاربرد فنون داده‌کاوی است. قابلیت‌های طبقه‌بندی و پیش‌بینی ابزارهای داده‌کاوی، موجب کاربرد این ابزارها به منظور مقاصد پیش‌بینی و رشکستگی، وضعیت تداوم فعالیت، پیش‌بینی درمانگی مالی، تشخیص و شناسایی تقلب مدیریت، هم‌چنین تخمین ریسک اعتباری و عملکرد واحد تجاری شده و بدین‌ترتیب داده‌کاوی را تبدیل به موضوعی با درجه اهمیت زیاد، در امور مالی و حسابداری نموده است. هدف این مقاله بررسی اجمالی این فناوری جدید و تحقیقات مرتبط با آن، به منظور نمایان ساختن کارکردهای داده‌کاوی در امور مالی و حسابداری و فرستهای تحقیقاتی آن می‌باشد.

واژه‌های کلیدی: حسابداری، امور مالی، روش‌های داده‌کاوی، ابزارهای داده‌کاوی، داده‌کاوی

تاریخ دریافت: ۱۳۹۰/۴/۱۰

تاریخ پذیرش: ۱۳۹۰/۵/۱۲

فناوری داده‌کاوی؛ رویکردی نوین در حوزه مالی

دکتر محسن دستگیر / استاد دانشگاه شهید چمران اهواز

مرتضی شفیعی سردشت / دانشجوی کارشناسی ارشد دانشگاه آزاد مشهد

مقدمة

داده‌کاوی، شاخه‌ای در خشان در علوم کامپیوتری است که در اواخر دهه ۸۰ با بکارگیری مفاهیم و روش‌های مرتبط با هوش مصنوعی^۱، شناسایی الگو^۲، سیستم‌های پایگاه داده^۳ و علم آمار^۴، پا به عرصه گذاشت.

در طول سال‌های اخیر، به تدریج با این واقعیت رشد یافته‌ایم که حجم عظیمی از داده‌ها وجود دارند که کامپیوترها، شبکه‌ها و در حقیقت تمام زندگی ما را فرا گرفته‌اند. سازمان‌های دولتی، مؤسسه‌های علمی و تجاری سرمایه هنگفتی را برای گردآوری و ذخیره این داده‌ها اختصاص داده‌اند. حجم اطلاعات جمع‌آوری و ذخیره شده رشد سریعی به خود گرفته است. برخی از تحقیقات نشان می‌دهند که میزان داده‌های جهانی تقریباً هر سال دو برابر می‌شوند، در همان حال نیز هزینه ذخیره‌سازی این داده‌ها به طور قابل توجهی از یک دلار برای هر مگابایت به چند «پن» کاهش یافته است. به گونه‌ای مشابه، توان محاسباتی در هر ۱۸ الی ۲۴ ماه دو برابر شده، در حالی که هزینه مرتبط با این توان محاسباتی در حال کاهش می‌باشد (Seifert, ۲۰۰۴).

با وجود این، فقط مقدار کمی از این داده‌ها مورد استفاده قرار می‌گیرند، زیرا در بسیاری از موارد، حجم داده‌های لازم برای سازماندهی، بسیار بالا بوده یا ساختار آنها به منظور تجزیه و تحلیل مؤثر و کارا به شدت پیچیده می‌باشد.

این مشکل چطور به وجود آمد؟ در پاسخ باید گفت اغلب تلاش‌ها برای ایجاد مجموعه‌ای از داده، بر روی موضوعاتی از قبیل کارآمدی و ذخیره‌سازی متمرکز می‌شوند و در حقیقت، هیچگونه طرح یا برنامه‌ای برای این وجود ندارد که در نهایت، داده‌ها چگونه استفاده و تجزیه و تحلیل شوند (علیخانزاده، ۱۳۸۵).

با ورود به عصر اطلاع رسانی دیجیتال، مسئله انفجار داده‌ها، اطلاعات و ضرورت نگهداری آنها هر روز بیش از پیش احساس می‌شود. تحلیل داده‌ها می‌تواند دانش مضاعفی در مورد تجارتی خاص ایجاد کند. این مهم با افزایش شفافیت داده‌ها به منظور استخراج دانش، دست یافتنی خواهد بود. مجموعه داده‌ها در شکل خام و پرداخت نشده خود از ارزش کمی برخوردار هستند، آنچه که در حقیقت ارزشمند است دانش و شناختی است که می‌تواند از داده‌ها حاصل شود و مورد استفاده قرار گیرد. به لحاظ نظری، داده‌های حجمی و بزرگ می‌توانند به نتایج قوی‌تری منجر شوند. جامعه تجارتی به خوبی از سریزی و حجم زیاد اطلاعات امروز آگاه است. نتایج بررسی‌های به عمل آمده نشان می‌دهد که:

داده‌کاوی و داده‌های مالی

در دنیای تجارت نیز داده‌های مالی به عنوان سرمایه راهبرد مطرح هستند. داده‌های مالی توسط موسساتی مثل بانک‌ها، بورس اوراق بهادار، سازمان‌های مالیاتی، پایگاه‌های داده ویژه حسابرسان و حسابداران و غیره جمع آوری و نگهداری می‌شوند. روش‌های داده‌کاوی در داده‌های مالی، می‌تواند در حل مشکلات طبقه‌بندی و پیش‌بینی و تسهیل فرآیند تصمیم‌گیری به کار رود. نمونه‌هایی از مسائل طبقه‌بندی مالی^{۱۴} شامل ورشکستگی شرکت‌ها^{۱۵}، تخمین ریسک اعتباری^{۱۶}، گزارش تداوم

- الف- ۶۱ درصد از مدیران معتقد به سریزی اطلاعات در شرکت‌ها می‌باشند.
 - ب- ۸۰ درصد اعتقاد به بدتر شدن وضعیت فوق دارند.
 - ج- بیش از ۵۰ درصد مدیران از داده‌ها در مراحل تصمیم‌گیری فعلی به واسطه سریزی اطلاعات صرف‌نظر می‌کنند.
 - د- ۸۴ درصد از مدیران، این اطلاعات را برای آینده ذخیره می‌کنند و از آنها برای تحلیل‌های جاری استفاده نمی‌کنند.
 - ه- ۶۰ درصد عقیده دارند که هزینه جمع‌آوری و ذخیره اطلاعات بیشتر از ارزش خود اطلاعات می‌باشد (همان).
- در این شرایط، داده‌کاوی^۱ ابزاری سودمند است که در اختیار تمام شرکت‌ها می‌باشد. داده‌کاوی و کشف دانش، ابزارهایی هوشمند هستند که به جمع‌آوری و پردازش داده‌ها و بهره‌مندی از آنها، منجر می‌شوند. در واقع، داده‌کاوی شامل استفاده از ابزارهای تحلیلی پیچیده به منظور کشف الگوهای معترض و ناشناخته وابستگی‌های موجود در مجموعه داده‌های حجمی می‌باشد. نتیجه این که، فناوری داده‌کاوی چیزی بیش از جمع‌آوری و مدیریت داده بوده و شامل تجزیه و تحلیل و پیش‌بینی نیز می‌باشد (۲۰۰۴، Seifert).

برخی از صاحب‌نظران براساس مفاهیم داده‌کاوی، معتقدند این فناوری تنها مرحله‌ای از فرآیندی فرآگیرتر به نام کشف دانش در پایگاه داده^۲ می‌باشد. مراحل دیگر در این فرآیند فرآگیر شامل پاک‌سازی^۳، یکپارچه‌سازی^۴، انتخاب داده‌ها^۵، تبدیل^۶ و ارزیابی الگو^۷ است. فناوری داده‌کاوی میان بسیاری از حوزه‌های فنی شامل پایگاه داده، آمار، یادگیری ماشین^۸ و رابط انسان و کامپیوتر^۹ ارتباط برقرار می‌سازد (al pechenizkiy et.al ۲۰۰۳).

فعالیت^{۱۷}، در ماندگی مالی^{۱۸} و پیش‌بینی عملکرد واحد تجاری^{۱۹} می‌باشد.^{۲۰} بدین ترتیب دامنه اهمیت داده کاوی در امور مالی و حسابداری می‌تواند طیفی گسترده از کشف تقلب تا افزایش سودآوری واحد تجاری باشد (Rumsey and et.al, ۲۰۰۴).

