

دسته‌بندی مشتریان بیمه با استفاده از داده‌کاوی



نویسنده: سید محمود ایزدپرست

– کارشناس ارشد مدیریت فناوری اطلاعات، دانشگاه پیام نور تهران

پیکیده

امروزه نقش مشتریان از حالت پیروی از تولیدکننده، به هدایت تولیدکنندگان مبدل گشته است، به همین دلیل دسته‌بندی مشتریان، در هدفمند ساختن سازمان‌ها در سفارشی‌سازی خدمات‌شان و نیز الویت‌بندی محصولات براساس میزان سودآوری آن محصول کمک شایانی می‌کند. روش داده‌کاوی برای دستیابی به قوانین تصمیم‌گیری و مدل پیش‌بینی رفتار مشتریان آتی در یکی از شرکت‌های بیمه استفاده شده است. در اجرای روش دسته‌بندی در داده‌کاوی، دو تکنیک درخت تصمیم و شبکه‌های عصبی در یکی از شرکت‌های بیمه فصولی به‌کاررفته است. هدف این پژوهش استفاده از روش شبکه‌های عصبی و تکنیک درخت تصمیم به منظور دسته‌بندی مشتریان بیمه و درنهایت ارزیابی نتایج به دست آمده است. به منظور ارزیابی مدل‌ها، نتایج حاصل از دو مدل را مقایسه کردیم که تطابق آنها نشان‌دهنده صحت عملکرد مدل‌هاست. البته بررسی‌ها نشان داده‌اند که روش درخت تصمیم نتایج بهتری را دربرداشته است و این بدان معنی است که روش درخت تصمیم، روش مناسب‌تری برای دسته‌بندی مشتریان بیمه ایجاد می‌کند. همچنین نتایج به دست آمده از این تحقیق، توسط فبرگان صنعت بیمه تأیید و نظرات آنها به شکل توفیقی از مشخصه‌ها یا متغیرهای ورودی تحقیق ارائه شده است.

واژگان کلیدی: داده‌کاوی، بیمه، فناوری اطلاعات، دسته‌بندی، درخت تصمیم، شبکه‌های عصبی



مقدمه

داده کاوی ارائه می‌کنیم و سپس روش‌های دسته‌بندی که به کار گرفته شده (درخت تصمیم^۲ و شبکه‌های عصبی^۳) و اینکه چگونه از این روش‌ها برای دسته‌بندی مشتریان استفاده کرده‌ایم را توضیح می‌دهیم. پس از آن به منظور ارزیابی نتایج به دست آمده، این دو روش را مقایسه می‌کنیم تا ببینیم کدام روش، نتایج دقیق‌تری حاصل می‌نماید. در نهایت پیشنهادهایی برای تحقیقات آتی در این زمینه ارائه می‌کنیم.

۱. مفهوم داده کاوی^۱

از زمانی که علم آمار به وجود آمد دانشمندان نیاز به کشف خصوصیات داده‌ها را احساس کرده بودند. با استفاده از آمار در آن زمان، خصوصیات داده‌ها از قبیل پراکندگی و تمرکز آنها بررسی می‌شد. با افزایش نیاز به استفاده از داده‌ها و درک ارزش اطلاعات، داده‌ها به سرعت در حال گردآوری و ذخیره شدن می‌باشند که این سرعت همه روزه در حال افزایش است (Chen, & Su, 2006). به موازات سرعت زیاد ذخیره شدن داده‌ها، ابزارهای محاسبه مکانیکی، الکتریکی و در نهایت کامپیوتری نیز به وجود آمده‌اند که نرخ افزایش سرعت محاسبه این ابزارها نیز نمایی است. حجم زیاد داده‌های ذخیره شده و نیز گستردگی ابعاد آن که معمولاً از منابع گوناگون تهیه شده بودند و بعضاً دارای قالب‌های متفاوتی نیز بودند،

در جهان کنونی که امکان تولید انبوه کالا و خدمات، زمینه لازم افزایش عرضه، نسبت به تقاضا را فراهم آورده است، برای تولیدکنندگان راهی جز جلب رضایت مشتری باقی نمانده و دیگر نمی‌توان حیطة بازار و عرضه را با ابزارهای محدود گذشته تعریف کرد (Berson et al, 2001). تجربه نشان داده است، سازمان‌هایی که از نظر سنتی به مفاهیم مشتری، کالا، بازار، فروش، خرید، رقابت، تبلیغات و کیفیت نگاه می‌کنند، علاوه بر عدم کسب موفقیت، سرمایه‌های خود را هم از دست می‌دهند. با ظهور اقتصاد رقابتی، مفاهیمی چون مشتری‌مداری و کسب رضایت مشتری، پایه و اساس کسب و کار تلقی شده و سازمانی که بدان بی‌توجه باشد و این مطلب را در نظر نگیرد از صحنه بازار حذف خواهد شد. در این تحقیق ابتدا داده‌های جمع‌آوری شده را توسط روش‌های پیش‌پردازش غربال می‌کنیم تا ناخالصی‌ها و داده‌های ناقص آن حذف گردد. سپس سعی می‌کنیم با استفاده از روش‌های داده کاوی و به خصوص روش‌های دسته‌بندی^۱، مشتریان بیمه را براساس ویژگی‌هایشان دسته‌بندی کنیم و بدین ترتیب مشتریان آتی بیمه را می‌توان در یکی از دسته‌ها قرار داد و با توجه به اطلاعاتی که از هر کدام از این دسته‌ها داریم می‌توان ویژگی‌های آنها را پیش‌بینی کرد. در این مقاله ابتدا مقدمه‌ای در زمینه

2. Decision Tree
3. Neural Network
4. Data Mining



دسته‌بندی، پیش‌بینی^۲، تخمین^۳ و خوشه‌بندی^۴ داده‌ها استفاده کرد. برای انجام این کارها تکنیک‌هایی توسعه یافته‌اند که با توجه به پیشرفت کامپیوترها و این علم همه روزه بر تعداد و کیفیت این تکنیک‌ها افزوده می‌شود. تعدادی از معروف‌ترین این تکنیک‌ها عبارت‌اند از: الگوریتم‌های خوشه‌بندی^۵، شبکه‌های عصبی، الگوریتم ژنتیک^۶، نزدیک‌ترین همسایگی^۷ و درخت تصمیم‌گیری (Clifton & Thuraisingham, 2001).

