

## استفاده از یک مدل داده کاوی برای پیش بینی رفتار مشتریان مشاغل خرد در بنگاه های اقتصادی و تعیین ریسک اعتباری در اعطای تسهیلات بانکی

سمیه مهاجری

کارشناس ارشد مهندسی کامپیوتر گرایش نرم افزار - Somayyeh\_Mohajeri@yahoo.com

### خلاصه

هدف بالقوه هر سازمان مالی حفظ مشتریان موجود و رسیدن به مشتریان جدید در بلند مدت است. رفتار اقتصادی مشتری و ماهیت سازمان با یک توصیفی با نام "شناخت مشتریان خود (KYC)"<sup>۱</sup> در نظام بانکداری کنترل می شود. زمانی که بانک ها قصد دارند به مشتریان وام یا تسهیلاتی پرداخت کنند، وضع اعتبار دارایی ها و خوش حسابی یا بدحسابی مشتری مورد نظر را می سنجند. مشتریان سپرده گذار در برخی از بخش ها با ریسک بالا هستند؛ در حالی که در برخی از بخش ها خطر متوسط و ریسک بسیار پائینی دارند. در حال حاضر ریسک اعتباری را می توان بطور گسترده تحت عوامل کمی و کیفی طبقه بندی کرد. اگر چه سیستم های موجود بسیاری در حفظ مشتری در بانک نهادینه شده است اما روش های سختگیرانه و بدون رویکرد روشن و تعریف شده ای برای پرداخت تسهیلات در بخش بنگاه اقتصادی های وجود دارد. در این مقاله از اطلاعات مشتریان مشاغل خرد بنگاه های اقتصادی در یکی از شعب بانک های ایران برای تجزیه و تحلیل و عوامل عمده تراکنش های مشتریان و پیش بینی رفتار ایشان در پرداخت تسهیلات بانکی استفاده شده است و بهره گیری از روش های داده کاوی در شعب بانک برای موضوع فوق و ساده تر شدن فعالیت های سنجش و گزینش درخواست تسهیلات مشتریان، از میان تعداد زیاد درخواست تسهیلات به عنوان راه حل مناسب اعطای تسهیلات به مشتریان بنگاه های اقتصادی پیشنهاد شده است و از درخت تصمیم گیری و الگوریتم C4.5 در مدل استفاده شده است و در نهایت عملکرد خود را با نتایج MATLAB مورد آزمایش قرار داده ایم.

کلمات کلیدی: داده کاوی، درخت تصمیم، مشتری، بانک، تسهیلات

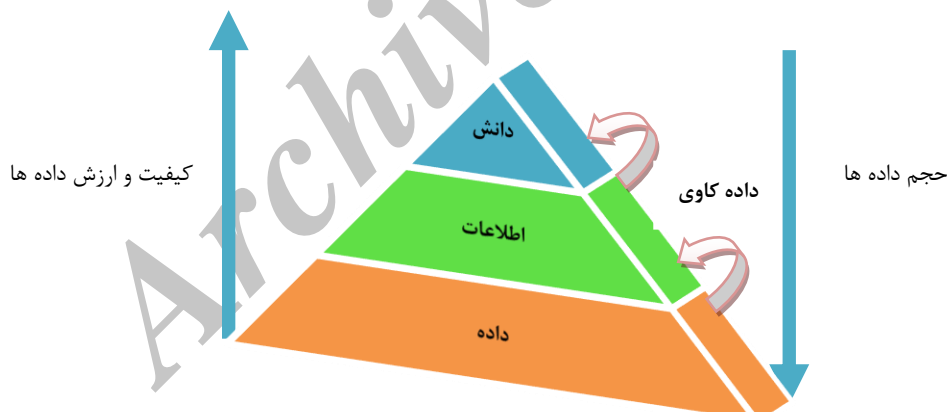
<sup>۱</sup> Know Your Customer

۱. مقدمه

وقوع رویدادهایی نظیر توسعه رقابت جهانی، پیشرفت فناوری اطلاعات و دسترسی به سیستم‌های اطلاعاتی ارزان در سالهای اخیر و نیز تلاش واحدهای اقتصادی در جهت احراز رتبه جهانی و ورود به بازارهای بین‌المللی، داشتن نگرشهایی همچون رضایت مشتریان، مدیریت اطلاعات را اجتناب‌ناپذیر کرده است. آنچه مسلم است موفقیت و تداوم فعالیت در محیط رقابتی جدید، مستلزم استفاده از روشه‌ای نوینی است که کسب‌وکار را در رده جهانی قرار می‌دهد که در این میان بانک‌ها نیز به عنوان یک بنگاه اقتصادی از این امر مستثنی نبوده و در سال‌های اخیر تعداد بانک‌ها و موسسات مالی که خدمات بانکی و محصولات بانکداری الکترونیکی را معرفی و گسترش می‌دهند افزایش یافته است.

اساس و پایه استراتژی این موسسات مالی توانایی در حفظ مشتریان موجود و رسیدن به مشتریان جدید و آینده‌نگر است. صنعت بانکداری در ایران به سرعت در حال رشد است و این در حالی است که همگام با رشد صنعت بانکداری از طریق پیشرفت‌های تکنولوژیکی، ارائه ایده‌های جدید برای بازار بین بانکی برای حفظ آن به ضروریات تبدیل شده است و لاجرم پرتفوی محصولات ارائه شده توسط بانک‌ها متنوع شده است. با توجه به این که موضوع تجهیز منابع و تخصیص آن به عنوان خط اصلی فعالیت بانکها شناخته شده به طوریکه حتی در تعریف بانکها این مهم گنجانده شده است؛ لذا به نظر می‌رسد هر چه از عمر یک بانک می‌گذرد می‌بایست منابع بیشتری جذب و به همان نسبت تسهیلات بیشتری پرداخت نموده و مقداری از منابع را نیز مصرف دارائی‌های غیر تسهیلاتی کرده باشد. در طول سال‌های اخیر، جذب مشتریان بیشتر از پیش در رسیدن به این هدف مهم شده است؛ تجمع داده‌های عملیاتی به ناچار از این رشد در این صنعت منتج می‌شود.