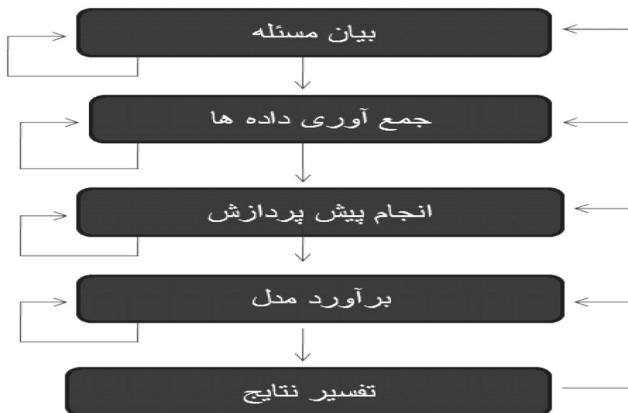
اهمیت داده کاوی توسط بسیاری از سازمان‌های حرفه‌ای تشخیص داده شده است. انجمن حسابداران رسمی آمریکا، داده کاوی را به عنوان یکی از ۱۰ فناوری برتر برای آینده معرفی کرده است. همچنین انجمن حسابرسان داخلی آمریکا نیز این فناوری را در فهرست یکی از چهار اولویت تحقیقاتی خود گنجانده است (kirkods and et.al, ۲۰۰۴).

کاربرد احتمالی داده کاوی در حسابداری مالی منجر به استفاده از صورت‌های مالی به جهت تعیین سودآوری، تطابق نسبت‌ها با میانگین صنعت و بررسی اثربخشی یک کسب و کار می‌گردد. در حسابداری مدیریت نیز داده‌های واحدهای تجاری به منظور تعیین چگونگی عملیات روزمره واحد تجاری، تجزیه و تحلیل بهای تمام شده و سودآوری بخش‌های مختلف، مورد استفاده قرار می‌گیرد. در حسابرسی، ابزارهای داده کاوی برای تجزیه و تحلیل داده‌های مرتبط با احتمال تقلب مدیریت می‌توانند به کار گرفته شوند. بدین ترتیب داده کاوی در تمام جنبه‌های حسابداری سودمند خواهد بود، زیرا حسابداری عبارت از جمع‌آوری، تجزیه و تحلیل واستفاده از اطلاعات مالی و غیر مالی به منظور پردازش برای اشخاص ذی‌نفع است (Rumsey and et.al, ۲۰۰۴). شرکت‌های تجاری می‌توانند از داده کاوی به منظور شناسایی فرصت‌های مالی و عملیاتی قابل توجه برای اطمینان از کارایی فرآیند زنجیره تامین^{۲۱} بهره‌مند شوند.

مفاهیم داده کاوی

الف- مراحل (فرآیند) داده کاوی

با نگاهی اجمالی به مراحل داده کاوی، می‌توان نحوه انجام آن را به صورت نمودار شماره یک بیان کرد:



نمودار شماره یک : مراحل داده‌کاوی

مراحل داده‌کاوی را می‌توان به گونه‌ای جزئی تر به صورت زیر ابراز کرد:

- ۱- **شناسایی هدف:** در این مرحله مشخص می‌شود کاربر چه چیزی را می‌خواهد و تا چه سطحی از اطلاعات را در نظر دارد که از پایگاه داده، اخذ نماید.
- ۲- **انتخاب داده‌ها:** در این مرحله باید داده‌ها بر مبنای معیارهای مشخص انتخاب گردد.
- ۳- **آماده‌سازی داده‌ها:** شکل قابل استفاده داده و شناسایی متغیرهای زائد وظیفه این مرحله از فرآیند خواهد بود.
- ۴- **ارزیابی داده‌ها:** چارچوب کلی این مرحله، معیارهایی از قبیل نوع توزیع داده‌ها، ویژگی‌ها و ساختار پایگاه داده و شرایط کلی داده‌ها... می‌باشد.
- ۵- **قالب‌بندی پاسخ:** خروجی این بخش، ارایه فرمت به شکل تصویر، نمودار، شبکه عصبی ... است.
- ۶- **انتخاب ابزار:** در این مرحله ابزارهای مناسب برای داده‌کاوی انتخاب می‌گردد.
- ۷- **الگوسازی:** فرآیند داده‌کاوی به صورت اصلی از این مرحله آغاز می‌گردد که شامل جستجوی الگوها در مجموعه داده، طبقه‌بندی و ارزشیابی داده‌ها... می‌باشد.
- ۸- **اعتبارسازی یافته‌ها:** این مرحله، شامل آزمون کردن الگوها است.
- ۹- **ارایه نتایج:** نتیجه این بخش، گزارش نهایی برای کاربر است.

۱۰- استفاده از نتایج: هدف اصلی داده کاوی استفاده از نتایج کشف شده برای تصمیم‌گیری، سیاست‌گذاری و پیش‌بینی به منظور ایجاد یک موقعیت بهتر و جدید می‌باشد (نجات و اکبری، ۱۳۸۷).

ب: ابزارهای داده کاوی^{۲۲}

ابزارهای داده کاوی نرم‌افزارهایی هستند که به کاربران اجازه استخراج اطلاعات از داده‌ها را می‌دهند. این ابزارها توانایی گردآوری داده‌ها و کاربرد آن‌ها به منظور تصمیم‌گیری در خصوص مصرف کننده خاص یا گروهی از مصرف‌کنندگان را، به سازمان‌ها و افراد می‌دهند. استخراج دستی داده‌ها از صدها سال پیش وجود داشته است. این در حالی است که ماشینی شدن فرآیند داده کاوی از هنگام ورود کامپیوتر شایع شده است. هدف نهایی این ابزارها ظاهر ساختن الگوهای پنهان^{۲۳} می‌باشد.

به هر حال ابزارهای داده کاوی می‌توانند شامل نرم‌افزارهای زیر باشد:

- ۱- Data to knowledge^{۲۴}
- ۲- sas^{۲۵}
- ۳- clementine^{۲۶}
- ۴- Intelligent-miner^{۲۷}
- ۵- Insightful miner^{۲۸}
- ۶- weka
- ۷- microsoft excel
- ۸- ACL^{۲۹}

ج- روش‌های داده کاوی^{۳۰}

اهداف داده کاوی شامل پیش‌بینی^{۳۱} و توصیف^{۳۲} یا ترکیبی از آنهاست. هدف "پیش‌بینی" تمرکز بر روی دقیقت در توانایی پیش‌بینی بوده و "توصیف" بر درک فرآیند تولید داده‌ها تمرکز دارد. در پیش‌بینی تازمانی که مدل قدرت پیش‌بینی دارد، کاربر توجهی به این ندارد که مدل انعکاس دهنده واقعیت است؛ مثلاً مدلی که شاخص‌های مالی را به شکل غیر خطی ترکیب می‌کند تا نرخ تبادل ارز را پیش‌بینی نماید. از سویی دیگر مدل توصیفی به مثابه واقعیت تفسیر می‌شود. برای مثال مدلی که متغیرهای اقتصادی و جمعیتی را به پیشرفت‌های آموزشی مرتبط می‌سازد، به عنوان مبنایی برای

توصیه‌های سیاست اجتماعی توصیف می‌شود. در عمل، اغلب کاربردهای اکتشاف دانش به درجه‌ای از هر دو مدل سازی توصیفی و پیش‌بینی نیاز دارند.

به هر ترتیب، اهداف داده‌کاوی با استفاده از روش‌های داده‌کاوی، محقق می‌شوند. اصطلاح روش‌های داده‌کاوی در واقع بیانگر جمع کثیری از الگوریتم‌ها^{۳۳} و فنون است که از علومی مانند آمار، یادگیری ماشین، پایگاه داده و تجسم‌سازی^{۳۴} استنتاج شده‌اند. برخی از این روش‌ها در مورد داده‌های مالی به کار رفته‌اند. روش‌های داده‌کاوی مشهوری که در این مقاله معرفی خواهند شد شامل شبکه‌های عصبی^{۳۵}، الگوریتم‌های ژنتیک^{۳۶} درختان تصمیم^{۳۷}، نظریه مجموعه اولیه^{۳۸}، استدلال مبتنی بر مورد^{۳۹} (استدلال موردگر) و برنامه‌نویسی ریاضی^{۴۰} هستند. بدیهی است که این فهرست از روش‌ها، فهرستی کامل نبوده و ترتیبی برای اولویت‌بندی در کاربرد این روش‌ها، پیشنهاد نمی‌شود (علیخانزاده، ۱۳۸۵).

