

# ارائه‌ی یک سامانه‌ی هوشمند با تلفیقی از درخت رگرسیونی و نقشه‌ی خودسازمانده بهینه‌شده برای تقسیم‌بندی بهینه‌ی مشتریان

علیرضا سروش\* (دانشجوی دکتری)

اردشیر بحرینی‌نژاد (استادیار)

محمدرضا امین‌ناصری (دانشیار)

دانشکده‌ی فنی و مهندسی بخش مهندسی صنایع، دانشگاه تربیت مدرس

تقسیم‌بندی بهینه‌ی مشتریان بر مبنای ویژگی‌های مرتبط می‌تواند به توسعه‌ی استراتژی‌های بازاریابی دقیق‌تر به منظور صرف کارا تر منابع کمک کند. اما ایجاد سامانه‌ی تقسیم‌بندی مشتریان که علاوه بر پیچیدگی کم از قابلیت تقسیم‌بندی بهینه‌ی برخوردار باشد، به دلیل حجم زیاد ویژگی‌ها کاری بسیار مشکل است. هدف این نوشتار، ارائه‌ی یک سامانه‌ی تلفیقی هوشمند مبتنی بر درخت رگرسیونی و نقشه‌ی خودسازمانده بهینه‌سازی شده است که از نظر محاسباتی کارا و دقیق باشد. نتایج نشان می‌دهد که درخت رگرسیونی ۹۳٪ از ویژگی‌ها را در حالت بهینه حذف می‌کند و لذا به کاهش قابل توجه هزینه‌ی محاسبات می‌انجامد. به علاوه، نتایج اعتبارسنجی نشان می‌دهد که این سامانه با دقت قابل توجهی خوشه‌ها را تفکیک کرده است و بدین طریق می‌توان منابع بازاریابی را برای جذب مشتریان مشابه با مشتریان بهترین خوشه‌ها صرف کرد.

واژگان کلیدی: مدیریت ارتباط با مشتری، انتخاب ویژگی، درخت رگرسیونی،

تقسیم‌بندی مشتریان، نقشه‌ی خودسازمانده.

a.soroush@modares.ac.ir  
bahreininejad@modares.ac.ir  
amin\_nas@modares.ac.ir

## ۱. مقدمه

با گسترش رقابت در عرصه‌های جهانی و محدودیت منابع، به‌کارگیری صحیح آن به‌عنوان یکی از چالش‌های عمده‌ی مدیریت مطرح شد. روابط بین کارکنان و مشتریان نیز به‌عنوان مهم‌ترین و کمیاب‌ترین منابع سازمان -- که نقش قابل توجهی در سودآوری آن دارند -- از چنان اهمیت و توجه روزافزونی برخوردار شد که می‌توان گفت مشتریان تنها مرکز سودده در شرکت‌ها هستند. بنابراین سازمان‌ها همواره سعی دارند براساس مفاهیم جدید بازاریابی نوین که به معنی دانش و هنر یافتن مشتری، و نگه‌داری و افزودن بر تعداد آن است، با خلق نیازها و خواسته‌های جدید و بدیع برای مشتریان و دوری جستن از اعمال قدرت و ضوابط خشک و به‌کارگیری مشارکت و تفاهم به مدیریت روابط با مشتریان پرداخته و آن‌ها را برای تضمین سودآوری و بقای خود در اختیار بگیرند.

در سال‌های اخیر با گسترش پایگاه‌های داده مشتریان و افزایش رقابت میان سازمان‌ها در جهان بیش از پیش به موضوع «مدیریت ارتباط با مشتری (CRM)»<sup>۱</sup> توجه کرده‌اند تا بتوانند ضمن شناسایی تقاضاهای مختلف مشتریان، به مزیتی رقابتی دست یابند.<sup>[۱]</sup> به‌طوری‌که سازمان‌ها به دنبال روش‌هایی برای انجام بازاریابی‌های

\* نویسنده مسئول

تاریخ: دریافت ۱۳۸۹/۶/۳۱، اصلاحیه ۱۳۸۹/۱۰/۲۵، پذیرش ۱۳۹۰/۳/۲۱.