نیاز فزاینده‌ای برای تبدیل داده‌ها به اطلاعات مفید و باارزش یکسان سازی شده به منظور حرکت رو به جلو و به دست آوردن یک مزیت رقابتی وجود دارد. داده کاوی<sup>۱</sup> نقش مهمی در این تلاش‌ها دارد.



شکل ۱: نقش داده کاوی در تبدیل داده‌ها

در این مقاله با استفاده از مدل درخت‌های تصمیم<sup>۲</sup> به عنوان یکی مدل‌های داده کاوی، سعی شده است اطلاعات ارزشمندی از مشتریان تسهیلات بانکی با توجه به اینکه پیدا کردن مشتریان خوب برای پرداخت تسهیلات واقعا یک مسئله چالش برانگیز در عصر بانکی است، استخراج شود. در مدل درخت‌های تصمیم از الگوریتم C4.5 برای استخراج اطلاعات مشتریان دارای ریسک اعتباری کمتر و سودمندتر برای پرداخت تسهیلات بانکی استفاده شده است.

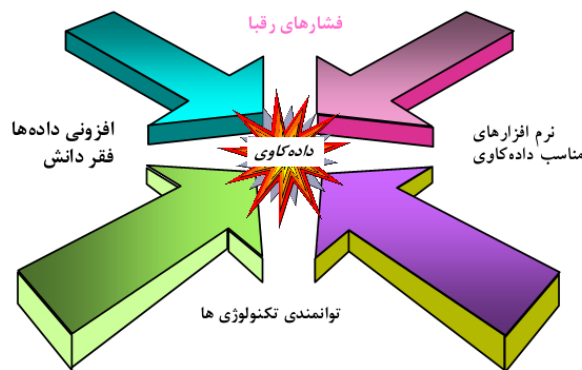
<sup>۱</sup>Data Mining

<sup>۲</sup>Decision Trees

## ۲. داده کاوی و نقش آن در بانک ها

در بسیاری از رشته های علمی و فنی، در نهایت ما با مجموعه ای از داده ها روبرو هستیم که حجم کم یا زیادی را دارند؛ اما مهم ترین کار، به دست آوردن چنین پایگاه داده ای نیست. بلکه باید بتوانیم، سطح بالاتری از دانش را با توجه به پایگاه داده مذکور به دست بیاوریم؛ این یعنی نتیجه گیری و جمع بندی تمام تلاش هایی که برای جمع آوری آن داده ها صرف شده است.

داده کاوی فرآیند تکرار شونده است که ترکیبی از دانش بنگاه اقتصادی، روش های یادگیری ماشین و ابزار و مقادیر زیادی از اطلاعات دقیق و مربوط به توانایی کشف بینش غیر بصری پنهان در اطلاعات حقوقی سازمان است. این اطلاعات می تواند اصلاح فرآیندهای موجود و کشف روندها بوده و به تدوین سیاست های مربوط به رابطه این شرکت با مشتریان و کارکنانش کمک کند. در بانک ها و موسسات مالی، داده کاوی به طور موفقیت آمیزی در تعیین کاندیدهای واجد شرایط احتمالی برای پرداخت تسهیلات بانکی، پیدا کردن مشتریان سودآور، محصولات بانکی، مشخص کردن عوامل ارائه خدمات مختلف استفاده شده است (Hu, ۲۰۰۵). تمامی این عوامل، بهبود روش های قدیمی انجام امور و وادار کردن بانک ها به در نظر گرفتن بازنگری فعالیت های خود را برای افزایش سهم از بازار در بر دارد.



شکل ۲: نقش داده کاوی در رقابت بین بانکی

تحقیقات متعددی در زمینه تجزیه و تحلیل برخورد و حفظ مشتری در بخش های بانکی ساخته شده است. برخی مطالعات نشان می دهد که مهم ترین متغیرهای موثر بر انتخاب مشتری خدمات ارائه شده موثر و کارآمد، سرعت و کیفیت خدمات، انواع خدمات عرضه شده، هزینه های خدمات الکترونیکی پایین، امکانات بانکداری آنلاین، ایمنی وجوه و صندوق امانات و در دسترس بودن خدمات مبتنی بر فن آوری، نرخ بهره پایین تسهیلات، محل مناسب شعبه، منظربانک، مدیریت خوب، محیط کلی بانک می باشد (Aregbeyen, ۲۰۱۱ - Rehman and Ahmed, ۲۰۰۸ - Siddiqi, ۲۰۱۱).