۲. تشریح تکنیک دسته‌بندی

دسته‌بندی از مسائل متعارفی است که به‌طور گسترده توسط متخصصان آمار و محققان فراگیری ماشینی مطالعه شده است. ارائه یک تعریف دقیق از روش دسته‌بندی، دشوار است اما مطابق با تعریف کلی، تکنیک دسته‌بندی، جداسازی یا قراردادن اجزا یا اشیا در تعدادی از کلاس‌هاست (Tan, & Yu, 2006).

اگر کلاس‌ها بدون آزمون داده‌ها (به‌طور غیر تجربی) ساخته شده باشند، دسته‌بندی استقرایی^۸ نامیده می‌شود. در مقابل، اگر کلاس‌ها به‌طور تجربی ساخته شده باشند (با آزمایش داده‌ها)، دسته‌بندی مؤخر^۹ نامیده می‌شود. در

2. Prediction
3. Estimation
4. Clustering
5. Cluster Detection Algorithm
6. Genetic Algorithm
7. Nearest Neighboring
8. Apriori Classification
9. Posteriori Classification

سبب شد که در بسیاری موارد روش‌های آماری به تنهایی قادر به کشف خصوصیات داده‌ها نباشند. زمانی که می‌خواهیم به بررسی تأثیر تعداد کمی از عوامل بر روی هدف پردازیم معمولاً روش‌های آماری مناسب است. ولی زمانی که تعداد این عوامل زیاد می‌شود دیگر این روش‌ها کارایی مناسبی ندارند. مثلاً در تحلیل داده‌های شبکه‌های زندگی افراد^۱ به دلیل اینکه این داده‌ها دارای ابعاد بسیار زیادی می‌باشند کمتر از روش‌های آماری استفاده می‌شود. دانشمندان برای رفع این مشکل تصمیم گرفتند که از سرعت بالای کامپیوترها استفاده کنند، همین امر سبب شد که روش‌های ابتکاری دیگری علاوه بر روش‌های آماری مثل شبکه‌های عصبی و الگوریتم ژنتیک برای کمک به این منظور ایجاد شود (SY, 2001). سه موضوع ذخیره‌سازی داده‌ها، افزایش سرعت کامپیوترها و پیدایش الگوریتم‌های جدید کار با داده‌ها، عامل ایجاد علمی به نام داده‌کاوی شده است. داده‌کاوی عبارت است از «استخراج اطلاعات و دانش و کشف الگوهای پنهان از یک پایگاه داده‌های بسیار بزرگ». که این الگوها و دانش‌ها معمولاً مستتر در داده می‌باشند (Chan & Lewis, 2002).

از داده‌کاوی می‌توان برای انجام کارهایی مثل

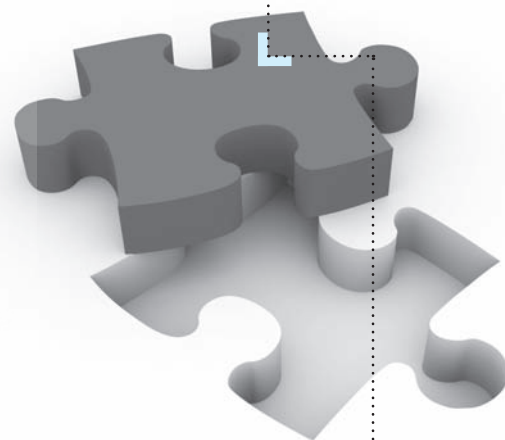
1. Demographic

بیشتر ادبیات حاکم بر دسته‌بندی فرض می‌شود که کلاس‌ها استقرایی باشند و دسته‌بندی که شامل آموزش سیستمی شده و در زمان معرفی یک شیء جدید به سیستم آزموده شده است، قابلیت اختصاص آن شیء به یکی از کلاس‌های موجود را دارد. این روش همچنین یادگیری هدایت‌شده نیز نامیده می‌شود، برخی دیگر از این شیوه‌ها برای دسته‌بندی مؤخر یا هدایت‌نشده در دسترس هستند (Gupta, 2006) که در آن کلاس‌ها براساس داده‌های موجود تعیین شده‌اند. داده کاوی کاربرد جدیدی را در دسته‌بندی ایجاد کرده است، از زمانی که مجموعه‌های داده در داده کاوی افزایش یافته‌اند، روش‌های جدید دسته‌بندی به منظور سروکار داشتن با میلیون‌ها شیء توسعه یافته‌اند که شاید دارای ده‌ها یا صدها مشخصه باشند.

یک فرآیند دسته‌بندی که در آن کلاس‌ها از پیش تعریف شده‌اند، به روشی نیاز دارد که سیستم دسته‌بندی را به منظور تخصیص اشیاء به کلاس‌ها آموزش دهد؛ آموزشی که بر پایه یک آموزش نمونه است، در واقع به صورت مجموعه‌ای از داده‌ها در جایی که برای هر نمونه، یک کلاس از قبل تعریف شده است. فرض می‌شود کلاس هر شیء با داشتن برخی مشخصه‌ها و براساس هر کدام از آنها به ما گفته می‌شود. این مشخصه برای داده‌های آموزشی شناخته شده است، اما برای داده‌های دیگر به جز داده‌های آموزشی (ما این داده‌های دیگر را داده‌های آزمایشی می‌نامیم) فرض می‌شود که ارزش هر مشخصه شناخته نشده و باید با روش‌های دسته‌بندی مشخص گردد. این مشخصه ممکن است به عنوان خروجی تمام مشخصه‌های دیگر در نظر گرفته شود و معمولاً به عنوان مشخصه خروجی یا مشخصه وابسته بیان شده است. مشخصه‌های دیگر به جز مشخصه خروجی، مشخصه‌های داخلی یا مشخصه‌های مستقل نامیده می‌شوند (Wu & Su, 2005).