برخی از ترکیبی از آن‌ها<sup>[۷]</sup> و یا پیمایش رضایت مشتری<sup>[۸]</sup> برای تقسیم‌بندی مشتریان استفاده کرده‌اند. هدف این محققین از تقسیم‌بندی مشتریان شناسایی و اولویت‌بندی مشتریان براساس سودآوری آن‌ها بوده است. برخی دیگر نیز، براساس احتیاجات مشتریان<sup>[۹]</sup> با هدف شناسایی انواع سلیقه‌های آنان به تفکیک مشتریان پرداخته‌اند. همچنین به منظور بهبود دقت پیش‌بینی مصرف روزانه‌ی ایستگاه‌های برق، براساس پروفایل روزانه‌ی مشتریان ایستگاه‌های برق را تقسیم کرده‌اند.<sup>[۱۱]</sup> در مطالعه‌ی دیگری، تقسیم‌بندی به‌گونه‌ی است که برای کاهش هزینه‌ی عملیات لجستیک نسبت به خوشه‌بندی مشتریان براساس ویژگی‌های تقاضای آن‌ها قبل از اجرای مسیریابی ناوگان در عملیات لجستیک اقدام کرده‌اند.<sup>[۱۲]</sup> بررسی‌ها نشان می‌دهد که بیشتر محققین شیوه‌ی شبکه‌ی عصبی نقشه‌ی خودسازمانده (SOM)<sup>[۲]</sup> یا حالت ساده‌تر آن، تکنیک k-means را به کار گرفته‌اند. اما هیچ‌کدام از نویسندگان اقدام به محاسبه و تعیین ویژگی‌های مؤثرتر پیش از تقسیم‌بندی نکرده‌اند و ویژگی‌ها به صورت شهودی انتخاب شده‌اند؛ این موضوع می‌تواند به عدم کارایی و دقت پایین تقسیم‌بندی منجر شود.

مسئله‌ی انتخاب ویژگی‌ها، مسئله‌ی مستقل در نظریه‌ی تشخیص الگو بوده و تاکنون حل نشده است. فرایند انتخاب ویژگی‌ها به‌عنوان مسئله‌ی از بهینه‌سازی ترکیبی کلی در یادگیری ماشین شناخته می‌شود که تعداد ویژگی‌ها را کاهش داده و داده‌های غیرمرتبط و زائد را حذف می‌کند. هدف اصلی انتخاب ویژگی، شناسایی زیرمجموعه‌ی از ویژگی‌هاست که تأثیر بیشتری بر یک متغیر پاسخ معلوم دارند.<sup>[۱۳]</sup> کشف زیرمجموعه‌ی بهینه‌ی از ویژگی‌ها معمولاً مشکل بوده و نشان داده شده که بسیاری از مسائل مرتبط NP-hard شناخته می‌شوند.<sup>[۱۴]</sup> انتخاب ویژگی‌ها موفقیت‌های بسیاری را در کاربردهای دنیای واقعی داشته است، زیرا غالباً می‌تواند ابعاد داده‌ها را برای به‌کارگیری الگوریتم‌های داده‌کاوی به‌طور چشمگیری کاهش دهد. در سال‌های اخیر، CRM یکی از زمینه‌های تحقیقاتی بوده است که در روش‌های انتخاب ویژگی به کار گرفته شده‌اند.<sup>[۱۵]</sup> پیاده‌سازی مناسب انتخاب ویژگی‌ها نه تنها اطلاعات مهم‌تر را برای تقسیم‌بندی بهینه‌ی مشتریان برمی‌گزیند، بلکه محاسبات مورد نیاز برای تحلیل داده‌های چندبعدی را کاهش داده و دقت آن را افزایش می‌دهد. به همین دلیل، پیاده‌سازی انتخاب ویژگی برای آماده‌سازی سامانه‌ی تقسیم‌بندی مشتریان بهینه دارای اهمیت است.

در سال‌های اخیر، چندین الگوریتم برای انتخاب مناسب ویژگی‌ها در حوزه‌ی CRM به کار گرفته شده و تعدادی مطالعه‌ی مقایسه‌ی نیز انجام شده است. به علاوه، از آنجا که هر ویژگی مورد استفاده به‌عنوان بخشی از یک رویه‌ی تقسیم‌بندی می‌تواند هزینه و زمان اجرای سامانه‌ی تقسیم‌بندی مشتریان را افزایش دهد، انگیزه‌ی قوی برای طراحی و پیاده‌سازی سامانه‌ی با مجموعه‌ی ویژگی‌های کم وجود دارد. هم‌زمان یک نیاز متضاد برای لحاظ مجموعه‌ی کافی از ویژگی‌ها به منظور دست‌یابی به نرخ‌های تشخیص بالا تحت شرایط مشکل وجود دارد. این موضوع منجر به توسعه‌ی تکنیک‌های متنوعی برای کشف یک زیرمجموعه‌ی بهینه از ویژگی‌ها شده است. ساده‌ترین روش جست‌وجوی بهینه یک روش جامع است، هرچند که با این روش تعداد زیرمجموعه‌های ممکن به سرعت رشد کرده و برای مجموعه ویژگی‌های با اندازه‌ی متوسط نیز غیرقابل استفاده است. روش‌های جست‌وجوی بهینه‌ی همچون الگوریتم شاخه و تحدید وجود دارند که مانع از رویکردی جامع می‌شود. دو روش ابتکاری مشهور نیز انتخاب پیشرو متوالی و انتخاب پسرو متوالی هستند. پویا شده این دو روش به نام‌های جست‌وجوی شناور پیشرو متوالی و جست‌وجوی شناور پسرو متوالی معرفی می‌شوند.