از سوی دیگر، مشتری هسته عملیات بانک ها می باشد، به طوری که پرورش و حفظ آنها برای موفقیت بانکها مهم است. با رشد چشمگیر سیستمهای اطلاعاتی در سامانه های فروش و جمع آوری داده های عملکردی مشتریان نیاز به ابزارهای مناسب جهت تحلیل داده های مشتریان و رفتارسنجی آنان جهت ارتباط موثر با مشتریان دوجندان شده است. بانک ها و موسسات مالی جهت ارتباط موثر با مشتریان خود ضمن ثبت مشخصات مشتریان، داده های تراکنش های مالی آنها مانند زمان انجام تراکنش، نوع کالا یا خدمات خریداری شده را در سامانه های اینترنتی یا کامپیوتری ثبت می کنند. ثبت این داده ها بدون پردازش آنها دانش موجود در رفتار آنها جهت ارائه خدمات من جمله تسهیلات بانکی را ندارد. داده کاوی ابزار

کاوش در داده‌های زیاد به دنبال کشف دانش از آنها است. لذا استفاده از تکنیکهای داده‌کاوی بر روی داده‌های مشتری در حل مشکل یاد شده و تحقق سودآوری و کاهش ریسک اعتباری مشتریان کمک خواهد کرد. بسیاری از مسائل داده‌کاوی را می‌توان به صورت یک مسأله کلاسه‌بندی<sup>۱</sup> بیان نمود، که در آن در نهایت یک عامل پیش بین تربیت می‌شود که می‌تواند با در دست داشتن دانش موجود برای طبقه‌بندی یک مجموعه از موارد، آن دانش را به کلاسه‌بندی سایر موارد تعمیم دهد. در واقع مسأله کلاسه‌بندی، یک مسأله یادگیری نظارت شده است.

### ۳. کلاسه‌بندی با درخت‌های تصمیم در داده‌کاوی

درخت تصمیم شیوه منحصر به فردی از ارائه یک سیستم است، که تصمیم‌گیری‌های آتی را تسهیل و سیستم را به نحو مناسبی تعریف می‌کند. با توجه به اینکه اکثر سیستم‌های مهندسی، اجرایی و محاسباتی را می‌توان در قالب یک سری داده (ویژگی یا ویژگی‌ها و خروجی منطبق با آنها) تعریف کرد، می‌توان با استفاده از یک الگوریتم (ایجاد درخت (ویژگی‌ها و خروجی‌ها را آنالیز کرد و سیستم را بر اساس این داده‌ها در قالب یک درخت تصمیم ارائه نمود. استفاده از درخت‌های تصمیم به عنوان کلاسه‌بندی‌کننده در حال حاضر کاملاً کاربردی شده و انواع مختلف پیاده‌سازی آن من جمله ID<sup>۳</sup> و C<sup>۴,۵</sup> عملاً در کاربردهای متنوع استفاده می‌شود. درخت تصمیم بر اساس آنالیز داده‌های ورودی و به منظور پیدا کردن یک ویژگی<sup>۲</sup> به عنوان مبنای تصمیم‌گیری هر گره بکار می‌رود. در هر گره ویژگی‌های مختلف داده‌ها بررسی شده و یک ویژگی که در صورت انتخاب باعث کاهش بی‌نظمی (آنترپی) می‌شود، گزینش می‌شود (Quinlan, ۱۹۹۶).

بعد از ایجاد یک درخت می‌توان از ساختار ایجاد شده برای کلاسه‌بندی داده‌های تست استفاده کرد؛ در اینصورت توسط این ساختار سریع، داده‌های تست به خوبی برای جداسازی در کلاس‌های مختلف تصمیم‌گیری می‌شوند. عموماً داده‌های تعلیم که برای ایجاد درخت استفاده می‌شوند، از داده‌های تست که برای ارزیابی درخت ایجاد شده استفاده می‌شوند متفاوت بوده و تعداد خطا در تشخیص کلاس‌های داده‌تست، معیار مناسب بودن الگوریتم می‌باشد. الگوریتم ID<sup>۳</sup> یک الگوریتم برای ساختن درخت تصمیم می‌باشد. در این الگوریتم از مفهوم بی‌نظمی<sup>۳</sup> برای دسته‌بندی داده‌ها استفاده شده است و الگوریتم در صدد آن است که میزان بی‌نظمی در گره‌های بالایی درخت حداقل باشد تا بتوان درختی با حداقل ارتفاع داشت. پس ابتدا برای تمامی ویژگی‌های داده‌های اولیه بی‌نظمی را محاسبه کرده و سپس آن ویژگی که بیشترین سودمندی را خواهد داشت به عنوان ریشه انتخاب می‌کند. این الگوریتم فقط قادر به دسته‌بندی داده‌ها با دامنه ویژگی‌های گسسته و محدود می‌باشد و در مورد داده‌های نویزی و مخدوش کارایی ندارد.

#### ۳.۱. الگوریتم C<sup>۴,۵</sup>

الگوریتم C<sup>۴,۵</sup> تکمیل شده الگوریتم ID<sup>۳</sup> می‌باشد. این الگوریتم قادر به دسته‌بندی داده‌های پیوسته و نویزی نیز می‌باشد. برای این منظور ابتدا داده‌ها را مرتب کرده سپس مقادیر سودمندی را برای تمامی حالت‌هایی که امکان جداسازی این داده‌های مرتب شده از هم وجود دارد را بدست آورده و جداساز متناظر با بزرگترین مقدار سودمندی را به عنوان جداکننده انتخاب می‌کنیم. به طور کلی کلاسه‌بندی داده‌ها یک فرآیند دو مرحله‌ای است:

<sup>۱</sup>Classification  
<sup>۲</sup>Feature  
<sup>۳</sup>Entropy

۱. در مرحله اول یک مدل ساخته می‌شود که مجموعه‌ای از کلاسهای داده‌ای یا مفاهیم را مشخص می‌کند. این مرحله را مرحله یادگیری<sup>۱</sup> گوئیم که در آن یک الگوریتم کلاسه‌بندی یک مدل را با تحلیل یک مجموعه آموزشی<sup>۲</sup> که مجموعه‌ای از تاپل‌های پایگاه است می‌سازد و برچسب کلاس‌های مربوط به این تاپل‌ها را مشخص می‌کند. یک تاپل  $X$  با یک بردار صفت  $X=(x_1, x_2, \dots, x_n)$  نمایش داده می‌شود. فرض می‌شود که هر تاپل به یک کلاس از پیش تعریف شده متعلق است و کلاس با یک صفت که به آن صفت برچسب کلاس می‌گوئیم مشخص می‌شود. مجموعه آموزشی به صورت تصادفی از پایگاه انتخاب می‌شود.