۳. درخت تصمیم

درخت تصمیم، روشی معروف برای دسته‌بندی است که نتایج آن در یک فلوچارت، شبیه ساختار درخت ارائه شده است که هر گره نشانگر یک تست بر روی ارزش مشخصه و هر شاخه، خروجی هر تست را نمایش می‌دهد، برگ‌های درخت نیز نمایانگر کلاس‌هاست. به‌طور عادی، پیچیدگی یک درخت تصمیم با افزایش تعداد مشخصه‌ها افزایش می‌یابد. اگرچه در بعضی از شرایط دیده شده است که تنها تعداد کمی از مشخصه‌ها می‌توانند کلاسی را که هر شیء به آن تعلق دارد، تعیین کنند و بقیه مشخصه‌ها کم یا بی‌تأثیرند (Tan & Yu, 2006).



می‌شوند. پیش‌بینی‌کننده‌ای که برای این جداکردن استفاده شده، معیار تعدیل دسته نامیده می‌شود.

- **معیار حداقل اندازه گره:** مشخص‌کننده این امر است که هر گره چند زیرگره می‌تواند داشته باشد. هرچه این مقدار بیشتر باشد، درخت کوچک‌تر خواهد شد.

- **معیار حد کثر خلوص:** می‌توان این معیار را انتخاب نکرد، اما با این انتخاب، مقدار آن باید بین صفر تا ۱۰۰٪ باشد. هرچه این مقدار بیشتر باشد، درخت بزرگ‌تر خواهد شد.

- **معیار حداکثر عمق:** مقدار آن باید عدد صحیحی بین صفر تا ۲۰ باشد. هرچه این مقدار بیشتر باشد، درخت بزرگ‌تر خواهد شد.

- **قطع درخت:** با این گزینه می‌توان تأثیر قطع درخت را بررسی کرد.

- **آموزش و آزمون داده‌ها:** مشخص‌کننده درصد داده‌هایی است که برای آموزش یا تست مدل استفاده می‌شود. اگر مجموعه آزمایشی به‌طور دستی انتخاب شود، ایده خوبی است برای بررسی اینکه تمام دسته‌های کلاس که در داده‌های آموزشی وجود دارد، مشابه داده‌های آزمایشی ایجاد شده باشند (Adda & Missaoui, 2007).

۴. طراحی معماری شبکه عصبی

در ساخت معماری شبکه عصبی، موارد زیر تعیین می‌شوند:

- تعداد ورودی‌ها؛
- تعداد لایه‌های پنهان شبکه؛
- پارامترهای یادگیری؛
- تعیین نوع انتخاب داده‌های آموزشی و آزمایشی؛
- سایز اولیه تابع وزنی؛
- تعداد سیکل‌های آموزشی داده‌ها؛
- افزودن تکانه: می‌توان قانون تغییر وزن‌ها را طوری در نظر گرفت که تغییر وزن در تکرار حدی، به اندازه تغییر وزن در تکرار قبلی بستگی داشته باشد. افزودن تکانه

3. Minimum Node Size Criterion
4. Maximum Purity Criterion
5. Maximum Depth Criterion
6. Pruning Option
7. Momentum

در ساخت درخت‌های تصمیم معمولاً داده‌ها را به دو دسته تقسیم می‌کنند:

- داده‌های آموزشی^۱ که برای ساخت مدل مورد استفاده قرار می‌گیرند؛

- داده‌های تست^۲ که برای تست و ارزیابی مدل ساخته‌شده کاربرد دارند.

کیفیت داده‌های آموزشی اغلب نقش مهمی در تعیین کیفیت درخت تصمیم دارد. در صورتی که آموزش سیستم زیاد شود یعنی داده‌هایی که برای آموزش و ساخت مدل به کار می‌رود درصد زیادی از داده‌ها باشد، دچار حالتی به نام «آموزش بیش از حد مدل» خواهیم شد که به دلیل وجود موارد غیرعادی در داده‌های آموزشی، خطا تولید می‌کند (Chan & Lewis, 2002).

اندازه‌گیری کیفیت یک درخت تصمیم مسئله مهمی است. دقت تعیین‌شده دسته‌بندی با استفاده از داده‌های آزمایشی به‌طور آشکار یک شاخص مطلوب است. اما ممکن است شاخص‌های دیگری نیز برای اندازه‌گیری به کار روند. این شاخص‌ها شامل هزینه متوسط و هزینه بدترین مورد در دسته‌بندی یک شیء است. یک درخت تصمیم ممکن است قادر به دسته‌بندی داده‌های آموزشی با دقت ۱۰۰ درصد باشد، اما بدین معنا نیست که درخت، روی داده‌های آزمایشی که بخشی از مجموعه آموزشی نبوده، بسیار دقیق باشد. در واقع ممکن است در صورتی که داده‌های آموزشی نمونه خوبی از گروه داده‌ها باشند، عملکرد درخت بر روی داده‌های آزمایشی نیز تا ۱۰۰ درصد بالا بیاید. در نتیجه ورودی‌های آموزشی اهمیت بالایی دارند. به همین دلیل معیارهایی را برای این ورودی‌ها در نظر می‌گیریم تا با رعایت آن بتوانیم نتایج بهتری را به دست آوریم. در ادامه چند مورد از این معیارها را بیان می‌کنیم:

- **معیار تعدیل دسته‌های یک پیش‌بینی‌کننده دسته‌ای:** زمانی که درخت رشد می‌کند، گره‌های جزئی (فرعی) با جداکردن گره‌های اصلی ساخته

1. Train Data
2. Test Data

باشد. باتوجه به توضیحات روش ارائه شده جهت اجرای نرم‌افزار دسته‌بندی با استفاده از درخت تصمیم در فصل دوم، خلاصه کار به این صورت است:

- ورود داده‌ها؛

- ساخت مدل؛

- دریافت نتایج مدل‌سازی؛

- تولید قوانین.