شبکه‌های عصبی

شبکه‌های عصبی روشی محاسباتی هستند که بر پایه اتصال چندین واحد پردازشی، ایجاد می‌شوند. شبکه مورد نظر از تعداد دلخواهی سلول یا گره^{۴۱} یا نرون^{۴۲} تشکیل می‌شود که مجموعه ورودی را به مجموعه خروجی ارتباط می‌دهند. در واقع شبکه‌های عصبی مصنوعی، تا حد زیادی الهام یافته از سیستم‌های یادگیری طبیعی هستند که در آنها یک مجموعه پیچیده از نرون‌های به هم متصل، در کار یادگیری دخیل می‌باشند. به همراه هر اتصال یک ارزش عددی که وزن^{۴۳} نامیده می‌شود وجود دارد. هر نرون سیگنال‌هایی^{۴۴} از نرون‌های متصل دریافت می‌کند. اگر جمع سیگنال‌های ورودی از آستانه^{۴۵} تجاوز یابد در این صورت نرون از کار می‌افتد. نرون‌ها به صورت گروهی لایه‌بندی می‌شوند.

یک شبکه لایه‌ای، لاقل شامل یک لایه ورودی^{۴۶} و یک لایه خروجی^{۴۷} می‌باشد. بین این لایه‌های ورودی و خروجی ممکن است لایه‌های پنهان^{۴۸} نیز وجود داشته باشد. وقتی سیگنال یا پالسی به یک لایه ارسال می‌شود، این سیگنال از لایه بالایی شروع به فعالیت می‌کند و توسط نرون‌های آن لایه بررسی و اصلاح می‌شود. در حقیقت هر نرون قدرت سیگنال را بالا می‌برد و آن پالس را به لایه بعدی انتقال می‌دهد (صفایی، ۱۳۸۵).

أنواع مختلف شبکه‌های عصبی تعداد متناوبی از لایه‌ها را دارند. نقشه‌های خود سازمان یافته^{۴۹}

تنها یک لایه ورودی و خروجی دارند در حالی که شبکه عصبی پس انتشار خطای^{۵۰} یک یا چند لایه پنهان اضافی دارد. برای اینکه به شبکه عصبی موجود توانایی آموختن داده شود، بعد از اینکه سیگنال از لایه اول شبکه به لایه پایینی می‌رود، بایستی اطلاعات هر نرون که روی سیگنال اثر می‌گذارد، به روزآوری و اصلاح شود، این روند را به اصطلاح BP یا پس انتشار خطای^{۵۱} گویند.

در شبکه‌های پس انتشار خطای، خروجی با نتیجه دلخواه مقایسه می‌شود و خطاهای ایجاد شده در لایه‌های قبلی در شبکه عصبی با میزان سازی وزن‌های اتصال‌ها، مشخص می‌شوند. این فرآیند تا زمانی که یک نرخ اشتباه قابل قبول بدست آید، تکرار خواهد شد. شبکه‌های عصبی پس انتشار خطای، به منظور پیش‌بینی و طبقه‌بندی مسائل از انواع معتبر تلقی می‌شوند (KIRKDOSE and et al. ۲۰۰۴).

نقشه‌های خود سازمان یافته، یک روش تصویرسازی و خوش‌بندی از نوع آموزش بدون مری^{۵۲} می‌باشند. دو کاربرد معمول توپولوژی نقشه‌های خود سازمان یافته شامل لاتیس (شبکه) مستطیل^{۵۳} جائی که هر نرون دارای ۴ مجاور است و لاتیس شش ضلعی^{۵۴} جائی که هر نرون ۶ مجاور دارد، می‌باشد. از معایب شبکه‌های عصبی این است که تفسیر روشی که این شبکه‌ها به تصمیمات خود می‌رسند، بسیار دشوار است. از دیگر انتقادهای واردۀ بر شبکه‌های عصبی می‌توان به تعدادی از عوامل نظیر توپولوژی شبکه اشاره کرد که باید به طور تجربی، تعریف شوند. شبکه‌های عصبی مصنوعی برخلاف روش‌های آماری کلاسیک، مفروضاتی در خصوص ویژگی‌های توزیع داده‌ها و متغیرهای مستقل ورودی ندارند.

الگوریتم ژنتیک

این الگوریتم بر پایه فرضیه تکامل داروین بوده و کاربرد آن بر ژنتیک طبیعی استوار است. اصول اولیه الگوریتم ژنتیک توسط هلند^{۵۵} و همکاران در سال ۱۹۶۲ ارایه شده است. بر این اساس موجودات دارای سازگاری بیشتر با محیط، با میزان بالاتری زنده مانده و تولید مثل می‌کنند، در نتیجه شانس حضور آنها در نسل‌های بعدی بیشتر است (موسی زادگان و ذکری، ۱۳۸۷).

کاربرد الگوریتم ژنتیک، مستلزم بیان هر یک از پاسخ‌های یک مسئله خاص، در قالب یک رشته از اعداد می‌باشد. به هر یک از این رشته‌ها که در واقع نشان دهنده یکی از جواب‌های مسئله خواهد بود؛ کروموزم و به اعداد تشکیل دهنده آن نیز ژن^{۵۶} گویند. الگوریتم ژنتیک با تعدادی از جواب‌های گوناگون آغاز می‌شود که جمعیت نسل اولیه را تشکیل می‌دهند. تعداد این جوابها در تکرارها و

تک مرحله‌ای ایجاد کند.

تفکیک می‌کند. این انعطاف‌پذیری ممکن است بهبودی در کارایی را نسبت به دسته‌بندی کننده‌های

درخت وجود دارد، به شکلی که زیر مجموعه انتخاب شده به شکل بهینه بین دسته‌های این گروه را

انعطاف‌پذیری انتخاب زیر مجموعه‌های مختلفی از صفات در گروههای داخلی مختلف

کننده درخت، فقط از زیر مجموعه‌ای از صفات، برای روش بین

دسته‌ها استفاده می‌شود که معمولاً با یک معیار بهینه سراسری انتخاب می‌شود. در دسته‌بندی

تک مرحله‌ای رایج که هر نمونه، روی تمام دسته‌ها امتحان می‌شود، در یک دسته‌بندی کننده

تصمیم محلی ساده‌تر در سطوح مختلف درخت تقریب زده شوند. برخلاف دسته‌بندی کننده‌های

نواحی تصمیم پیچیده سراسری (خصوصاً در فضاهایی با ابعاد زیاد) می‌توانند با اجتماع نواحی

تصمیم محلی را در سطوح مختلف درخت تقریب زده شوند. برخلاف دسته‌بندی کننده‌های

درخت، یک نمونه روی زیر مجموعه‌های خاصی از دسته‌ها امتحان شده و محاسبات غیر لازم حذف

می‌شوند. در دسته‌بندی کننده تک مرحله‌ای، فقط از زیر مجموعه‌ای از صفات، برای روش بین

دسته‌ها استفاده می‌شود که معمولاً با یک معیار بهینه سراسری انتخاب می‌شود. در دسته‌بندی

کننده درخت، انعطاف‌پذیری انتخاب زیر مجموعه‌های مختلفی از صفات در گروههای داخلی مختلف

درخت وجود دارد، به شکلی که زیر مجموعه انتخاب شده به شکل بهینه بین دسته‌های این گروه را

تفکیک می‌کند. این انعطاف‌پذیری ممکن است بهبودی در کارایی را نسبت به دسته‌بندی کننده‌های

تک مرحله‌ای ایجاد کند.

نسل‌های بعدی ثابت هستند، ولی کیفیت آن‌ها به تدریج بهبود می‌یابد. در هر نسل جواب‌ها به کمک تابع هدف^{۵۶} سنجیده شده و بر آن اساس، نسل بعد شکل می‌گیرد. بدین ترتیب برای تولید نسل بعد، تعدادی از جواب‌ها مستقیماً به نسل بعدی منتقل می‌شود و برخی دیگر از جواب‌ها با ایجاد جهش در کروموزوم آن‌ها، به جواب‌های دیگری تبدیل شده و به نسل بعد خود، منتقل می‌گردند. در واقع الگوریتم ژنتیک، جستجو در علم رایانه برای یافتن راه حل بهینه و مسائل مرتبط با جستجو است.

درخت تصمیم‌گیری

درخت تصمیم‌گیری، ساختاری باز رخداد نگر (بازگشتی) برای بیان یک فرآیند طبقه‌بندی متناوب می‌باشد که به وسیله مجموعه‌ای از صفات^{۵۷} تشریح گردیده و یک وضعیت را به مجموعه‌ای گسسته از طبقات تخصیص می‌دهد (QUINLAN، ۱۹۸۷).