۲. در مرحله دوم، یادگیری از طریق یک تابع  $y=f(X)$  انجام می‌شود که می‌تواند برچسب کلاس هر تاپل  $X$  از پایگاه را پیش بینی کند. این تابع به صورت قواعد کلاسه‌بندی، درخت‌های تصمیم‌گیری یا فرمول‌های ریاضی است. آنچه ما در اینجا برای کلاسه‌بندی از درختان تصمیم‌استفاده می‌کنیم.

عملکرد الگوریتم فوق به شرح ذیل است:

- الگوریتم با پارامترهای مجموعه آموزشی و برچسب کلاس‌های متناظر با آنها ( $D$ )، لیستی از صفات موجود در تاپل‌ها روال انتخاب ویژگی فراخوانی می‌شود.  $D$  در واقع یک بخش<sup>۳</sup> داده‌ای است و روال انتخاب ویژگی یک روال ابتکاری<sup>۴</sup> است که بهترین صفت را برای جدا کردن تاپل‌ها براساس کلاس‌ها می‌دهد. این متد از یک معیار انتخاب صفت مانند بهره اطلاعاتی<sup>۵</sup> یا شاخص جینی<sup>۶</sup> استفاده می‌کند.

- درخت در گام اول با یک گره تنها ( $N$ ) که مجموعه آموزشی را نشان می‌دهد ایجاد می‌شود.

- اگر تاپل‌های  $D$  همه از یک کلاس باشند گره  $N$  یک برگ خواهد بود و با آن کلاس برچسب می‌خورد.

- در غیر این صورت روال انتخاب ویژگی فراخوانی می‌شود تا معیار شکاف<sup>۷</sup> را مشخص کند. معیار شکاف مشخص می‌کند که کدام صفت باید در گره  $N$  مورد آزمون قرار گیرد. معیار شکاف همچنین بیان می‌کند که چه شاخه‌هایی باید از گره  $N$  با توجه به آزمون مربوطه، خارج شوند. به عبارت دیگر معیار شکاف، صفت یا نقطه شکاف را تعیین می‌کند. نقطه شکاف،  $D$  را به یکسری بخش تبدیل می‌کند. این بخش‌ها باید تا حد ممکن خالص<sup>۸</sup> باشند به این معنی که همه تاپل‌های موجود در یک بخش باید مربوط به یک کلاس باشند.

- گره  $N$  با معیار شکاف برچسب می‌خورد. یک شاخه از گره  $N$  به هر یک از خروجی‌های معیار شکاف می‌رود. تاپل‌های  $D$  متناظر بخش بندی می‌شوند.

- الگوریتم فرایند مشابهی را به صورت بازگشتی در هر یک از بخش‌های حاصل شده دنبال می‌کند.

- بخش بندی بازگشتی در صورتی که یکی از شرایط زیر بوجود آید متوقف می‌شود:

۱. اگر تمام تاپل‌ها در بخش  $D$  متعلق به یک کلاس باشند.

۲. صفتی برای بخش بندی بیشتر وجود نداشته باشد. در این حالت گره  $N$  به یک برگ تبدیل می‌شود و برچسب کلاس آن کلاس متداول در  $D$  خواهد بود.

۳. تاپلی برای یک شاخه وجود نداشته باشد، در واقع اگر یکی از بخش‌های  $D$  مانند  $D$  تهی باشد در این موارد یک برگ با برچسب کلاس متداول در  $D$  ایجاد می‌شود (Han and Kamber, ۲۰۰۱)

<sup>۱</sup>Learning

<sup>۲</sup>training set

<sup>۳</sup>Partition

<sup>۴</sup>Heuristic

<sup>۵</sup>Information Gain

<sup>۶</sup>Gini Index

<sup>۷</sup>Splitting criterion

<sup>۸</sup>pure

معیارهای مختلفی برای انتخاب صفتی که شکاف باید بر اساس آن انجام شود وجود دارد که ما از بهره اطلاعاتی استفاده می‌کنیم. اطلاعات مورد نیاز برای کلاسه بندی یک تاپل در  $D$  برابر رابطه (۱) است:

$$Info(D) = -\sum_{i=1}^m p_i \log_2(p_i) \quad (1)$$

که در آن  $P_i$  احتمال آن است که یک تاپل دلخواه در  $D$  متعلق به کلاس  $C_i$  باشد که این احتمال به صورت  $|C_i, D|/|D|$  تخمین زده می‌شود ( $|D|$  و  $|C_i, D|$  تعداد تاپل‌ها در  $D$  و  $C_i, D$  را نشان می‌دهد). تعداد کلاسه‌های موجود  $m$  است. در واقع رابطه (۱) نشان دهنده‌ی نظم‌ی می‌باشد.