یک رویه‌ی آماری ساده، انتخاب ویژگی رو به جلو مبتنی بر امتیاز مجذور

خی‌دواست، اگرچه برای دست‌یابی به حل بهینه با محدودیت‌هایی مواجه‌اند. روش دیگر، رویکرد مبتنی بر جست‌وجوی تصادفی الگوریتم ژنتیک (GA) است که از گام‌های احتمالی یا تکنیک‌های نمونه‌گیری استفاده می‌کند. این روش وزن‌ها را به ویژگی‌ها تخصیص می‌دهد. ویژگی‌های وزن داده شده با پیشروی تا حد آستانه‌ی که توسط کاربر تعریف شده انتخاب می‌شوند.<sup>[۱۷]</sup> به منظور هدف‌گذاری مشتریان<sup>[۱۸]</sup> و نیز به منظور بهبود پیش‌بینی میزان خرید مشتری<sup>[۱۹]</sup> از یک الگوریتم ژنتیک (GA) استفاده می‌شود. به‌کارگیری GA برای تعیین ویژگی‌های مؤثر به زمان بسیار زیادی برای محاسبات نیاز خواهد داشت. محققین یک سامانه‌ی استدلال مبتنی بر مورد همراه با تکنیک کاهش دوبعدی برای رضایت مشتریان پیشنهاد کرده‌اند.<sup>[۲۰]</sup> که رفتار خرید مشتریان را برای یک محصول خاص با استفاده از مشخصه‌های آماری آن‌ها پیش‌گویی می‌کند. برای انتخاب ویژگی، ترکیبی از الگوریتم افزاینده تودرتو<sup>[۲]</sup> و روش ذوب شبیه‌سازی شده<sup>[۴]</sup> به منظور تشخیص مشتری<sup>[۲۱]</sup> و نظریه‌ی مجموعه‌ی ناهنجار<sup>[۵]</sup> به منظور پیش‌گویی رفتار خرید مشتری<sup>[۲۲]</sup> به کار برده شده‌اند. همچنین، ماشین برداری پشتیبان<sup>[۶]</sup> برای رتبه‌بندی ویژگی‌ها در کاربردهای CRM دنیای واقعی پیشنهاد شده است.<sup>[۲۳]</sup> مهم‌ترین نقطه ضعف این روش‌ها این است که خطی بوده و روابط غیرخطی بین هر ویژگی و متغیر هدف را لحاظ نمی‌کنند. برخی از محققین رویکرد درخت تصمیم را برای انتخاب ویژگی پیشنهاد، و بیان کرده‌اند که این روش بسیار دقیق بوده و سابقه‌ی عملکرد خوبی دارد.<sup>[۱۵]</sup>

مقصود این نوشتار توسعه‌ی یک سامانه‌ی هوشمند تلفیقی است که به لحاظ محاسباتی کارا و دقیق است. رویکرد انتخابی بدین صورت است: ابتدا یک درخت رگرسیونی هرس شده با سرعت محاسباتی بالا و دقتی قابل توجه برای انتخاب ویژگی‌های بهینه طراحی می‌شود. سپس یک شبکه‌ی عصبی SOM برای تقسیم‌بندی بهینه‌ی مشتریان مبتنی بر شاخص دیویس-بولدین<sup>[۷]</sup> ایجاد شده و بهینه‌ترین خوشه‌ها به ترتیب اولویت تعیین شده برای ویژگی‌های مشتریان در هر خوشه، به منظور تدوین راهکارهای بازاریابی توصیف می‌شوند. بدین ترتیب، سازمان می‌تواند بر خوشه‌هایی که احتمال تأمین هدف مورد نظر سازمان در آن‌ها بیشتر است متمرکز شود تا هزینه‌ی بازاریابی کم‌تری متوجه آن شود. از طریق به‌کارگیری چنین رویکردی، دست‌یابی به هر دو نیاز -- پیچیدگی کم و عملکرد تقسیم‌بندی بهینه -- ممکن می‌شود. پس از تعیین تعداد خوشه‌های بهینه، نتایج آن با داده‌هایی جدا از داده‌های آموزشی برای طراحی سامانه، اعتبارسنجی می‌شود. در ادامه سامانه‌ی هوشمند تلفیقی طراحی شده روی داده‌های یک شرکت بیمه پیاده‌سازی شده و قابلیت آن بررسی می‌شود.