هر برگ درخت تصمیم‌گیری، نماینده یک طبقه می‌باشد. درخت تصمیم روش کارآمد ویژه‌ای برای ایجاد دسته‌بندی کننده^{۵۸}‌ها از داده‌ها است. مهم‌ترین خصوصیت درخت‌های تصمیم، قابلیت آنها در شکستن فرآیند پیچیده تصمیم‌گیری به مجموعه‌ای از تصمیمات ساده‌تر است که براحتی قابل تفسیر هستند (یزدانی و کنگاوری، ۱۳۸۳).

نواحی تصمیم پیچیده سراسری (خصوصاً در فضاهایی با ابعاد زیاد) می‌توانند با اجتماع نواحی تصمیم محلی ساده‌تر در سطوح مختلف درخت تقریب زده شوند. برخلاف دسته‌بندی کننده‌های تک مرحله‌ای رایج که هر نمونه، روی تمام دسته‌ها امتحان می‌شود، در یک دسته‌بندی کننده تک مرحله‌ای استفاده می‌شود که معمولاً با یک معیار بهینه سراسری انتخاب می‌شود. در دسته‌بندی کننده درخت، انعطاف‌پذیری انتخاب زیر مجموعه‌های مختلفی از صفات در گروههای داخلی مختلف درخت وجود دارد، به شکلی که زیر مجموعه انتخاب شده به شکل بهینه بین دسته‌های این گروه را تفکیک می‌کند. این انعطاف‌پذیری ممکن است بهبودی در کارایی را نسبت به دسته‌بندی کننده‌های تک مرحله‌ای ایجاد کند.

نظریه مجموعه اولیه

تئوری مجموعه اولیه در سال ۱۹۹۱ بوسیله پلاک^{۱۵۹} ارایه گردید. نظریه مجموعه اولیه تعمیمی از نظریه دوتایی^{۱۶۰} می‌باشد که از نظریه‌های معروف در ریاضیات بازه هست (عربانی، کلانتری، ۱۳۸۱). اگرچه مسئله کمیت اندازه‌گیری، کیفیت یادگیری و کمال دانش ناظران در مواجهه با اطلاعات ناقص هنوز به طور کامل حل نشده است، اما در هر حال نظریه مجموعه اولیه، ابزار قدرتمندی برای استخراج دانش از اطلاعات پر ابهام و ناقص می‌باشد. این نظریه می‌تواند برای شرح وابستگی بین صفات و به منظور سنجش میزان اهمیت صفات و ارتباط با داده‌های متناقض نیز به کار برده شود (۲۰۰۴، kirkdos and et.al).

استدلال مبتنی بر مورد (استدلال مورد گرا)

استدلال مبتنی بر مورد، حل مسئله به روش استنتاجی^{۱۶۱} است. این روش براساس استفاده از پاسخ مسائل قبلی، برای حل مسائل مشابه جدید، شکل گرفته است. استدلال مبتنی بر مورد، به عنوان روشی شناخته می‌شود که از نحوه رفتار انسان‌ها در برخورد با مسائل جدید الگوبرداری کرده است (فائز و قدسی‌پور و غضنفری، ۱۳۸۵) به این ترتیب که از تجربیات کسب شده در حل مسائل گذشته به عنوان راهنمایی برای حل مسئله جدید بهره می‌برد.

حل مسئله به روش CBR در برگیرنده ۴ عمل عمدی است:

- ۱- بازیابی مورد^{۱۶۲} مشابه بازیابی شده با مسئله جدید.
- ۲- استفاده از پاسخ مسئله مشابه بازیابی برای تهیه پاسخ پیشنهادی برای مسئله جدید.
- ۳- بازبینی در پاسخ پیشنهادی در صورت وجود مغایرت در شرایط مسئله جدید و مسائل بازیابی شده.
- ۴- نگهداری مورد جدید (مسئله جدید و پاسخ آن) برای استفاده در آینده. (همان، پیشین: ۵۷۱)

حوزه‌های کاربردی داده‌کاوی

به واسطه قابلیت‌های طبقه‌بندی و پیش‌بینی روش‌های داده‌کاوی، این فنون به منظور تسهیل فرآیند حسابرسی، پیش‌بینی عملکرد شرکت و کمک به ارزیابی ریسک اعتباری مورد استفاده واقع می‌شوند.

از حسابرسی به عنوان رشته‌ای قدیمی یاد می‌شود و متناوبًاً حسابرسان را به عنوان افرادی خشک

و رسمی و ایرادگیر مورد قضاوت قرار می‌دهند. در واقع این موضوع دیگر حقیقت ندارد. در سال‌های گذشته حسابرسان پی به افزایش شدید حجم معاملات و پیچیدگی‌های داده‌های مالی و غیر مالی صاحبکارشان برده‌اند (sirikulvaadhana, ۲۰۰۲).

به منظور حسابرسی این شرکت‌ها، حسابرسان پی در پی باید با داده‌های حجمی و ساختار پیچیده این داده‌ها سر و کار داشته باشند. در نتیجه، حسابرسان بیش از این نمی‌توانند تنها به ابزارهای گزارشگری یا خلاصه کردن در فرآیند حسابرسی متکی باشند. ابزارهای دیگری نظیر داده‌کاوی که به طور خودکار اطلاعات را از میان حجم داده‌های زیاد، استخراج می‌کنند، می‌توانند مفید باشند، اگر چه تطبیق فنون داده‌کاوی در فرآیند حسابرسی در واقع زمینه‌ای نسبتاً جدید می‌باشد، با وجود این، داده‌کاوی نشان داده است که در بسیاری از کاربردهای تجاری مرتبط با حسابرسی نظریه کشف تقلب، حسابداری دادگاهی^{۶۳} دارای مزیت‌های فراوانی خواهد بود. از طرفی دیگر رویدادهای اخیر نیز نشان دهنده مشکلات قابل توجهی در فرآیند حسابرسی است. سقوط انرون^{۶۴} و آرتور اندرسون^{۶۵} و سایر موارد مشابه، همچنین کاربرد رو به گسترش پختن دفاتر^{۶۶} در حرفه حسابداری شواهد دیگری برای تغییر در فرآیند حسابرسی ایجاد نموده است.

بر اساس استاندارد حسابرسی شماره ۵۶ منتشر شده توسط انجمن حسابداران رسمی امریکا، حسابرس باید انتظارات خود را گسترش داده و آنها را با مبالغ ثبت شده یا نسبت‌ها مقایسه کند. رویه‌های تحلیلی اجازه بررسی صحت مانده یک حساب را بدون آزمایش و بررسی معاملات انفرادی مرتبط می‌دهد (Kirkdos and et.al., ۲۰۰۴).

گرایش جدید حسابرسی شامل مفاهیم ریسک اقتصادی^{۶۷} بوده که بر اهداف راهبردی واحد تجاری تاکید دارد. در این رویکرد از بالا به پایین^{۶۸} حسابرس این اهداف راهبردی را شناخته و در راستای فرآیند تجاری انجام وظیفه می‌کند. روش‌های داده‌کاوی مانند NN، GA، CBR و منطق فازی^{۶۹} می‌تواند این رویکرد جدید حسابرسی مبتنی بر ریسک را تسهیل سازند.

در هر حال مقالات مرتبط با حوزه‌های کاربردی در حسابرسی غالباً اشاره بر پیش‌بینی ورشکستگی، پیش‌بینی تداوم فعالیت، مضیقه مالی و تقلب مدیریت دارند. از سویی دیگر، هسته مرکزی سایر مقالات شامل پیش‌بینی عملکرد شرکت و تخمين و ارزیابی ریسک اعتباری واحدهای تجاری است.

پیش‌بینی و رشکستگی

به نظر می‌رسد پیش‌بینی و رشکستگی مشهورترین عنوان برای کاربرد روش‌های داده‌کاوی در داده‌های مالی است. رشکستگی شرکت‌ها باعث خسارات اقتصادی فراوانی برای مدیریت، سرمایه‌گذاران، اعتباردهندگان و کارمندان شده و هزینه‌های اجتماعی^{۷۰} بسیاری نیز در برخواهد داشت. رویکردهای آماری مانند تحلیل تشخیصی و تحلیل لاجیت^{۷۱} و پروبیت^{۷۲} و همچنین رویکرد فنون هوشمند محاسباتی شامل سیستم‌های خبره^{۷۳}، شبکه‌های عصبی مصنوعی و... در فرآیند پیش‌بینی و رشکستگی شرکت‌ها قابل استفاده می‌باشد.