فرض می‌کنیم صفت  $A$  دارای  $v$  مقدار متمایز بصورت  $\{a_1, a_2, \dots, a_v\}$  باشد یا عبارت دیگر  $A$  یک صفت گسسته است. اگر بخواهیم  $D$  را بر حسب صفت  $A$  بشکافیم  $v$  بخش یا زیرمجموعه مانند  $\{D_1, D_2, \dots, D_v\}$  حاصل می‌شود که در آن  $D_j$  شامل تاپل‌هایی از  $D$  است که مقدار صفت  $A$  در آنها برابر  $a_j$  است. اگر فرض کنیم که  $D$  در گره‌ای چون  $N$  قرار داشته باشد آنگاه این بخش‌ها متناظر با شاخه‌هایی هستند که از  $N$  خارج می‌شوند. اطلاعات مورد نیاز برای کلاسه‌بندی یک تاپل از  $D$  بر حسب صفت  $A$  برابر با:

$$Info_A(D) = \sum_{j=1}^v |D_j|/|D| \times Info(D_j) \quad (2)$$

است. عبارت  $|D_j|/|D|$  در واقع وزن بخش  $j$  را نشان می‌دهد.

اطلاعات حاصل از انشعاب بر حسب صفت  $A$  را به صورت زیر تعریف می‌کنیم:

$$Gain(A) = Info(D) - Info_A(D) \quad (3)$$

هر چه مقدار بهره صفت  $A$  بیشتر باشد یا به عبارت دیگر هر چه  $Info_A(D)$  کمتر باشد صفت  $A$  به عنوان صفت شکاف انتخاب می‌شود.

#### ۴. سیستم پیشنهادی

سیستم پیشنهادی در شکل ۳ نشان داده شده است. اطلاعات مربوط به تراکنش‌های مشتری از پایگاه داده گرفته شده و تکنیک پیش پردازشی داده کاوی روی این داده‌ها استفاده می‌شود. بعد از چندتجزیه و تحلیل آماری، یک کلاس سلسله مراتبی از پیش تعریف شده (مقادیر هدف و یا ویژگی‌های اختیاری) در هر اطلاعات حساب اختصاص داده می‌شود. سپس روی این داده‌ها یک روش طبقه بندی مناسب انجام می‌گیرد. روش کلاسه بندی درخت تصمیم گفته شده برای استخراج قواعد برای پیش بینی رفتار مشتریان استفاده شده است.



شکل ۳: نمایی از سیستم پیشنهادی

## ۵. پیاده سازی و طراحی سیستم

### ۵.۱. شناسایی دوره و گردش داده

برای اطلاعات اولیه از پایگاه داده اطلاعات مشتریان شعبه مرکزی یکی از بانکهای خصوصی ایران استفاده کردیم. با توجه به موقعیت شعبه که در مرکز شهر واقع بوده و اندازه شعبه که بزرگترین شعبه بانک فوق است و درجه ممتاز را در درجه بندی شعب استان ها از سال ۱۳۹۱ اخذ نموده است؛ تنوع و تعدد مشتریان اطمینان از همه جانبه بودن و صحت نتایج خروجی را به ما می دهد.

با توجه به در دسترس بودن داده ها و دوره تناوب زمان خود، تصمیم به شروع انتخاب داده با استخراج یک زیر مجموعه از تراکنش های حساب ها که در وضعیت "فعال" برای مدت یک سال قبل از پرداخت تسهیلات بود، نمودیم. دلیل این انتخاب، وضعیت حساب مشتری قبل از پرداخت تسهیلات می باشد که در آن معدل موجودی شش ماهه و یک ساله حساب مشتری در زمان پرداخت تسهیلات مشخص می شود.

با توجه به اینکه درخواست هایی که از سوی سازمان ها و ادارات جهت پرداخت تسهیلات معرفی می شوند و تحت عنوان تسهیلات تکلیفی در بانک ها شناخته می شود و بانک موظف به پرداخت این تسهیلات می باشد اطلاعات این مشتریان در پایگاه داده استفاده نشد از این جمله به پرداخت تسهیلات قرض الحسنه ازدواج می توان اشاره نمود و تسهیلات پرداختی شش ماهه سال از ابتدای سال ۱۳۹۴ یعنی ۱۳۹۴/۰۱/۰۵ تا ۱۳۹۴/۰۶/۳۱ انتخاب شده و با توجه به وضعیت انتقال تسهیلات به سرفصل مطالبات سررسید گذشته که پس از دو ماه از سررسید اقساط پرداخت نشده و به طبقه معوق پس از گذشت شش ماه از اولین قسط پرداخت نشده بازه زمانی را حداکثر ۱۳۹۴/۰۸/۱۵ در نظر گرفتیم. "از تاریخ" و "تا تاریخ" زیر پرس و جوی مورد استفاده در پرس و جو اصلی برای بازیابی اطلاعات از پایگاه داده هستند.

برای انواع تسهیلات، از انواع سرفصل و معین هایی که در سیستم بانکی استفاده می شود، بهره گرفتیم هر سرفصل تسهیلاتی یک عدد نهایت سه رقمی بوده که نوع تسهیلات را نشان می دهد. معین طبقه ریز تری از سرفصل بوده و اطلاعات کامل تری از نوع تسهیلات به ما می دهد برای مثال سرفصل ۴۳ مربوط به تسهیلات مشارکتی بوده و معین ۱ بخش صنعت و معدن و معین ۵ مربوط به بخش کشاورزی می باشد که ما ترکیب سرفصل و معین را در تحقیق خود استفاده نمودیم ( همانند ۴۳۱ و ۴۳۵ در نمونه گفته شده). همچنین برای بدست آوردن ریسک اعتباری هر شخص از فرمول زیر استفاده نمودیم:

$$\text{Risk} = m * 100 / M$$

(۴)

که در آن Risk ریسک اعتباری مشتری تسهیلاتی و m میزان اقساط سررسید شده پرداخت نشده و M میزان تسهیلات می باشد. نهایتاً اطلاعات به دست آمده در یک جدول موقت برای مورد استفاده قرار گرفتن با پیوستن با تمام جداول مربوطه واکشی شده در بخش های بعدی ذخیره می شوند.