خروجی اولیه درخت تصمیم شامل گره‌های به‌دست آمده به‌صورت نمودار ۱ است. نتایج درخت تصمیم به‌صورت جزئی و تفکیک شده در مدل دسته‌بندی در نمودار ۲ آورده شده است.

براساس این مدل نهایی، خروجی نرم‌افزار CTree، تعداد کل داده‌های آزمایشی و داده‌های آموزشی ارائه شده است. تعداد پیش‌بینی‌کننده‌ها به تعداد ورودی‌های نرم‌افزار هستند که برابر ۶ است. تحقیق شامل سه کلاس A، B و C به ترتیب بیمه عمر و پس‌انداز، بیمه حوادث خانواده و بیمه حوادث فردی است. کلاس حداکثر، A بوده و از این کلاس به‌عنوان پیش‌بینی‌کننده استفاده شده است که دسته‌بندی ۳۲ درصد خطای دسته‌بندی را دارد. به این معنا که مدل به میزان ۳۲ درصد کلاس را به اشتباه پیش‌بینی کرده است. باتوجه به اینکه درصد اشتباه پایین بوده است، رکوردها تا حدود ۷۰ درصد، کلاس‌های خروجی پیش‌بینی شده را پشتیبانی می‌کنند. معیار کلی پایین یا بالا بودن به‌طور معمول ۵۰ درصد محسوب می‌شود که در تحقیق‌های کاربردی درخت تصمیم نیز از همین ملاک استفاده شده است.

ماتریس نهایی شامل اطلاعات داده‌های آزمایشی و آموزشی است. بدین‌صورت که تعداد موجود در هر یک از کلاس‌ها را مشخص نموده و در خانه‌های ماتریس قرار می‌دهد. محور افقی ماتریس، نشانگر کلاس پیش‌بینی‌کننده و محور عمودی، معرف کلاس صحیح هر یک از رکوردهاست. مثلاً کلاس A به‌عنوان پیش‌بینی‌کننده دارای ۱۵۸۳ داده آموزشی است، کلاس B شامل ۵۷۵ رکورد صحیح در B و کلاس C

باعث می‌شود تا با حرکت در مسیر قبلی در سطح خطا از گرفتار شدن در مینیمم محلی و قرار گرفتن در سطوح صاف پرهیز شود و با افزایش تدریجی مقدار پله تغییرات، سرعت جستجو افزایش یابد (Belhadjali & Whaley, 2004).

۵. دریافت نتایج مدل‌سازی

خروجی اول، یک مدل شبکه عصبی است که عمدتاً مجموعه‌ای از وزن‌ها بین لایه‌های شبکه است. بعد از اتمام اجرا، باید مجموعه نهایی وزن‌ها نیز در صفحه محاسبه ذخیره شود. خروجی دیگر مدل، محاسبه مقادیر خطای نسبی مطلق^۱ و میانگین مربعات خطا^۲ بر روی مجموعه داده‌های آموزشی می‌باشد. نهایتاً اعتبارسنجی صورت گرفته روی مدل ارائه خواهد شد. میانگین مربعات خطا معیار مناسبی برای تعیین این امر است که مدل ایجاد شده تا چه حدی می‌تواند واقعیات را پیش‌بینی کند. در واقع خروجی نهایی مدل‌سازی، وزن‌های تعیین شده برای یال‌ها و نیز میزان خطایی است که این وزن‌ها در پیش‌بینی دارند (Sivanandam & Sumathi, 2006).

۶. تجزیه و تحلیل نتایج درخت تصمیم

درخت تصمیم، روشی شناخته شده به‌منظور دسته‌بندی است که نتایج آن در یک فلوچارت شبیه ساختار درخت ارائه می‌شود (بدین‌صورت که هر گره نشانگر یک آزمون بر روی ارزش مشخصه و هر شاخه، خروجی آن را نمایش می‌دهد، برگ‌های درخت نیز نمایانگر کلاس‌هاست).

برای تحلیل داده‌ها از نرم‌افزار قدرتمند Ctree استفاده می‌کنیم، که ابزار و امکانات زیادی برای ایجاد درخت تصمیم و بررسی نتایج آن و تحلیل در اختیار ما قرار می‌دهد. همان‌طور که بیان شد، نیاز است بین تعداد نمونه‌های آزمایشی و تعداد مشخصه‌های مستقل توازن باشد. به‌طور کلی اگر تعداد مشخصه‌های مستقل کم باشد، بهتر است تعداد نمونه‌های آموزشی مورد نیاز کم باشد و به‌همین ترتیب زمانی که تعداد مشخصه‌ها زیاد است، بهتر است تعداد نمونه‌های آموزشی مورد نیاز نیز زیاد

1. Absolute Relative Error
2. Mean Squared Error

شامل تمام ۵۷۲ رکورد کلاس خود و ۲ رکورد متعلق به کلاس B است. در داده‌های آزمایشی نیز کلاس A شامل ۵۰۵ رکورد صحیح، کلاس B شامل ۱۷۰ رکورد صحیح در B و ۳ رکورد در کلاس C و کلاس C شامل تمام ۱۹۶ رکورد صحیح است. در نتیجه تنها ۵ رکورد به‌طور کلی به اشتباه در کلاس اولیه خود حضور داشته‌اند و دیگر رکوردها در خروجی صحیح خود قرار دارند. به‌عبارت‌دیگر خطای دسته‌بندی اشتباه در داده‌های آموزشی ۰/۰۷ درصد و در داده‌های آزمایشی ۰/۳۴ درصد است که قابل توجه نیستند. همچنین درخت تصمیم نیز شامل ۳۳ گره، ۱۹ برگ و ۱۰ سطح است. جدول خلاصه قوانین، کیفیت قوانین واحد را ارائه داده است که این کیفیت توسط سه روش متفاوت زیر اندازه‌گیری می‌شود:

- **پشتیبانی^۱**: درصدی از داده‌های آموزشی که سمت چپ قانون برای آنها صدق می‌کند. اگر در یک مشاهده، قاعده سمت چپ^۲ قانون صدق کند، می‌گوییم که قانون برای آن مشاهده به‌کار می‌رود. این شاخص بیان می‌کند که قانون مورد نظر تا چه حد قابل کاربرد است.