اغلب تحقیقات پیش‌بینی و رشکستگی، عملکرد هر یک روش‌های فوق را مقایسه می‌کنند، اما به طور قطعی نمی‌توان اظهار داشت که کدام یک از این روش‌ها، از دقت کلی بیشتری در مقایسه با سایر روش‌ها برخوردار است (Wermter and et.al ۲۰۰۷). پیش‌بینی و رشکستگی به وسیله داده‌های صورت‌های مالی نخستین بار از آثار آلتمن (۱۹۶۸) سرچشمه گرفته است. آلتمن استدلال می‌کند و رشکستگی فرآیندی بلندمدت است، از این‌رو صورت‌های مالی، می‌تواند حاوی هشدارهایی از رخداد یک و رشکستگی قریب‌الواقع باشد. وی با به کارگیری تحلیل تشخیصی چندگانه^{۷۴} مدلی را برای پیش‌بینی و رشکستگی شرکت‌ها ایجاد کرد. در همین حال بسیاری از محققان نیز مدل‌های جایگزین را با استفاده از فنون آماری گسترش می‌دادند.

لین^{۷۵} و مک لین^{۷۶} (۲۰۰۱) تلاش کردند تا با استفاده از ۴ روش متفاوت، رشکستگی شرکتها را پیش‌بینی نمایند. ۲ مورد از این روش‌ها، روش‌های آماری (تحلیل تشخیصی و رگرسیون لجستیک^{۷۷}) و دو مورد دیگر نیز از فنون یادگیری ماشین (درختان تصمیم‌گیری و شبکه‌های عصبی) بوده‌اند. همچنین آنها براساس روش‌های پیش‌گفته، الگوریتمی ترکیبی^{۷۸} را نیز پیشنهاد کردند. نمونه مورد استفاده شامل ۱۱۳۳ شرکت تجاری انگلیسی بوده است. ۶۹۰ شرکت موفق و ۱۰۶ شرکت ناموفق به عنوان مجموعه آموزشی^{۷۹} و ۲۸۹ شرکت ناموفق و ۴۸ شرکت موفق به عنوان مجموعه آزمایشی^{۸۰} استفاده شدند. تلاشی برای تطبیق کمپانی‌های موفق و ناموفق صورت نگرفت. ۳۷ مورد از نسبت‌های مالی، برگرفته از ترازنامه و صورت سود و زیان به عنوان متغیرهای ورودی^{۸۱} انتخاب شد. دو روش انتخاب ویژگی^{۸۲} بکار رفته، متغیرهای ورودی را با اعمال قضاوت انسان^{۸۳} به ۴ و با اعمال تحلیل واریانس^{۸۴} به ۱۵ مورد کاهش داده است. نتایج حاصل از روش‌های شبکه‌های عصبی مصنوعی و درختان تصمیم برای هر دو انتخاب ویژگی مبتنی بر قضاوت انسان و تحلیل واریانس، قابل قبول تر

از دیگر روش‌ها بوده است. سرانجام نویسنده‌گان مقاله، الگوریتمی ترکیبی را پیشنهاد می‌کنند که رای‌گیری وزن دار^{۸۵} طبقه‌بندی کننده‌های مختلف را به کار می‌برد. تانگ^{۸۶} و همکاران (۲۰۰۴) یک روش ترکیبی را با یکپارچه‌سازی شبکه‌های عصبی مصنوعی و سیستم‌های فازی به کار برده‌اند. این مدل که شبکه عصبی فازی خود سازمانده عمومی نامیده می‌شود، یک سیستم مبتنی بر قاعده است که مبتنی بر قانون فازی «اگر-آنگاه»^{۸۷} بوده که نهایتاً منجر به خود تعديلی (اصلاح) پارامترهای قوانین فازی با استفاده از الگوریتم‌های یادگیری می‌گردد. این الگوریتم‌ها از الگوی شبکه عصبی به دست می‌آیند. مزیت اصلی شبکه‌های عصبی فازی، توانایی در مدل‌سازی مسئله‌ای خاص با استفاده از یک مدل زبانی سطح بالای قابل فهم، بر خلاف عبارات ریاضی پیچیده می‌باشد. این مدل به منظور پیش‌بینی ورشکستگی بانک به کار رفته است. متغیرهای ورودی شامل ۹ متغیر مالی بوده‌اند که در بررسی‌های قبلی اهمیت آن‌ها محرز شده است. نمونه مورد استفاده شامل داده‌هایی در خصوص ۲۵۵۵۵ بانک موفق و ۵۴۸ بانک ناموفق بوده‌اند. ۲۰ درصد از این داده‌ها به عنوان مجموعه آموزشی و ۸۰ درصد مابقی به عنوان مجموعه آزمایش انتخاب شدند. همچنین برای کاهش خطای نوع اول^{۸۸}، نمونه نسبت به تعداد بانک‌های موفق و ناموفق، به میزان مساوی تعديل گردید. نویسنده‌گان هنگام استخراج داده‌ها از آخرین صورت‌های مالی، عملکرد مدل را ۹۳ درصد، هنگامی که این داده‌ها از صورت‌های مالی سال قبل بدست آمد میزان عملکرد را ۸۵ درصد و با دریافت داده‌ها از صورت‌های مالی دو سال قبل عملکرد را ۷۵ درصد، گزارش نمودند. مدل مورد نظر، در حدود ۵۰ قاعده فازی اگر-آنگاه، تولید کرد که تعاملات بین ۹ متغیر ورودی منتخب و اثرشان بر سلامت مالی بانک‌های مورد بررسی را، تشریح نموده است.

شین^{۸۹} و لی^{۹۰} (۲۰۰۲) نیز مدلی مبتنی بر الگوریتم ژنتیک پیشنهاد دادند. آنها بر این واقعیت تاکید کردند که برخلاف شبکه‌های عصبی مصنوعی، الگوریتم‌های ژنتیک می‌توانند قوانینی قابل فهم ایجاد کنند. در این مورد، الگوریتم‌های ژنتیک به منظور یافتن آستانه‌های یک یا چند متغیر در بالا یا پایین سطحی به کار رفته‌ند که موقعیت یک واحد تجاری را بحرانی می‌سازد. این مدل از ساختاری قاعده‌مند بهره می‌برد که دارای ۵ شرط است که هر کدام به یک متغیر خارج از ۹ نسبت مالی، اشاره دارد. این شرایط با عملکرد عطف منطقی^{۹۱}، ترکیب خواهند شد. مجموعه داده‌ها مرکب از ۲۶۴ شرکت موفق و ۲۶۴ شرکت ناموفق بوده و ۹ نسبت مالی به عنوان متغیرهای ورودی انتخاب گردیدند. ۹۰ درصد نمونه به عنوان مجموعه آموزشی و ۱۰ درصد نیز به منظور اعتبار سنجی به کار

رفته‌اند. دقت کلی روش مزبور در حدود ۸۰ درصد گزارش شد.^{۹۲} دیمیتراس^{۹۳} و همکاران (۱۹۹۸) به بررسی پیش‌بینی ورشکستگی با استفاده از نظریه مجموعه اولیه پرداختند. مجموعه آموزشی مورد استفاده متشکل از ۴۰ شرکت ناموفق و ۴۰ شرکت موفق از میان واحدهای تجاری یونانی در یک دوره ۵ ساله و مجموعه آزمایش نیز از ۱۹ شرکت ناموفق و ۱۹ شرکت موفق تشکیل شده بود. ۱۲ نسبت مالی بر اساس نظر مدیر اعتباری بانکی یونانی به منظور ورود به جدول اطلاعات انتخاب گردید. سرانجام قواعد تصمیم^{۹۴}، استنتاج شد و نتایج این روش با نتایج تحلیل تشخیصی و تحلیل لاجیت مقایسه و برتری این روش نسبت به سایر روش‌ها مشخص گردید. پارک^{۹۵} و هان^{۹۶} نیز در سال ۲۰۰۲ در یک تحقیق براساس روش استدلال مبتنی بر مورد، مدلی برای پیش‌بینی ورشکستگی بانک‌ها طراحی کردند.