## ۵.۲. انتخاب داده ها

اطلاعات استخراج شده از پایگاه داده بانک ، یک ردیف از هر حساب برای داده های اساسی بود مانند تعداد حساب ، عنوان حساب، شغل ، توانایی پرداخت و غیره . نتیجه داده های اولیه در جدول ۱ نشان داده شده است مقادیر این متغیر های ثابت در طول زمان تغییر نمی کند. مقادیر میانگین داده های تراکنشی متغیرهای حساس به زمان روزانه تغییر پیدا می کند و حفظ تمام این مقادیر مختلف برای شش ماه در جهت پیدا کردن رفتار فصلی مرتبط به فعالیت های حساب امری ضروری است. برای همین منظور معدل موجودی سه ماهه و شش ماهه و یکساله حساب استخراج شده است. داده های تراکنشی یعنی حداقل مقدار اعتبار یک تک تراکنش از هر حساب با استفاده از پرس و جوی زیر از پایگاه داده که در آن "از تاریخ" و "تا تاریخ" در بخش قبلی توضیح داده شد واکشی شده است . برخی اطلاعات دیگر از جمله هزینه استعلام از بانک مرکزی، کارمزد ارائه خدمات ، سود متعلقه ، هزینه پست، ۱/۵در هزار هزینه کارشناسی ، حق ماموریت، حق بیمه، هزینه ارزیابی وثایق ملکی و سایر هزینه های دیگر در پرس و جوی تراکنش های تولید شده نادیده گرفته شد. خروجی در بررسی دیگر برای مقادیر صفر مورد آزمون قرار گرفت.

جدول ۱: برخی از داده های اولیه پایگاه داده مشتریان تسهیلاتی شعبه بانک

شماره قرارداد	سرفصل	تاریخ ایجاد	مبلغ تسهیلات	نوع وثیقه	ارزش وثیقه	شغل
۹۴۸۸۸۴۰۱۶۹۸	۴۷۱	۹۴۰۵۲۹	۱۰۰	۶۰	۶	فروشنده لوازم یدکی اتومبیل
۹۴۸۲۹۰۲۰۶۸۱	۴۳۱	۹۴۰۶۱۶	۲۰۰	۶۰	۱۱	شرکت خدماتی
۹۴۳۹۶۶۴۴۳۴۰	۴۳۱	۹۴۰۵۱۱	۱۰۰	۶۰	۱۲	راننده تریلر
۹۴۳۱۲۳۶۹۱۱۷	۱۹۳۸	۹۴۰۵۰۸	۶۳	۶۰	۱۳	کارگر شرکت
۹۴۸۸۶۸۶۴۲۸۸	۴۳۱	۹۴۰۵۱۸	۱۱۰	۶۰	۱۵	مهندس معدن
۹۴۸۱۹۴۷۱۸۳۱	۴۳۱	۹۴۰۵۲۸	۱۱۵	۶۰	۱۷	پیمانکار آب و فاضلاب روستایی
۹۴۳۴۵۴۹۹۱۳۷	۴۳۱	۹۴۰۵۰۸	۱۰۰	۶۰	۲۰	کارگر کارخانه
۹۴۸۸۶۳۴۵۳۷۲	۴۳۱	۹۴۰۶۰۳	۱۰۰	۶۰	۲۰	کارمند دولتی
۹۴۳۹۴۴۶۴۲۱۲	۴۳۱	۹۴۰۵۱۱	۲۰	۸۰	۲۳	منشی
۹۴۸۱۱۴۰۹۱۱۵	۴۳۱	۹۴۰۵۲۵	۱۳۰	۶۰	۲۴	راننده تریلر
۹۴۸۱۹۱۲۳۷۲۱	۱۹۳۸	۹۴۰۵۱۷	۷۵	۶۰	۲۵	فرشباغ
۹۴۸۱۵۰۵۲۹۱۹	۴۳۱	۹۴۰۵۱۳	۱۲۵	۶۰	۲۵	نمایشگاه اتومبیل
۹۴۸۲۷۴۳۱۳۹۱	۴۳۱	۹۴۰۶۳۱	۲۴	۸۰	۲۹	کارمند دولتی
۹۴۰۷۲۷۳۴۹۳۷	۴۳۱	۹۴۰۳۲۰	۲۶	۸۰	۳۱	دامدار
۹۴۸۱۹۸۳۳۳۳۸	۱۹۳۸	۹۴۰۶۰۹	۱۰۰	۶۰	۳۳	راننده تریلر
۹۴۰۵۹۲۸۲۵۷۴	۴۳۱	۹۴۰۲۲۷	۴۰	۸۰	۵۰	کارمند دولتی

## ۶. تحلیل آماری



بخش قابل توجهی از مجموعه داده‌ها فیلتر کردن این فیلدها در مراحل اولیه به طور قابل توجهی می‌تواند به کاهش زمان پردازش و بهبود مدل دقت کمک کند (Hu, 2005). تعدادی از رکوردها که یک ویژگی تک مقداری داشتند و برخی از رکوردها نیز مقادیر صفر داشتند، آماری معنی‌دار تلقی نشدند و ارزش آماری نداشتند. این فیلدها از جدول داده‌های هدف به منظور کاهش زمان محاسبات برای حصول اطمینان از فرایند مدل‌سازی سریع‌تر حذف شدند.