اطمینان^۳: درصد رکوردهای آزمایشی است که قاعده سمت چپ در آنها صادق است و علاوه بر آن قاعده سمت راست^۴ نیز در آنها صدق می‌کند. به‌عبارت‌دیگر، این شاخص نشان می‌دهد قانون برای چه درصدی از مشاهداتی که قانون در آنها به‌کار می‌رود، صدق می‌کند. این شاخص صحت قوانین را بررسی می‌کند.

- **دریافت^۵**: چه درصدی از رکوردهای مورد نظر به‌طور صحیح توسط یک قانون دریافت می‌شوند، در واقع نوعی انعکاس از ساختار مسئله است. اگر یک قانون با دریافت نزدیک به ۱۰۰ درصد وجود داشته باشد، به این معنی است که در فضای پیش‌بینی‌کننده، تمام مشاهدات با این کلاس به‌طور نزدیکی کنار هم قرار گرفته‌اند و قانون توانسته است آن بخش از فضای پیش‌بینی‌کننده را

به‌خوبی دریافت کند.

برای قضاوت در رابطه با کیفیت قانون، هر سه معیار باید در نظر گرفته شوند. از طرفی نمی‌توان گفت یک قانون به‌عنوان مثال با ۱۰۰ درصد اطمینان و پشتیبانی ۵ درصد، لزوماً ضعیف است، چرا که ممکن است قانون پیش‌بینی‌کننده در واقع تنها ۵ درصد از زمان‌ها در داده‌های آموزشی رخ دهد و همان مقدار را قانون به‌طور مطلوبی دریافت کرده باشد و کیفیت خوبی داشته باشد. کیفیت این قوانین به‌صورت بصری در نمودار ۳ ارائه‌شده که محور افقی آن، تعداد رکوردهای مشاهدات و محور عمودی، قوانین تولیدشده است.

هر نقطه در این نمودار نشانگر یک مشاهده و رنگ آن نشان‌دهنده کلاس صحیح آن است. رنگ شماره ۱ کلاس A، رنگ شماره ۳ کلاس B و رنگ شماره ۲ کلاس C است. اولین خط از پایین قانون ۱ را نشان می‌دهد که این ردیف با توجه به کلاس پیش‌بینی‌شده توسط این قانون رنگ خاصی می‌گیرد.

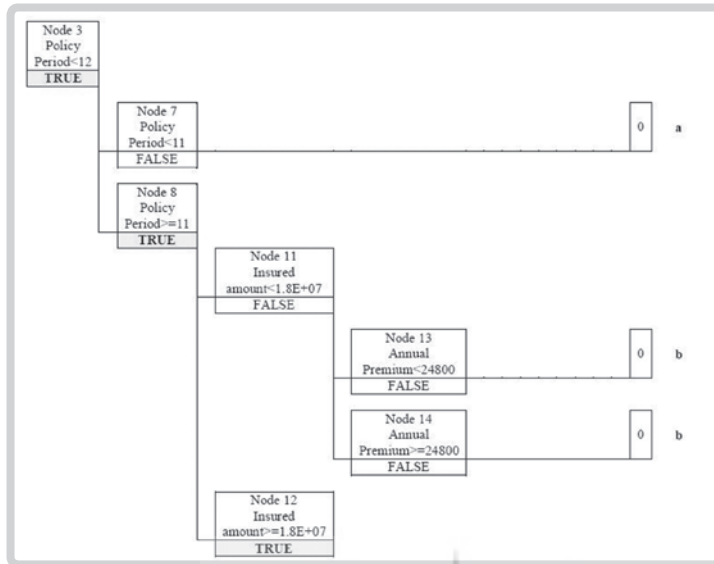
تعداد نقطه‌ها در یک ردیف، میزان پشتیبانی یک قانون را می‌رساند که در اینجا تنها قوانین ۷ تا ۱۰، پشتیبانی کمی داشته یا به عبارتی کمتر رخ می‌دهند.

اگر هر یک از نقاط را در ردیف قوانین با نقاط پایین‌ترین ردیف مقایسه کنیم و رنگ یکسانی را داشته باشیم، قانون، کلاس را به‌درستی پیش‌بینی کرده است. در اینجا ۸ قانون با توجه به شکل، کلاس‌ها را به‌درستی پیش‌بینی کرده‌اند که قابل توجه است. در عین حال قانون ۱۴ و ۱۵ پیش‌بینی درستی نداشته و درصد اطمینان پایینی را کسب کرده‌اند.

به‌همین ترتیب آن میزان از رنگ قانون که به رنگ کلاس صحیح می‌باشد، در واقع میزان درصدی از کلاس صحیح است که توسط قانون مورد نظر پیش‌بینی شده است. با توجه به شکل، قانون‌های اول تا چهارم دریافت بالایی داشته‌اند.