تمامیت و مضیقه مالی

رسایی‌هایی حسابرسی در سال‌های اخیر، اهمیت تجزیه و تحلیل‌های حسابرسی‌های کامل را در چندان ساخته است. براساس استاندارد حسابرسی شماره ۵۹^{۹۷} حسابرس باید توانایی صاحبکار خود را در مورد امکان ادامه فعالیت حداقل یک سال بعد از تاریخ ترازنامه ارزیابی کند. چنان‌چه نشانه‌هایی از وجود مشکلات مالی دیده شود که احتمالاً منجر به ورشکستگی خواهد شد، حسابرس می‌باید در خصوص تمامیت صاحبکار اظهارنظر کند. ارزیابی وضعیت تمامیت شرکت‌ها کارگندان ساده‌ای نیست. مطالعات نشان می‌دهند بخشی نسبتاً کوچک از واحدهای تجاری ورشکسته، وارد شرایط مبنای تمامیت شرکت هستند (KIRKDOS and et.al. ۲۰۰۴). به هر حال جهت تسهیل وظیفه حسابرسان و وظایف آنها در مورد وضعیت تمامیت شرکت‌ها، روش‌های آماری و یادگیری ماشین پیشنهاد شده‌اند.

کاه^{۹۸} (۲۰۰۴) در بررسی موردی^{۹۹}، پیش‌بینی تمامیت روش‌هایی نظری شبکه‌های عصبی مصنوعی، درختان تصمیم و رگرسیون لجستیک را با یکدیگر مقایسه کرده است. نمونه مورد استفاده شامل ۱۶۵ شرکت دارای فعالیت و ۱۶۵ شرکت فاقد فعالیت بوده است. ۶ نسبت مالی برگزیده به عنوان متغیرهای ورودی انتخاب شدند. نهایتاً، نتایج نشان می‌دهد روش درختان تصمیم نسبت به دو روش دیگر از عملکرد بهتری برخوردار بوده است.

тан^{۱۰۰} و دیهار دجو^{۱۰۱} (۲۰۰۱) براساس تحقیقی از تان، سعی در پیش‌بینی مضیقه مالی

تقلب مدیریت

تحقیقات نشان داده‌اند که پیش‌بینی تقلب، تقریباً غیر ممکن است. در واقع بسیاری از خصوصیات و ویژگی‌های مرتبط با جرم کارمندی^{۱۰۹}، دقیقاً همان‌هایی هستند که سازمان‌ها به هنگام استخدام کارمندان خود در جستجوی آنها هستند (ردگر و همکاران، ۲۰۰۳). تقلب مدیریت، تقلبی فraigir است که مدیریت بواسطه تحریف صورت‌های مالی مرتکب می‌شود. تقلب مدیریت می‌تواند منافع مسئولین مالیاتی، سهامداران و اعتباردهندگان را با مخاطره مواجه کند.

اسپاتیس^{۱۱۰} (۲۰۰۲) از دو مدل برای شناسایی تحریف صورت‌های مالی از طریق داده‌های قابل دسترس عمومی استفاده کرد. متغیرهای ورودی برای نخستین مدل شامل ۹ نسبت مالی بوده و برای دومین مدل امتیاز Z^{۱۱۱} به عنوان متغیرهای ورودی به منظور تطبیق رابطه میان درماندگی مالی و دستکاری ارقام مالی، اضافه گردید. این روش از رگرسیون لجستیکی بهره می‌برد. داده‌های نمونه

اتحادیه‌های اعتباری استرالیا^{۱۱۲} براساس روش شبکه‌های عصبی مصنوعی داشته‌اند. در تحقیق قبلی از تان، داده‌های مالی سه ماهه استفاده شد و تلاش شد که درماندگی مالی براساس رویکردی سه ماهه پیش‌بینی گردد. آنها مدل مزبور را با وارد کردن مفهومی به نام پیش‌یاب^{۱۱۳} ارتقاء دادند. بدین ترتیب وقتی مدل، پیش‌بینی می‌کند یک اتحادیه اعتباری در سه ماهه‌ای خاص دچار درماندگی مالی خواهد شد و واقعاً آن اتفاق در سه ماهه بعدی می‌افتد، به آن سه ماهه پیش‌یاب گفته می‌شود. این روش نسبت به روش قبلی به لحاظ میزان خطای نوع دوم^{۱۱۴} از عملکرد بهتری برخوردار بوده است. در این مطالعه^{۱۱۵} نسبت مالی به عنوان متغیر ورودی، مورد استفاده واقع شده و مشاهدات نمونه شامل ۲۱۴۴ مورد بوده است. نتایج بدست آمده با نتایج مدل پروبیت مقایسه و به طور نامحسوسی قابلیت آن برای میزان خطای نوع اول بیشتر تخمین زده شد. کنو^{۱۱۶} و کوباشی^{۱۱۷} (۲۰۰۰) روشی را برای رتبه‌بندی واحدهای تجاری براساس فنون برنامه‌نویسی ریاضی ارایه کردند. این مدل هیچ فرضی در خصوص توزیع تحمیل نمی‌کند.

سه روش جایگزین شامل تشخیص مبتنی بر تقابل مربعی^{۱۱۸} (رویه درجه دوم)، تشخیص مبتنی بر ابر صفحه^{۱۱۹} و تشخیص مبتنی بر رویه بیضوی^{۱۱۱} بوده‌اند. ۶ نسبت مالی بدست آمده از صورت‌های مالی به عنوان متغیرهای ورودی و داده‌های نمونه شامل ۴۵۵ واحد تجاری بوده‌اند. در نهایت مدل برای هر واحد تجاری امتیازی خاص محاسبه کرده است.

مرکب از ۳۸ مورد تحریف صورت‌های مالی و ۳۸ مورد عدم تحریف بوده است. نتایج هر دو مدل نشان داد سه متغیر، با ضریب معنا داری^{۱۱۲} وارد مدل شده‌اند.

پیش‌بینی عملکرد واحد تجاری

با وجود چالش‌های ناشی از تفکیک مالکیت از مدیریت و تضاد منافع بین گروههای ذی نفع، شاهد هستیم از معیارهای حسابداری نظری سود هر سهم، نسبت بازدهی حقوق صاحبان سهام و یا بازدهی سرمایه‌گذاری به منظور ارزیابی عملکرد واحدهای تجاری استفاده می‌شود (ایزدی نیا، ۱۳۸۴).^{۱۱۳} لم (۲۰۰۳) مدلی را به منظور پیش‌بینی نرخ بازده حقوق سهامداران عادی بکار گرفته است. وی از شبکه‌های عصبی پس انتشار خطاط در این خصوص استفاده کرد. بردار ورودی شامل ۱۵ نسبت مالی و یک متغیر از نوع تحلیل تکنیکال^{۱۱۴} بود.

بک و همکاران^{۱۱۵} (۲۰۰۰) اقدام به ارایه دو مدل به منظور گروه بندی شرکت‌ها براساس عملکرد نمودند. در هر دو مدل از نقشه‌های خود سازمان یافته استفاده شده بود. نخستین مدل با داده‌های مالی ۱۶۰ شرکت سر و کار داشته و بر اساس به کارگیری روش متن کاوی^{۱۱۶} انجام یافت. دومین مدل نیز به تحلیل و بررسی گزارش‌های سالانه ارایه شده به مدیران ارشد اجرایی می‌پردازد. نویسنده‌گان تفاوت‌هایی میان نتایج گروه بندی در دو روش مشاهده می‌کنند.

ارزیابی ریسک اعتباری

تجزیه و تحلیل ریسک اعتباری به واسطه افزایش تعداد شرکت‌های ورشکسته و پیشنهادهای رقابتی اعتبار دهنده‌گان، تبدیل به یکی از محبوب‌ترین حوزه‌های مالی شده است. روش‌های داده کاوی به منظور تسهیل پیش‌بینی ریسک اعتباری می‌تواند به کار برد شود. هانگ و همکاران^{۱۱۷} (۲۰۰۳) تجزیه و تحلیل رتبه‌بندی ریسک اعتباری^{۱۱۸} را با استفاده از ماشین بردار پشتیبان^{۱۱۹} که یک روش یادگیری ماشین است، اجرا کردند. دو مجموعه داده مورد استفاده در این تحقیق، شامل ۷۴ شرکت کره‌ای و دیگری شامل ۲۶۵ شرکت آمریکایی بوده است. برای هر مجموعه داده، ۵ نوع رتبه‌بندی تعريف شد. همچنین، دو مدل برای داده‌های شرکت‌های کره‌ای و دو مدل نیز برای شرکت‌های آمریکایی هر کدام با یک بردار ورودی به کار رفت. ماشین بردار پشتیبان و شبکه‌های عصبی پس از انتشار خطاط برای اهداف پیش‌بینی رتبه‌بندی اعتباری استفاده شدند. مدل گارسون^{۱۲۰} نیز برای