بسیاری از داده‌ها در جدول نهایی عددی بودند به عنوان مثال نوع وثیقه از جمله داده‌های با اهمیت یافته اصلی از داده‌های تراکنشی از مشتریان تسهیلاتی شعبه بانک بود. در تحقیق از نرمال‌سازی مقیاس دهدهی<sup>۱</sup> استفاده شد همچنین ما برای یکسان‌سازی ارقام مبالغ را در واحد میلیون ریال در نظر گرفتیم.

#### ۷. تعریف مقادیر هدف

بخشی از روش داده‌کاوی تعریف درست هدف مناسب است که با توجه به اهداف پرداخت تسهیلات برای تجزیه و تحلیل داده‌کاوی با کمک از کارشناسان حوزه اعتباری که در بخش وام مشغول هستند، مقادیر هدف از نظر داده‌های موجود تعریف شد و با این، مقدار متغیر هدف، یعنی متغیری که ویژگی‌های اختیاری را تعیین می‌کند، تعریف شد. بدین وسیله به عنوان برچسب کلاس تعریف شد. شرایط استفاده از تعریف: عالی (برابر و بالاتر از ۹۰)؛ خیلی خوب (برابر و بالاتر از ۸۰ و کمتر از ۹۰)؛ خوب (برابر و بالاتر از ۷۰ و کمتر از ۸۰)؛ مرزی (برابر و بالاتر از ۶۰ و کمتر از ۷۰)؛ و بد (کمتر از ۶۰) این است.

#### ۸. درخت تصمیم

درخت تصمیم از یک درخت برای ساخت یک مدل پیش‌بینی (تخمین) استفاده می‌کند که مشاهدات درباره یک آیتم را به نتیجه‌گیری‌هایی درباره مقدار هدف آن آیتم نگاشت می‌دهد. در واقع درخت تصمیم‌گیری یک ساختار درختی مانند فلوچارت است که مجموعه‌ای از قوانین را برای استفاده مدل پیش‌بینی ایجاد می‌کند. مزیت این روش این است که قوانین برای درک بسیار آسان است و آنها غالباً برای کشف فرآیندهای پرداخت تسهیلات اصولی مفید هستند. الگوریتم با یک مجموعه آموزش و برچسب کلاس مرتبط با خود را شروع می‌شود. مجموعه آموزش به صورت بازگشتی به زیر مجموعه‌های کوچکتر مانند درخت که ساخته شده است تقسیم می‌شود.

<sup>۱</sup>Decimalscaling normalization.

```

1 function test_targets = C4_5(train_patterns, train_targets, test_patterns, inc_node)
2
3 [N1, M] = size(train_patterns);
4 inc_node = inc_node*M/100;
5 Nu = 10;
6
7 discrete_dim = zeros(1,N1);
8 for i = 1:N1,
9     Ub = unique(train_patterns(i,:));
10    Nb = length(Ub);
11    if (Nb <= Nu),
12        discrete_dim(i) = Nb;
13        dist = abs(ones(Nb,1)*test_patterns(i,:) - Ub'*ones(1, size(test_patterns,2)));
14        [m, in] = min(dist);
15        test_patterns(i,:) = Ub(in);
16    end
17 end
18 %Build the tree recursively
19 disp('Building tree')
20 tree = make_tree(train_patterns, train_targets, inc_node, discrete_dim, max(discrete_dim), 0);
21 view;
22 %Classify test samples
23 disp('Classify test samples using the tree')
24 test_targets = use_tree(test_patterns, 1:size(test_patterns,2), tree, discrete_dim, unique(train_targets));
25 %END
26 function targets = use_tree(patterns, indices, tree, discrete_dim, Uc)
27 %Classify recursively using a tree
28 targets = zeros(1, size(patterns,2));
29 if (tree.dim == 0)

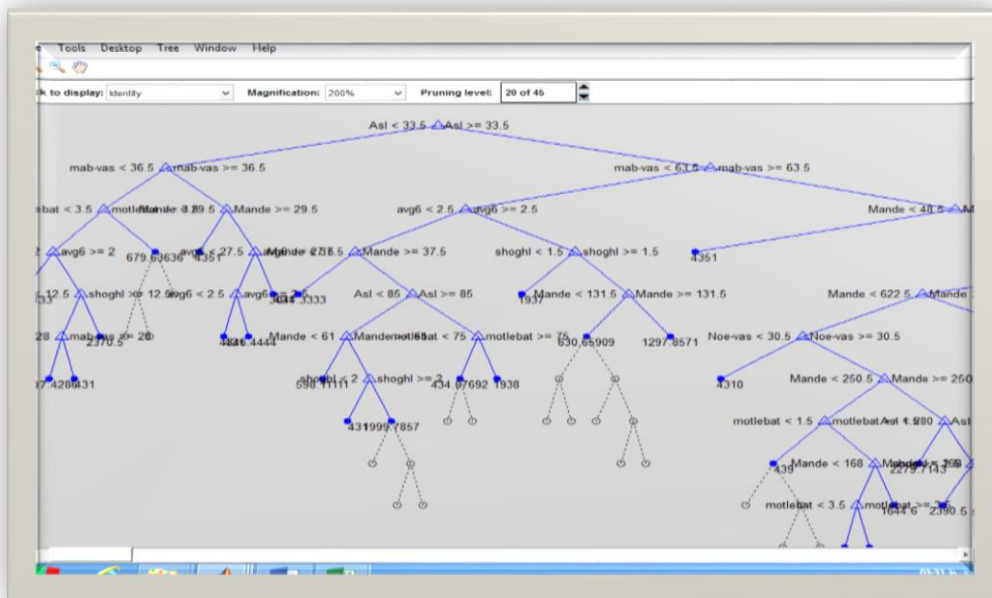
```