در ادامه‌ی این تحقیق به توصیف موردکاوی تحقیق می‌پردازیم و سپس از طریق درخت رگرسیونی بهینه‌شده به انتخاب ویژگی خواهیم پرداخت. پس از آن، در بخش چهارم، تقسیم‌بندی بهینه‌ی مشتریان با توسعه‌ی SOM صورت گرفته و نتایج مورد ارزیابی قرار می‌گیرند.

## ۲. توصیف داده‌ها: موردکاوی یک شرکت بیمه

یکی از پرکاربردترین زمینه‌هایی که می‌توان برای شناسایی مشتریان از آن استفاده کرد، صنعت بیمه است. به‌ویژه آن که، شناسایی ویژگی‌های مشتریان یکی از محصولات شرکت بیمه از طریق رفتاری که در قبال سایر محصولات بیمه‌ی شرکت از خود نشان داده‌اند، می‌تواند جالب توجه باشد. از این‌رو، در این تحقیق مجموعه‌ی داده‌ها که مربوط به کسب‌وکاری در دنیای واقعی است، از طریق یک شرکت بیمه تهیه شده است. این شرکت بیمه قصد دارد مشتریان بالقوه برای یک محصول معین را

قیاسی (اسمی)، قابلیت پشتیبانی پیشگویی‌کننده‌های بسیار زیاد (حداکثر ۸۰۰۰ عدد) و انتخاب یک به یک متغیرهای مرتبط در طول فرایند. از این رو، درخت رگرسیونی می‌تواند رویکردی مناسب برای شناسایی ویژگی‌های بهینه برای مسئله‌ی تقسیم‌بندی مشتریان باشد.<sup>[۲۶،۲۵]</sup>

### ۱.۳. درخت رگرسیونی هرس شده

درخت‌های تصمیم یکی از ساده‌ترین و موفق‌ترین الگوریتم‌های یادگیری در داده‌کاوی و یادگیری ماشین هستند. شهرت این درخت‌های تصمیم از آن روست که به‌سادگی قابل تفسیر و از نظر محاسباتی کم‌هزینه‌اند. در صنعت و محیط کسب‌وکار از شیوه‌هایی برای ارزیابی اعتبار، کشف کلاه‌برداری و مدیریت ارتباط با مشتری استفاده شده‌اند.<sup>[۲۵]</sup> معمولاً هدف پیدا کردن درخت تصمیم بهینه با کمینه‌سازی خطای تعمیم است. ویژگی‌های ورودی و مقادیر خروجی می‌توانند گسسته یا پیوسته باشند. درخت تصمیم دو نام دیگر نیز دارد: ۱. درخت تصمیم با محدوددهی از برچسب‌های دسته‌بندی گسسته (قیاسی یا طبقه‌بندی)، ۲. درخت طبقه‌بندی. این در حالی است که درخت تصمیم با محدوددهی از مقادیر خروجی پیوسته (عددی) را «درخت رگرسیونی» می‌نامند. به‌طور کلی، این درخت‌های تصمیم را درخت طبقه‌بندی و رگرسیونی (CART)<sup>۸</sup> می‌نامند.<sup>[۲۷]</sup>

میزان اهمیت یک ویژگی براساس مجموع بهبودها در کلیدی‌گره‌ها تعیین می‌شود که ویژگی نقش دو نیم‌کننده (وزن‌دهی شده توسط بخشی از داده‌های آموزشی در هر شکاف‌گره) را دارد. جانشین‌ها نیز در محاسبات اهمیت لحاظ می‌شوند، بدین معنا که حتی به متغیری که هرگز یک‌گره را دو نیم نمی‌کند، ممکن است امتیاز اهمیت بزرگی اختصاص داده شود. این موضوع، رتبه‌بندی اهمیت متغیر برای مشخص‌سازی پوشانه‌های متغیر و همبستگی غیرخطی میان ویژگی‌ها را ممکن می‌سازد. امتیازهای اهمیت را می‌توان به‌صورت ارادی به دونیم‌کننده‌ها محدود کرد؛ مقایسه‌ی دونیم‌کننده‌ها و رتبه‌بندی‌های اهمیت کل تشخیصی مفید است.<sup>[۲۵]</sup>