1. Support
2. Left Hand Side
3. Confidence
4. Right Hand Side
5. Capture

نمودار ۱. درخت تصمیم حاصل از تحلیل داده توسط نرم‌افزار CTree



نمودار ۲. خلاصه اطلاعات حاصل شده از تحلیل داده توسط نرم‌افزار CTree

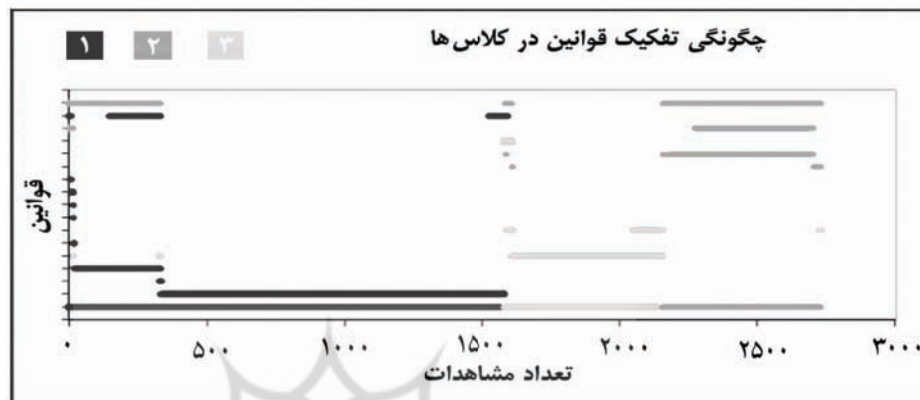
Classification Tree Model		Tree Information	
Number of Training observations	2732	Total Number of Nodes	33
Number of Test observations	874	Number of Leaf Nodes	19
Number of Predictors	6	Number of Levels	10
Class Variable	Class	% Missclassified	
Number of Classes	3	On Training Data	0.07%
Majority Class	a	On Test Data	0.34%
% Missclassified if Majority Class is used as Predicted Class	32%	Time Taken	
		Data Processing	56 Sec
		Tree Growing	3 Min : 36 Sec
		Tree Pruning	0 Sec
		Tree Drawing	1 Sec
		Classification using final tree	1 Min : 5 Sec
		Rule Generation	24 Sec
		Total	6 Min : 3 Sec

Confusion Matrix								
Training Data	Predicted Class			Test Data	Predicted Class			
	a	b	c		a	b	c	
True Class a	1583			True Class a	505			505
True Class b		575		True Class b		170	3	173
True Class c		2	572	True Class c			196	196
	1583	577	572		505	170	199	874

۷. فرضیه شبکه مصنوعی

صحت نتایج می‌توان براساس آن، مشخصه‌های مشتریان جدید را وارد مدل کرده و در نهایت دسته مورد نظر مشتری را پیش‌بینی کرد. آزمون شبکه‌های عصبی ابتدا با در نظر گرفتن داده‌های آموزشی و آزمایشی و بررسی انطباق آنها در طول دوره آموزشی شبکه عصبی مصنوعی استفاده از شبکه‌های عصبی مصنوعی، مشتریان بالقوه محصولات جدید بیمه‌ای و خدمات متناسب آنها شناسایی می‌شوند. با استفاده از روش شبکه‌های عصبی مصنوعی، مدل دسته‌بندی مشتریان ارائه شده و پس از اطمینان از

نمودار ۳. نمایش بصری کیفیت داده‌ها توسط نرم‌افزار CTree



نمودار ۴. نمایش اطلاعات شبکه عصبی حاصل شده توسط نرم‌افزار R

Neural Network Model for Classification

% MissClass.(Training) 14.51% % MissClass.(Validation) 13.47%

Number of Hidden Layers 2
 Layer Sizes 11 2 2 3

True Output (if available) A پژوهشگاه علوم انسانی و مطالعات فرهنگی

Model Output

نمودار ۵. نمایش داده‌های تبدیل شده برای شبکه عصبی با استفاده از نرم‌افزار R

Transformed Input	Bias	Income level.lo w	Income level.avera ge	Income level.ve ry low	Income level.hi gh	Income level.ve ry high	Gender.m ale
Hdn1_bias	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
Hdn1_Nrn1	-15.5528	-3.1816	-2.0158	-7.7086	-1.2370	-1.7267	-6.6269
Hdn1_Nrn2	-3.8321	8.3451	8.3627	35.828	7.3738	8.1817	-2.6685
	1.0000	#N/A	#N/A				
Hdn2_bias	0.0000	0.0000	0.0000	0.0000			
Hdn2_Nrn1	-4.1095	-1.4429	10.4773	#N/A			
Hdn2_Nrn2	2.6959	-5.1714	7.6796	#N/A			

و سپس مقایسه نتایج آن با نتایج روش درخت تصمیم انجام می‌شود.

۸. تجزیه و تحلیل نتایج شبکه‌های عصبی

همان‌طور که بیان شد، اجزای یک شبکه عصبی عبارت‌اند از:

- ورودی‌ها: ورودی‌ها می‌توانند خروجی سایر لایه‌ها بوده یا این که به حالت خام در اولین لایه و به صورت‌های داده‌های عددی و رقمی، متون ادبی، فنی و تصویر یا شکل باشند.

- وزن‌ها: میزان تأثیر ورودی بر خروجی توسط وزن اندازه‌گیری می‌شود.

- تابع جمع: خروجی مسئله را تا حدودی مشخص می‌کند.

- تابع تبدیل: این تابع توسط طراح مسئله انتخاب

می‌گردد و براساس انتخاب الگوریتم یادگیری، پارامترها

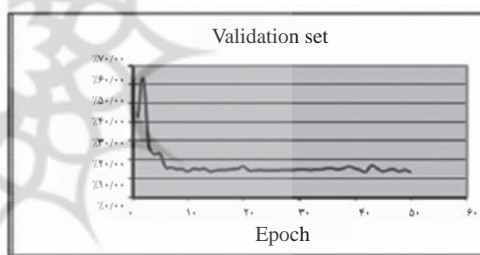
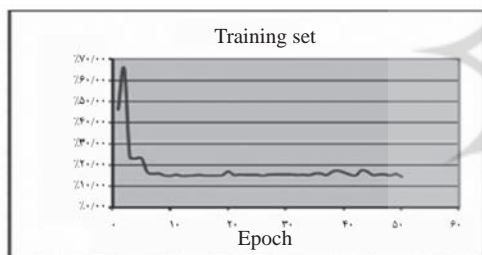
تنظیم می‌گردند.

- خروجی: منظور از خروجی، پاسخ مسئله است.

برای تحلیل شبکه‌های عصبی، ابزارهای گوناگونی وجود دارد که از نظر قدرت و امکانات متفاوت‌اند. در اینجا ما از نرم‌افزار R استفاده می‌کنیم که گزارش‌های بسیار کارایی از تحلیل داده‌ها در اختیار ما قرار می‌دهد و نیز قابلیت‌های فراوانی دارد. براساس ورودی‌های تعیین شده و معیارهای انتخاب‌شده در معماری شبکه عصبی و اجرای نرم‌افزار، مدل دسته‌بندی ساخته شده با استفاده از شبکه‌های عصبی به صورت نمودار ۴ است.