شکل ۴: پیاده سازی الگوریتم C4.5 در نرم افزار Matlab

### ۹. نتیجه و بحث

داده های آموزشی از پایگاه داده تراکنشی بعد از ۱۳۹۳/۰۱/۰۵ و قبل از ۱۳۹۴/۰۸/۱۵، جمع آوری شد. این مدل با برخی از نمونه داده ها برای قبل از ۱۳۹۳/۰۱/۰۵ و بعد از آن ۱۳۹۴/۰۸/۱۵ از پایگاه داده مورد آزمایش قرار گرفت. این مدل همچنین با داده های جمع آوری شده از نزدیک مورد آزمایش قرار گرفت. یافته های اولیه برخی نتایج جالب را نشان می دهد.

به منظور آزمایش اثر از مدل، مدل با نرم افزار MATLAB تست شده و متوجه شدیم که خروجی آنها تقریباً با اطلاعات



اولیه یکسان هستند. خروجی نتیجه MATLAB در شکل ۵ نشان داده شده است. از خروجی می توان گفت که مدل به اندازه کافی دقیق برای پیش بینی وضعیت مشتریان بنگاه اقتصادی در دریافت تسهیلات از یک شعبه بانک تجاری است.

شکل ۵: خروجی از نرم افزار MATLAB برای داده های مشتریان تسهیلاتی

مطالعات اولیه تعدادی از روابط بین متغیرها را پرده برداری کرد که بحث و تجزیه و تحلیل بیشتری را تصدیق می کرد. صداقت و امانتداری متقاضی دریافت تسهیلات و سابقه دریافت تسهیلات، عملکرد مثبت ۵ ساله، وضعیت رقابتی بازار و شرایط و مشخصات آن از جمله عواملی بود که می توانست بررسی ما را در پیش بینی رفتار مشتریان خرد بهبود بخشد.

#### ۱۰. نتیجه گیری

ما عوامل و بخش های چشم انداز اعطای تسهیلات به بنگاه های اقتصادی را در شعبه بانک با اجرای رویه ای از وظایف داده کاوی مورد بحث قرار دادیم. ما تعداد مطلوب خوشه ها را (عالی، روش های طبقه بندی خیلی خوب، خوبی، مرزی و بد) از پیش تعریف کردیم و از متد طبقه بندی (درخت تصمیم) برای پیش بینی رفتار مشتریان بخش بنگاه اقتصادی برای دریافت تسهیلات بانکی بر روی داده های رفتاری تراکنشی مشتریان موجود استفاده نمودیم. یافته های اولیه برخی از نتایج جالبی را نشان می دهد. آزمون ما را راهنمایی و متوجه کرد که مدل پیش بینی داده کاوی برای چشم انداز بخش های بنگاه اقتصادی برای پرداخت وام بسیار دقیق است.

تمرکز اصلی این مقاله بر صفات کتبی از مشتری بنگاه اقتصادی بود اما مشتری ویژگی های بسیاری از رفتار غیرکتبی و شفاهی دارد. پس موضوع از برخی از ویژگی های مهم نانوخته همچون صداقت، شخصیت، شرایط بازار مستثنی است. این مدل بر اساس روش های داده کاوی طبقه بندی درخت تصمیم گیری به عنوان یک گام اولیه توسعه یافته است. ما یک طرح برای گسترش مدل برای روشهای داده کاوی دیگر در آینده داریم.

#### ۱۱. مراجع

- [۱] Xiaohua Hu, (۲۰۰۵) A Data Mining Approach for Retailing Bank Customer Attrition Analysis. Applied Intelligence. Vol. ۲۲, pp. ۴۷-۶۰.
- [۲] OmoAregbeyen, Ph.D, (۲۰۱۱) The Determinants of Bank Selection Choices by Customers: Recent and Extensive Evidence from Nigeria. International Journal of Business and Social Science. Vol.۲, No. ۲۲, pp.۲۷۶-۲۸۸.
- [۳] Hafeez Ur Rehman and Saima Ahmed, (۲۰۰۸) An Empirical Analysis of the determinants of bankselection in Pakistan; A customer view. Pakistan Economic and Social Review. Vol. ۴۶, no. ۲, pp. ۱۴۷-۱۶۰.
- [۴] Kazi Omar Siddiqi, (۲۰۱۱) Interrelations between Service Quality Attributes, Customer Satisfaction and Customer Loyalty in the Retail Banking Sector in Bangladesh. International Journal of Businessand Management. Vol. ۶, No. ۳, pp.۱۲-۳۶.
- [۵] J. R. Quinlan, (۱۹۹۶) Improved use of continuous attributes in c۴.۵. Journal of Artificial Intelligence Research, ۴:۷۷-۹۰.

[۶] Jiawei Han and Micheline Kamber, (۲۰۰۱) "Data Mining: Concepts and Techniques", Morgan Kaufmann, ۲۰۰۱.

Archive of SID