از آنجا که درخت‌های تصمیم با پیچیدگی کم‌تر جامع‌ترند، معمولاً تصمیم‌گیرندگان آن‌ها را ترجیح می‌دهند. به‌علاوه، پیچیدگی درخت به‌دلیل یادگیری جزئیات خاص تأثیر شدیدی بر عملکرد دقت آن دارد و منجر به عملکرد تعمیمی ضعیفی خواهد شد. تکنیک درخت رگرسیونی دونیم‌کننده‌هایی را جست‌وجو می‌کند که مربعات خطای (انحراف حداقل مربعات) پیشگویی را کمینه می‌کند. پیشگویی در هر گره پایانی بر مبنای میانگین وزنی برای گره تعیین می‌شود. متداول‌ترین رویکرد مورد استفاده، ابتدا رشد یک درخت تا یک اندازه‌ی بزرگ و سپس هرس کردن گره‌ها براساس معیار هرس است. معیار توقف مورد استفاده و روش هرس به کار برده شده پیچیدگی را کنترل می‌کند.<sup>[۲۸]</sup> استفاده از معیار ارزیابی عملکرد به‌طوری که درخور مسئله‌ی تحت بررسی باشد، بسیار مهم است. هزینه‌ی میانگین مربعات خطا باید مورد محاسبه قرار گیرد، زیرا آن‌ها در انتخاب روش تأثیرگذارند. فرض برابری هزینه‌های میانگین مربعات خطا در موارد بسیار کمی مناسب است. معمولاً — حتی اگر این هزینه‌ها دقیقاً معلوم نباشند — تا اندازه‌ی می‌توان در خصوص هزینه‌ها صحبت کرد.

### ۲.۳. الگوریتم هرس<sup>۹</sup>

کارهای اولیه در حوزه‌ی درخت‌های تصمیم امکان هرس را نمی‌دادند. در آن شرایط، درخت‌ها آن‌قدر رشد می‌کردند تا با شرط توقف برخورد کنند، و نهایتاً درخت منتج به‌عنوان درخت نهایی در نظر گرفته می‌شد. «هرس» فرایند کاهش یک درخت از طریق تبدیل برخی گره‌های شاخه‌ی به گره‌های پایانی و حذف گره‌های پایانی

شناسایی کند. داده‌ها مربوط به ۹۸۲۲ مشتری این شرکت بیمه است که تعدادی از آن‌ها اقدام به خرید بیمه‌نامه‌ی خودرویی به‌نام کاروان کرده‌اند. این داده‌ها امکان ارزیابی قابلیت سامانه‌ی هوشمند در تقسیم‌بندی بهینه‌ی مشتریان را فراهم می‌سازد. در این مقاله، دو مجموعه داده‌ی جداگانه به کار برده می‌شود: یک مجموعه‌ی آموزشی با ۵۸۲۲ مشتری و یک مجموعه‌ی ارزیابی با ۴۰۰۰ مشتری. هر رکورد متشکل از ۸۶ ویژگی شامل داده‌های آمارگیری اجتماعی (۴۳ ویژگی) و مالکیت محصول (۴۲ ویژگی) می‌شود. ویژگی آخر یعنی «تعداد بیمه‌نامه‌ی کاروان»، متغیر هدف است. از داده‌های آموزشی برای آماده‌سازی سامانه‌ی هوشمند استفاده می‌شود و نتیجه براساس مجموعه‌ی ارزیابی اعتبارسنجی می‌شود. از ۵۸۲۲ مشتری احتمالی در مجموعه داده‌ی آموزشی، ۳۴۸ نفر بیمه‌نامه‌ی کاروان را خریده‌اند که از آن نرخ هدف  $5.97\% = 348/5822$  برای مشتریان مورد نظر حاصل می‌شود. همچنین، از ۴۰۰۰ مشتری احتمالی در مجموعه داده‌ی ارزیابی، ۲۳۸ نفر بیمه‌نامه‌ی کاروان را خریده‌اند که نرخ هدف  $5.95\% = 238/4000$  حاصل می‌شود. چنان که مشاهده می‌شود نرخ هدف تقریباً در هر دو مجموعه برابر است. سامانه‌ی هوشمند برای تقسیم‌بندی مشتریان و شناسایی ویژگی‌های ۲۰ درصد اول مشتریانی طراحی می‌شود که انتظار می‌رود در مجموعه داده ارزیابی محتمل‌ترین افراد برای خرید بیمه‌نامه‌ی کاروان باشند. ذکر این نکته ضروری است که تنها اطلاعات در مجموعه داده‌ی آموزشی در توسعه‌ی سامانه‌ی هوشمند استفاده می‌شود و مجموعه داده‌ی ارزیابی منحصراً برای اعتبارسنجی استفاده می‌شود.