همان‌طور که مشاهده می‌شود درصد خطای دسته‌بندی نادرست ۱۴ و ۱۳ درصد در داده‌های آموزشی و آزمایشی است که مطابق با ۱۴ درصد خطای مشابه در خروجی نهایی درخت تصمیم است.

نمودار ۶. نمایش شیب و عرض از مبدأ برای داده‌های آموزشی و آزمایشی با استفاده از نرم‌افزار R



۹. نتیجه‌گیری از خروجی‌های درخت تصمیم و شبکه‌های عصبی

همان‌طور که خروجی‌های نرم‌افزار درخت تصمیم و شبکه‌های عصبی یکدیگر را تأیید کردند و در راستای یکدیگر قرار گرفتند، نتایجی در رابطه با مشخصه‌های این پژوهش (متغیرهای ورودی) و بررسی قوانین تصمیم‌گیری تولیدشده به دست آمده که توسط خبرگان این صنعت نیز تأیید شده است. نتایج به دست آمده از هر دو روش به کار گرفته شده، طی جلسه‌ای به مدیر بیمه درمان و کارشناس ارشد بیمه عمر و درمان شرکت بیمه مورد نظر ارائه و نظرات بیان شده مستند گردید.

۱۰. نتیجه‌گیری

نتایج به دست آمده از این تحقیق بیان می‌کنند که فضای خالی بین تئوری و اجرا هنوز باید کمتر شود، به‌ویژه در ترغیب شرکت‌های بیمه به منظور تغییر دادن مفاهیم استفاده از داده‌کاوی در جایی که اطلاعات امنیتی نیز درگیر شده‌اند. این تحقیق به این دلیل که شواهد کاربردی موجب جلب توجه به نقش مهم داده‌کاوی به کاررفته در صنعت بیمه می‌شود، بسیار با ارزش است. همچنین ممکن است مسائل دیگری هم از قبیل در نظر گرفتن دیگر مشخصه‌های مشتریان (مانند شغل، نوع زندگی، تحصیلات و ...) و طراحی کردن تمام مراحل برای یک توسعه پیچیده (مانند ایجاد ماجول^۱ و به روز کردن نتایج کشف شده در هر زمان) مانند توسعه نرم‌افزاری مهم باشند.

با اهمیت یافتن بازاریابی فردی جهت پاسخ به مشتری، داده‌کاوی به صورت مؤثر در دسته‌بندی مشتریان هدف استفاده شده است. بازاریابی فردی بر این مطلب تأکید می‌کند که استراتژی بازاریابی باید بر مشخصات فردی مشتری متمرکز شود. مطابق با همین موضوع، مدل تجربه‌گرا که بیشتر شرکت‌های بیمه از آن استفاده می‌کنند، به دلیل اینکه باید بسیاری از مشتریان و همچنین بسیاری از خصوصیات آنها در نظر گرفته شود، دیگر برای اجرای امور پیچیده در دنیای امروز کارآمد نخواهد بود.

در این مسئله، داده‌های رکوردی مختلف را به شبکه داده و نام گروه هر رکورد را به عنوان خروجی مشخص می‌کنیم، پس از آموزش، شبکه توانسته است با دریافت داده‌های مربوط به نمونه‌های جدید مشخص کند که این نمونه به کدام دسته متعلق است. حال بخشی از داده‌های تبدیل شده در نرم‌افزار با دو لایه پنهان در شبکه به صورت نمودار ۵ است:

وجود Bias در تابع تبدیل که به صورت تابع مقابل است:

$$y=f(n)=f(wx+b)$$

W باید به نوعی کارایی داشته باشد که با ضرب در X مقدار را مشخص کند، اما در واقعیت نمی‌تواند، به همین دلیل مقداری چولگی دارد که باید به آن اضافه شود. مانند ترازویی که در یک مقدار ثابت مشکل دارد و باید به یک طرف آن وزنه‌ای را افزود. در این خروجی میزان b صفر است و وزن‌ها تا حد زیادی نزدیک به صحیح انتخاب شده‌اند. اعداد منفی در این جدول نشان‌دهنده حرکت در جهت کاهش شیب است که در نمودار ۵ قابل مشاهده است.

شیب و عرض از مبدأ تابع تبدیل به میزان تمام داده‌های آموزشی و آزمایشی در نمودار ۶ فراهم آمده است:

با توجه با اینکه دوره آموزشی در معماری شبکه عصبی این تحقیق ۵۰ انتخاب شده است، محور افقی، شامل ۵۰ آموزش شبکه است و محور عمودی، میزان خطای دسته‌بندی نادرست را در هر س کل نشان می‌دهد. همان‌طور که در نمودار مشاهده می‌شود اختلاف بسیار ناچیزی میان دوره‌های آموزش داده‌های آزمایشی و داده‌های آموزشی وجود داشته که تا حد زیادی براساس گفته‌های قبلی در روش تحقیق، اعتبار تحقیق را تأیید می‌کند. برای بررسی دقیق‌تر می‌توان درصدها را از روی جدول با مقادیر عددی مقایسه کرد. همچنین هرچه شبکه بیشتر آموزش می‌بیند، میزان خطای دسته‌بندی اشتباه کاهش می‌یابد، به طوری که در دوره‌های پایانی به ۱۳ تا ۱۴ درصد رسیده و شبکه موفق شده است تا تابع هدف را بیاموزد.

با رشد سریع در حجم داده‌هایی که سیستم‌های اطلاعاتی جمع‌آوری می‌کنند، بازاریاب‌های بیشتری تمایل به استفاده از ابزارهای پشتیبانی تصمیم‌گیری مبتنی بر داده برای بالابردن و بهبود کارایی و مؤثر بودن تصمیمات بازاریابی خود خواهند داشت.