نتیجه گیری

فاوری داده کاوی از قابلیت های پیش بینی و طبقه بندی فراوانی برخوردار است. بنابراین می تواند فرآیند تصمیم گیری در مسایل مالی را تسهیل کند. کاربرد روش های داده کاوی با توجه به مطالعات مرتبط و ماهیت آنها، می تواند طیف گسترده ای شامل پیش بینی و روشکستگی، تخمین ریسک اعتباری، وضعیت تداوم فعالیت، مضیقه مالی، پیش بینی عملکرد واحد تجاری و تقلب مدیریت را در بر گیرد. در این بین، پیش بینی و روشکستگی مشهور ترین این عناوین است. همچنین، روش های داده کاوی به کار رفته در تحقیقات مختلف نیز شامل شبکه های عصبی مصنوعی، الگوریتم ژنتیک، درخت تصمیم گیری، نظریه مجموعه اولیه، استدلال مبتنی بر مورد و برنامه نویسی ریاضی هستند. بر اساس بررسی های صورت گرفته، روش شبکه های عصبی مصنوعی تاکنون یکی از محبوب ترین روش های داده کاوی بوده است. اگرچه تعداد زیادی از تحقیقات در ارتباط با کاربرد روش های داده کاوی در امور مالی بوده اند، لیکن حوزه های مستعد دیگری به منظور گسترش دامنه تحقیقات وجود دارند. معرفی مدل های ترکیبی، بهبود مدل های کنونی، استخراج قواعدی جامع از شبکه های عصبی مصنوعی، بهسازی و یکپارچه سازی سیستم های ERP با ابزارهای داده کاوی می توانند گستره ای نوین و متفاوت از تحقیقات داده کاوی باشند. ارزیابی هزینه های مناسب به خطاهای نوع اول و دوم، توسعه بردار ورودی بوسیله اطلاعات کیفی و همچنین بهره برداری از روش های معتبر انتخاب ویژگی، می توانند از جمله سایر روش های تحقیق در حوزه داده کاوی در آینده باشند. بدین ترتیب با ارتقای مدل ها و

روش‌های مورد استفاده، می‌توان تبدیل شدن داده‌کاوی به ابزاری مفیدتر در فرآیند تصمیم‌گیری به خصوص در حوزه‌های مالی را، به نظراره نشست.

پی‌نوشت‌ها:

- 1- Artificial Intelligence
- 2- Pattern Recognition
- 3- Database System
- 4- Statistics
- 5- Data Mining(DM)
- 6- Knowledge Discovery In Database (KDD)
- 7- Data Cleaning
- 8- Data Integration
- 9- Data Selection
- 10- Data Transformation
- 11- Pattern Evaluation
- 12- Machine Learning
- 13- Human-Computer Interface
- 14- Financial Classification Problem
- 15- Corporate Bankruptcy
- 16- Credit Risk Estimation
- 17- Going Concern Reporting
- 18- Financial Distress
- 19- Corporate Performance Prediction

۲۰- در این مقاله سعی شده به بیان کاربردهایی از داده‌کاوی پردازیم که غالبا هسته مرکزی بیشتر پژوهش‌ها بوده‌اند. در هر حال سایر کاربردهای داده‌کاوی مطابق مقالات موجود می‌تواند شامل موارد زیر باشد:
- کاربرد داده‌کاوی در تحلیل‌های مرتبط با پولشویی
- کاربرد داده‌کاوی در حسابداری دادگاهی
- استفاده از تکنیک‌های داده‌کاوی در مدیریت ارتباط با مشتری
- استفاده از داده‌کاوی در بخش عمومی

- 21- Supply Chain
- 22- Data Ming Tool
- 23- Hidden Pattern
- 24- Www.Ncsa.Uiuc.Edu/Divisions/Dmv/Alg/D2k
- 25- Www.Sas.Com/
- 26- Www.Spss.Com/Spssbi/Clementine/
- 27- Www.3.Ibm.Com/Software/Data/Miner/
- 28- Www.Insightful.Com/Products/Product.Asp?Pid=26
- 29- Www.Acl.Com/Product/
- 30- Data Ming Technique
- 31- Prediction
- 32- Description

۳۳- مجموعه‌ای از قواعد و مراحل منظم که در حل مسائل ریاضی به صورت ثابت و متحدد الشکل بکار می‌رond
34- Visualization
35- Neural Network (NN)

- 36- Genetic Algorithms(GA)
- 37- Decision Tree
- 38- Rough Set Theory(RST)
- 39- Case Base Reasoning(CBR)
- 40- Mathematical Programming
- 41- Node
- 42- Neuron
- 43- Weight
- 44- Signal
- 45- Threshold
- 46- Input Layer
- 47- Output Layer
- 48- Hidden Layer
- 49- Self-Organizing Map(SOM)
- 50- Back Propagation

۵۱- در شبکه خود سازمان یافته، از روش یادگیری رقابتی برای آموزش استفاده می‌شود و مبتنی بر مشخصه‌های خاصی از مغز انسان توسعه یافته است. سلول‌ها در مغز انسان در نواحی مختلف طوری سازمان دهی شده‌اند که در نواحی حسی مختلف، با نقشه‌های محاسباتی مرتب و معنی‌دار ارایه می‌شوند. برای نمونه ورودی‌های حسی لامسه شنوایی، و... با یک ترتیب هندسی معنی دار به نواحی مختلف مرتبط هستند. در یک شبکه خود سازمان یافته، واحدهای پردازشگر در گره‌های یک شبکه یک بعدی، دو بعدی یا بیشتر قرارداده می‌شوند. واحدهای در یک فرآیند یادگیری رقابتی نسبت به الگوهای ورودی، منظم می‌شوند. محل واحدهای تنظیم شده در شبکه به گونه‌ای نظم می‌یابد که برای ویژگی‌های ورودی، یک‌ستگاه مختصات معنی دار روی شبکه ایجاد شود. لذا نقشه خود سازمان یافته، یک نقشه توبوگرافیک از الگوهای ورودی را تشکیل می‌دهد که در آن محل قرار گرفتن واحدهای، متناظر ویژگی‌های ذاتی الگوهای ورودی است. یادگیری رقابتی که در این قبیل شبکه‌های بکار گرفته می‌شود بدین صورت است که در هر قدم یادگیری، واحدهای برای فعل شدن با یکدیگر به رقابت می‌پردازند. در پایان یک مرحله رقابت تنها یک واحد برنده می‌شود، که وزن‌های آن نسبت به وزن‌های سایر واحدها به شکل متفاوتی تغییر داده می‌شود. این نوع از یادگیری را یادگیری بدون مربی یا بی‌نظرارت گویند.

- 52- Rectangular Lattice
- 53- Hexagonal Lattice
- 54- Holland
- 55- Gene

۵۶- تابع هدف جهت تعیین اینکه افراد چگونه در محدوده یک مساله ایفای نقش می‌کنند، استفاده می‌شود.

۵۷- Attribute

58- Classifier

59- Pawlak

60- Twin

۶۱- به عنوان زیر مجموعه‌ای گسترده از هوش مصنوعی، یادگیری ماشین وظیفه طراحی و توسعه الگوریتم‌هایی را به عهده دارد که به کامپیوتر اجازه یادگیری می‌دهند. به طور کلی همواره با دو نوع یادگیری استنتاجی و استقرایی روبرو هستیم.

- 62- case
- 63- Forensic Accounting
- 64- Enron
- 65- Arthur Andersen

۶۶- به معنی یک‌ست نشان دادن سود می‌باشد. (Book Cooking)

۶۷- احتمال خطر اینکه واحد اقتصادی از پوشش هزینه‌های عملیاتی خود عاجز باشد.

68- Top Down Approach

69- Fuzzy Logic

70- Social Cost

۷۱- مدل لاجیت با اختصاص وزن هایی به متغیرهای مستقل، رتبه هریک از شرکتهای نمونه را پیش بینی می کند. از این رتبه برای تعیین احتمال عضویت در یک گروه معین (برای مثال موفق یا ناموفق) استفاده می شود.

72- Probit

73- Expert System

74- Multiple Discriminant Analysis

75- lin

76- mcclean

۷۷- رگرسیون لوچستیک یک مدل آماری رگرسیون برای متغیرهای وابسته دودویی است. این مدل را می توان به عنوان مدل خطی تعمیم یافته ای که ازتابع لوچستیک به عنوانتابع پیوند استفاده می کند و خطایش از توزیع چند جمله ای پیروی می کند، به حساب آورد.

78- Hybrid Algorithm

79- Training Set

80- Testing Set

81- Input Variable

82- Feature Selection

83- Human Judgment

84- Analysis Of Variance(ANOVA)

۸۵- یک سیستم رای گیری است بر این پایه که تمام رای دهنده‌گان لزوماً با یک دیگر یکسان نمی باشند.

86- Tung

87- If-Then

۸۸- رد کردن فرضیه صفر که در واقع درست است.

89- shin

90- Lee

91- Logical AND Operator

92- Dimitras

93- Decision Rule

94- Park

95- Han

96- Statement Of Auditing Standards No.59: "The Auditors Consideration Of An Entity's Ability To Continue As A Going Concern"

97- KOH

98- Case study

99- Tan

100- Dihardjo

101- Australian Credit Union

102- Early Detector

۱۰۳- عدم رد فرضیه صفر که در واقع غلط است.

104- konno

105- Kobashi

106- Discrimination By Hyperplane

107- Discrimination By Quadratic Surface

108- Discrimination By Elliptic Surface

109- White –Collar Criminal

110- Spathis

منابع و مأخذ:

- ۱۱۱- Z Score
 ۱۱۲- Significant Coefficient
 ۱۱۳- Lam
 ۱۱۴- Technical Analysis Variable
 ۱۱۵- Back and et.al
 ۱۱۶- امروزه بانک‌های متنی بواسطه افزایش اطلاعات به شکل الکترونیکی مانند نشر الکترونیکی، کتابخانه‌های دیجیتال، پست الکترونیکی و وب جهانی به سرعت در حال رشد می‌باشند. متن کاوی در واقع نوعی از روش‌های داده کاوی است که به جستجو و استخراج اطلاعات جدید در مجموعه بزرگ داده‌های متنی، می‌پردازد.
- ۱۱۷- Hung and et.al
 ۱۱۸- Credit Rating Analyzing
 ۱۱۹- Support Vector Machine(SVM)
 ۱۲۰- Garson Model
 ۱۲۱- Hung and et.al
 ۱۲۲- Decision Diagram
 ۱۲۳- اتحادیه گمرکی بلژیک و هلند و لوکزامبوزگ که به تدریج معنی شخصیت‌های حقوقی این سه کشور را نیز به خود گرفته است.
 ۱۲۴- Rule Extraction

- ۱- اصغری اسکویی، محمد رضا(۱۳۸۱) کاربرد شبکه‌های عصبی در پیش‌بینی سریهای زمانی، فصلنامه پژوهش‌های اقتصادی ایران، شماره ۱۲.
- ۲- ایزدی‌نیا، ناصر(۱۳۸۴) نقدی بر معیارهای حسابداری ارزیابی عملکرد و پیشنهاد معیارهای ارزش افزوده اقتصادی و جریان‌های آزاد نقدی برای گزارشگری ارزش‌های واحد تجاری، مجله دانشکده علوم اداری و اقتصاد دانشگاه اصفهان.
- ۳- بهرام‌فر، نقی و جواد ساعی(۱۳۸۵) ارایه مدل برای پیش‌بینی عملکرد شرکت‌های پذیرفته شده در بورس اوراق بهادار تهران با استفاده از اطلاعات مالی منتشره، بررسی‌های حسابداری و حسابرسی، شماره ۴۵.
- ۴- صادقی مقدم، محمدرضا، امیر افسر و بابک سهرابی(۱۳۸۴) مدل‌سازی جریان مواد زنجیره تأمین با رویکرد الگوریتم ژنتیک، مجله علمی پژوهشی مدرس (علوم انسانی).
- ۵- صفائی، امین(۱۳۸۵) چگونه یک شبکه عصبی هوشمند بسازیم؟ - مثالی از برنامه‌نویسی شیء‌گرا در شبکه‌های عصبی و هوش مصنوعی، ماهنامه شبکه، شماره ۷۱.
- ۶- صورتگر، ابوالفضل و علیرضا داوودآبادی(۱۳۸۵) تشخیص اشغال کاذب در تجهیزات مخابراتی توسط شبکه‌های عصبی، مجله الکترونیکی پژوهشگاه اطلاعات و مدارک علمی ایران، شماره اول دوره ششم، قابل دسترس به صورت آنلاین در: davoodabadi_abs.htm/http://www.irandoc.ac.ir/data/E_J/vol
- ۷- عربانی، مهیار و فرزین کلانتری(۱۳۸۱) ارزیابی استفاده از تئوری Rough Set در سنجش عوامل ناپایداری شبیه، سومین همایش بین‌المللی مهندسی ژئوتکنیک و مهندسی خاک ایران.
- ۸- فائز، فرهاد، حسن قدسی‌پور و مهدی غضنفری(۱۳۸۵) ارایه یک مدل تصمیم یار برای انتخاب فروشنده با استفاده از روش استدلال مبتنی بر مورد در محیط فازی، نشریه دانشکده فنی، شماره ۴.
- ۹- محمد کانتار‌دزیک(۲۰۰۲) داده کاوی، ترجمه امیر علیخانزاده. چاپ اول، بابل: انتشارات علوم رایانه.
- ۱۰- موسی زادگان، حسنعلی و حسام الدین ذگردی(۱۳۸۷) مدلی جدید برای حل مسئله موازنۀ خط مونتاژ هزینه گرا، نشریه بین‌المللی علوم مهندسی دانشگاه علم و صنعت ایران، شماره ۱۹.
- ۱۱- نجات، امیر رضا و آرش علی‌اکبری(۱۳۸۷) داده کاوی راهی به سوی ناشناخته‌ها. دو ماهنامه توسعه انسانی پلیس، شماره ۱۸.

- ۱۲- همایون بور، مهدی و داریوش حکیم‌زاده (۱۳۸۱) سنجش وضعیت در ماشین‌های الکترونیکی با استفاده از شبکه‌های عصبی، هشتمین کنفرانس انجمن کامپیوتر ایران.
- 13- Abd Rahman (2008). Utilization Of Data Mining Technology Within The Accounting Information System In The Public Sector: A Country Study – Malaysia. Degree of Doctor of Philosophy, University of Tasmania.
- 14- Dass.R(2006).Data Mining In Banking And Finance: A Note For Bankers. Available on line at: <Http://CiteSeerx.Ist.Psu.Edu/Viewdoc/Download?DOI=10.1.1.102.7502&Rep=Rep1&Type=Pdf>
- 15- Kirkos S. and Manolopoulos Y., (2004), 'Data Mining in Finance and Accounting: A Review of Current Research Trends', Proceedings of the 1st International Conference on Enterprise Systems and Accounting (ICESAcc), Thessaloniki,Greece, pp. 63-78.
- 16- Rui. Hort, Beatriz De S. L. Pires De Lima And Carlos Cristiano H. Borges (2009).Data Preprocessing Of Bankruptcy Prediction Models Using Data Mining Techniques. Available on line At: Www.Blog.Campe.Com.Br/Wp-Content/Uploads/2009/03/Witpress_Conf-2.Pdf.
- 17- James A. Rodger, Parag C. Pendharka, David J. Pape And George Letcher(2003). Utilization Of Data Mining Techniques To Detect And Predict Accounting Fraud: A Comparison Of Neural Networks And Discriminant Analysis. Available on line At: <Http://Portal.Acm.Org/Citation.Cfm?Id=778616>.
- 18- Gary M. Weiss, Maytal Saar-Tsechansky And Bianca Zadrozny(2005). Utility-Based Data Mining. Available on line At: <Www.Springer.Com/Journal/10618>.
- 19- Gersper.m(2007). Data Mining Helps. Available on line At: www.Gdmllc.Com/Pdf/Data_Mining_Help_Financial_Executives.Pdf.
- 20- Hung,C , Chen,J, Wermter, S. (2007). Hybrid Probability-Based Ensembles For Bankruptcy Prediction. International Conference on Business And Information, July 11-13, 2007, Tokyo, Japan.
- 21- Pechenizkiy, S. Puuronen And A. Tsymbal(2006). On The Use of Information Systems Research Methods in Data Mining. Information Systems Development. 487-499
- 22-Quinlan.J. Generating Production Rules from Decision Trees (1987). IJCAI'87 Proceedings of the 10th international joint conference on Artificial intelligence - Volume 1
- 23- Sirikulvadhana. S (2002). Data Mining As A Financial Auditing Tool. Available on line At: <Pafis.Shh.Fi/Graduates/Supsir01.Pdf>.
- 24- Jeffrey W. Seifert (2004). Data Mining: An Overview. Available on line At: <Www.Fas.Org/Irp/Crs/Rl31798.Pdf>.