مشتریان و میزان ریسک آنان، عامل اصلی در سودآوری شرکت‌های بیمه است و شناخت آنها مهم‌ترین مسئله در این مورد است. استفاده از ابزارهای دسته‌بندی به منظور جداسازی مشتریان و شناخت آنها بسیار مفید است. داده‌کاوی، ابزار بسیار هوشمندی به منظور دسته‌بندی مشتریان است. با استفاده از دسته‌بندی مشتریان می‌توان مشتریان بیمه را بهتر شناخت و براساس گروهی که مشتری با توجه به خصیصه‌هایش در آن قرار گرفته است سیاست‌گذاری‌های مناسبی برای آنها در نظر گرفت. در این پژوهش با استفاده از روش‌های درخت تصمیم و شبکه عصبی، مشتریان بیمه را دسته‌بندی کرده‌ایم و سپس نتایج این دو روش را مقایسه کرده‌ایم. همان‌طور که پیشتر ذکر شد نتایج به دست آمده از این دو روش یکدیگر را تأیید می‌کنند که این امر بر صحت نتایج می‌افزاید. همچنین از نظرات کارشناسان برای ارزیابی نتایج استفاده کرده‌ایم. نظرات کارشناسان نیز نتایج به دست آمده از روش‌های داده‌کاوی را تأیید می‌کند. لذا با استفاده از داده‌کاوی توانسته‌ایم مشتریان بیمه را در دسته‌هایی قرار داده و با توجه به ویژگی‌های این دسته‌ها شناخت بهتری از مشتریان به دست آوریم. همچنین می‌توان مشتریان آتی را با توجه به ویژگی‌هایشان در یکی از این دسته‌ها قرار داد. دسته‌ای که مشتری در آن قرار گرفته، نشان‌دهنده خصایص مشتری و میزان ریسک آن برای شرکت بیمه است. بدین صورت می‌توان در مورد میزان ریسک مشتری برای شرکت بیمه پیش‌بینی داشته باشیم که این امر می‌تواند کمک شایانی در سودآوری شرکت‌های بیمه داشته باشد.

۱. پیشنهاد برای تحقیقات آتی

مباحث این تحقیق عمدتاً بر سه نکته متمرکز شده است که شامل تجربه کاربرد داده‌کاوی در دسته‌بندی

مشتریان هدف برای صنعت بیمه، قابلیت استفاده از داده‌کاوی در پشتیبانی از عملیات تصمیم‌گیری و جهشی از مباحث تئوری به سمت نتایج عملی است (Adda & Missoui, 2007). فناوری داده‌کاوی عمدتاً یک نمونه از پایگاه داده است که به‌عنوان ابزار مدیریتی در پشتیبانی تصمیم‌گیری استفاده شده است، اگرچه سبک تصمیم‌گیری تجربه‌گرا در حقیقت برای یک مدت طولانی به کاررفته و تا نقش مهمی در حمایت از مدیریت بازاری کند، ادامه پیدا خواهد کرد. همچنین با افزایش اطلاعات، سبک پشتیبانی تصمیم به‌طور تصاعدی از تجربه‌گرایی به اطلاعات‌گرایی تغییر کرده است.

هرچند که این تحقیق تنها در سطح یکی از واحدهای شرکت بوده، اما دانش کشف‌شده، توجه زیادی را در شرکت مورد نظر در جلسه با مدیر بیمه عمر و درمان جذب کرده است. این امر بدین دلیل است که نتایج نشان داده‌اند در شرکت مورد نظر، هرگز به رابطه بین محصولات و مشتریان توجه نکرده‌اند. قابلیت استفاده از دانش کشف‌شده و نمایش داده‌شده در این تحقیق، نشان داده است که در این مرحله، داده‌کاوی می‌تواند یک راه ممکن به سمت افزایش دسته‌بندی مشتریان هدف باشد. برای تأثیر دانش کشف‌شده در مرحله دوم، به‌منظور جمع‌آوری و مقایسه داده‌ها برای شرکت مورد نظر، زمان زیادی مورد نیاز است؛ برای مثال، آیا دانش کشف‌شده به‌طور قابل ملاحظه‌ای به افزایش پاسخ‌های مشتری برای سال بعد (دو، سه یا چهار سال بعد) کمک می‌کند، این بخش می‌تواند یکی از مهم‌ترین تحقیقات آینده باشد.

منابع

1. Adda, M & Missaoui, R 2007, 'Relation rule mining', *International Journal of Parallel, Emergent and Distributed Systems*, vol.22, no.6, pp. 439- 49.
2. Belhadjali, M & Whaley, GL 2004, 'A data mining approach to neural network training', *Information Management & Computer Security*, vol. 12, no. 1, pp. 117-24.
3. Berson, A; Smith, S & Therling, K 2001, *Building data mining applications for CRM*, Mc GrawHill.
4. Brockett, P & Xiaohua, X 1999, 'Operations research in insurance, A review', *Journal of Mathematics and Economics*, vol. 19, no. 2, p. 154.
5. Chalmeta, R 2005, 'Methodology for customer relationship management', *The Journal of System and Software*, vol. 15, pp. 192-201.
6. Chan, C & Lewis, B 2002, 'A basic primer on data mining', *Information Systems Management*, vol. 19, no. 4, pp. 56-60.
7. Chen, YH & Su, CT 2006, 'A kanoCKM model for customer knowledge discovery', *Total Quality Management & Business Excellence*, vol. 17, no.5, pp.589-608.
8. Clifton, C & Thuraisingham, B 2001, 'Emerging standards for data mining', *Computer Standards and Interfaces*, vol. 23, pp. 187-93.
9. Gupta, GK 2006, *Data mining with case studies*, New Delhi, Prentice Hall India.
10. Sivanandam, SN & Sumathi, S 2006, *Neural networks using MATLAB 6.0*, McGrawHill.
11. SY, BK 2001, 'Information statistical pattern based approach for data mining', *Journal of Statistical Computation and Simulation*, vol. 69, no. 2.
12. Tan, KC & Yu, Q 2006, 'A coevolutionary algorithm for rules discovery in data mining', *International Journal of Systems Science*, vol. 12, no. 37, pp. 835-64.
13. Wu, Ch & Su, Y 2005, 'Targeting customers via discovery knowledge for the insurance industry', *Expert Systems with Applications*, vol. 29, pp. 291 -99.