افزون بر این، در طراحی چنین سامانه‌ی سه موضوع ضرورت می‌یابد: ویژگی‌های منتخب، الگوریتم خوشه‌بندی، ارزیابی خروجی خوشه‌بندی. این سه موضوع کارایی، اثربخشی و کیفیت تقسیم‌بندی را تعیین می‌کنند.

### ۳. انتخاب ویژگی مبتنی بر یک درخت رگرسیونی

چنان که اشاره شد، یکی از مشکلات عمده در خصوص مسائل تشخیص الگو، ابعاد زیاد و انتخاب متغیرهای غیر سودآور است. بدین معنا که معمولاً تعداد ویژگی‌های در اختیار طراح سامانه‌ی طبقه‌بندی یا خوشه‌بندی بسیار زیاد است و احتمال انتخاب ویژگی‌های غیرمرتبط بسیار زیاد است. انتخاب ویژگی به فرایند انتخاب مرتبط‌ترین ویژگی‌ها گفته می‌شود که بیشترین ارتباط را با متغیر(های) خروجی داشته باشند.<sup>[۲۴]</sup> این فرایند باید بسیار دقیق انجام شود؛ اگر ویژگی‌های مهم به درستی انتخاب نشوند، ویژگی‌های غیرمرتبط می‌تواند بر پیچیدگی مسئله بیفزاید و دقت مدل را کم کند. هدف شناسایی کم‌ترین زیرمجموعه از متغیرهاست که بالاترین دقت را ارائه می‌کند. این موضوع ممکن است ساده به نظر آید، که البته چنین نیست. زیرا اولاً، برای یک پایگاه داده تنها با ۲۰ ویژگی، بیش از ۱ میلیون زیرمجموعه‌ی ممکن وجود دارد. ثانیاً برای ارزیابی هر زیرمجموعه، دسترسی به یک مدل و ارزیابی آن از طریق اندازه‌گیری خطا ضروری خواهد بود.

از آنجا که اندازه‌ی فایده‌ی داده‌های مشتری در حال افزایش است، گرایش به مسئله‌ی انتخاب متغیر به‌شدت رشد داشته است. تنها حفظ مرتبط‌ترین ویژگی‌ها، اجرای کاهش ابعاد با کم‌ترین هزینه‌ی اطلاعات داده‌ی برای کاهش پیچیدگی محاسبات و افزایش دقت تقسیم‌بندی را ممکن می‌سازد. امتیازات استفاده از درخت رگرسیونی به‌عنوان یک پیش‌پردازشگر برای خوشه‌بندی مشتریان با استفاده از SOM عبارت است از: غیرخطی بودن (با توجه به غیرخطی بودن رفتار مشتری)، زمان آموزش بسیار سریع، عدم نیاز به تبدیل یا آماده‌سازی داده‌ها، به‌کارگیری خودکار پیشگویی‌کننده‌های

زیرشاخه‌ی اصلی است. سازوکار هرس اکیداً بر داده‌های آموزشی مبتنی می‌شود و با یک سنججه هزینه‌ی پیچیدگی آغاز می‌شود.

الگوریتم هرس به‌طور کلی برای درخت‌هایی به کار گرفته می‌شود که لزوماً طبقه‌بندی نشده‌اند بلکه درخت‌های رگرسیونی‌اند. فرض کنید که  $R(t)$  اعداد حقیقی همبسته با هر گره  $t$  از یک درخت معلوم  $T$  باشد. مقدار  $R(t)$  مطابق رابطه‌ی ۱ تعیین شود:

$$R(t) = e(t)p(t) \quad (1)$$

که در آن  $e(t)$  میانگین مربعات خطا معلوم می‌کند که یک مورد در گره می‌افتد. اگر ما فرض کنیم که  $N(t)$  تعداد نمونه‌های  $L = \{(x_i, y_i), i = 1, \dots, n\}$  را نشان دهد، آنگاه  $p(t)$  را می‌توان چنین تعریف کرد:

$$p(t) = \frac{N(t)}{n} \quad (2)$$

از این رو، اگر  $t$  به‌عنوان یک گره پایانی در نظر گرفته شود،  $R(t)$  سهم آن گره به خطای کل است. فرض کنید که  $T_t$  زیردرختی با ریشه‌ی  $t$  باشد. اگر  $R_\alpha(T_t) < R_\alpha(t)$ ، آنگاه سهم به هزینه‌ی پیچیدگی زیردرخت کم‌تر از سهم گره  $t$  است. این موضوع برای  $\alpha$  کوچک اتفاق می‌افتد. همچنان که  $\alpha$  افزایش می‌یابد، تساوی حاصل می‌شود، زمانی که:

$$\alpha = \frac{R(t) - R(T_t)}{N_d(t) - 1} \quad (3)$$

که در آن  $N_d(t)$  تعداد گره‌های پایانی در  $T_t$  است، یعنی  $N_d(t) = |\tilde{T}_t|$  و حذف درخت در  $t$  برتری می‌یابد. بنابراین:

$$g(t) = \frac{R(t) - R(T_t)}{N_d(t) - 1} \quad (4)$$

به‌عنوان سنججه‌ی از شدت پیوند گره  $t$  تعریف می‌شود. اولین مرحله‌ی الگوریتم، گره با کوچک‌ترین مقدار  $g(t)$  را جست‌وجو می‌کند. آن گره تبدیل به یک گره پایانی می‌شود و مقدار  $g(t)$  برای کلیه‌ی اجزای محاسبه می‌شود. این فرایند تکرار می‌شود و ادامه می‌یابد تا به گره ریشه برسیم. بدین ترتیب، الگوریتم هرس، یک رشته از درخت‌ها را تولید می‌کند.<sup>[۱۷]</sup>

### ۳.۳. بررسی اعتبار<sup>۱۰</sup>

بررسی اعتبار  $(CV)$ ، روشی برای برآورد میزان خطاست که دارای ایده‌ی ساده است. مجموعه داده‌ها به اندازه‌ی نمونه به دو قسمت تفکیک می‌شوند. پارامترهای مدل با استفاده از یک مجموعه (با کمیته‌سازی چند معیار بهینه‌سازی) برآورد شده و معیار خوبی برازش<sup>۱۱</sup> بر روی مجموعه‌ی دوم ارزیابی می‌شود. نسخه‌ی معمول بررسی اعتبار یک روش صرف‌نظر از یکی<sup>۱۲</sup> است که در آن مجموعه‌ی دوم متشکل از تنها یک نمونه است. آنگاه برآورد بررسی اعتبار معیار خوبی برازش  $(CV)$  متوسط کلیه‌ی مجموعه‌های آموزشی ممکن به اندازه‌ی  $n - 1$  است. خطای بررسی اعتبار  $(CV)$  به‌عنوان وسیله‌ی برای تعیین یک مدل مناسب، برای هر عضو از خانواده مدل‌های کاندید  $\{M_k, k = 1, \dots, K\}$  و مدل  $M_{\hat{k}}$  انتخاب شده محاسبه می‌شود، که:

$$\hat{k} = \arg \min CV(k) \quad (5)$$

بررسی اعتبار هنگام انتخاب یک مدل صحیح گرایش به برازش بیش از حد دارد، به‌طوری که، برای مجموعه داده یک مدل بسیار پیچیده انتخاب می‌کند. شواهدی

وجود دارد که بررسی اعتبار چندتایی، زمانی که  $d > 1$  نمونه از مجموعه آموزشی حذف می‌شود، انتخاب مدل بهتر از بررسی اعتبار صرف‌نظر از یکی انجام می‌شود. برای  $n$  بزرگ، میزان محاسبات زیاد بوده و به طراحی  $n$  طبقه‌بندی‌کننده نیاز دارد. با این که با وجود هزینه‌ی یک افزایش در واریانس برآوردکننده، به‌صورت تخمینی ناریب است.<sup>[۱۹]</sup>

### ۴.۳. مشخصه‌های درخت رگرسیونی هرس شده

در قسمت‌های قبل نحوه‌ی عملکرد درخت رگرسیونی هرس شده به‌عنوان ابزاری برای انتخاب ویژگی‌ها توصیف شد. در این بخش نیز به بیان مشخصه‌های آن برای انتخاب ویژگی‌ها در موردکاوای می‌پردازیم. مشخصات درخت رگرسیونی هرس شده‌ی مورد استفاده عبارت است از:

- نوع درخت رگرسیونی است، زیرا نوع متغیر هدف عددی است (۰ یا ۱).
- از بررسی اعتبار  $5^0$  تکه‌ی به‌منظور برآورد خطای واقعی برای درخت‌ها در اندازه‌های مختلف استفاده می‌شود؛ بدین معنا که تابع نمونه به  $5^0$  زیرنمونه -- که به‌صورت تصادفی انتخاب شده و تقریباً دارای اندازه‌ی برابرند -- تفکیک می‌شود. برای هر یک از زیرنمونه‌ها، یک درخت به داده‌های باقی‌مانده برازش می‌یابد و از آن برای پیش‌بینی زیرنمونه استفاده می‌شود. سپس اطلاعات تمامی زیرنمونه‌ها برای محاسبه‌ی هزینه‌ی کل نمونه با یکدیگر ترکیب می‌شود. همچنین برداری شامل خطای معیار هر مقدار هزینه، برداری شامل تعداد گره‌های پایانی برای هر زیردرخت، و اسکالری شامل بهترین سطح برآوردشده‌ی هرس محاسبه می‌شود. بهترین سطح، کوچک‌ترین درختی را که در محدوده‌ی خطای معیار از زیردرخت کم‌ترین هزینه است، تولید می‌کند. متغیرهای منتخب آن‌هایی خواهند بود که بهترین عملکرد را دارند؛ یعنی درخت دارای کم‌ترین هزینه انتخاب می‌شود.
- برای هرس کردن درخت، معیار میانگین مربعات خطا به کار برده می‌شود.
- هزینه‌ی درخت، مجموع کلیه‌ی گره‌های پایانی شامل احتمال برآوردشده‌ی هر گره در هزینه‌ی گره است. از آنجا که درخت یک درخت رگرسیونی است، هزینه‌ی یک گره میانگین مربعات خطا در کلیه‌ی مشاهدات آن گره است.
- خطا برای هر گره، واریانس مشاهدات تخصیص داده شده به آن گره است.
- احتمال یک گره از طریق نسبتی از مشاهدات داده‌های اصلی که شرایط گره را برآورده می‌کنند، محاسبه می‌شود.
- اندازه یک گره به‌صورت تعدادی از مشاهدات داده‌های مورد استفاده به‌منظور ایجاد درختی که شرایط را برای گره برآورده کند، تعریف می‌شود.

### ۵.۳. نتایج انتخاب ویژگی‌ها

براساس مشخصه‌های یادشده، درخت رگرسیونی برای ۵۸۲۲ مشتری ترسیم می‌شود که شامل ۸۵ ویژگی است. با توجه به تعداد ویژگی‌ها و مشتریان می‌توان دریافت که اندازه و پیچیدگی درخت بسیار زیاد است؛ بنابراین مجموعه‌قواعد تولیدشده نمی‌توانند خیلی شهودی باشند. برای دست‌یابی به توصیف فشرده‌تر داده‌ها، برخی از شاخه‌های درخت را هرس می‌کنیم به‌طوری که درختی کوچک‌تر را خواهیم دید. بنابراین، درختی انتخاب می‌شود که کم‌ترین هزینه با تعداد ویژگی کم‌تر و تعداد گره‌های پایانی کم‌تر را دارد. در شکل ۱ نقطه‌ی بهینه برای هرس درخت با کم‌ترین هزینه نمایش داده شده است.

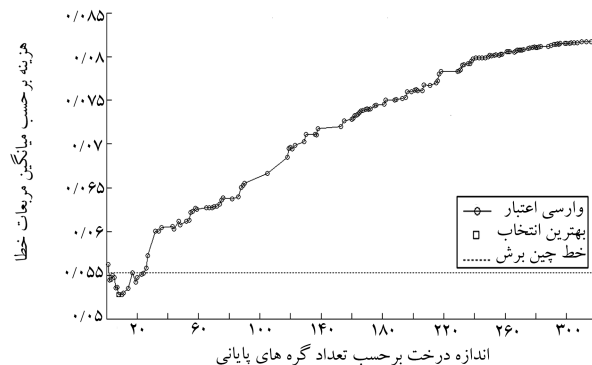
می‌رسند، زیرمجموعه‌ی منتخب بهینه می‌شوند. جدول ۱، قواعد درخت، اندازه درخت و مقادیر  $R(t)$ ،  $e(t)$ ،  $p(t)$  در هر گره را بعد از مرحله‌ی هرس بهینه نشان می‌دهد.

چنان که در جدول ۱ مشاهده می‌شود، پس از مرحله‌ی هرس بهینه ۱۵ گره باقی می‌ماند. شماره‌ی گره‌های پایانی ۲، ۶، ۱۰، ۱۱، ۱۲، ۱۳، ۱۴ و ۱۵ هستند و هر یک نشان می‌دهند که چه تعداد مشتری در ستون اندازه براساس چه قاعده‌یی و چه احتمالی در نظر گرفته می‌شوند. خطا برای هر گره در ستون پنجم نشان‌گر واریانس مشاهدات تخصیص داده شده به آن گره است، و ریسک برای هر گره در ستون ششم نشان‌گر خطای گره وزن‌دهی شده توسط احتمال گره است.

نهایتاً، شش ویژگی به نام‌های بیمه‌نامه‌ی خودرو، نوع مشتری، بیمه‌نامه‌ی آتش‌سوزی، بیمه‌نامه‌ی قاق، کارگر بی‌تجربه و بیمه‌ی شخص ثالث به عنوان تأثیرگذارترین ویژگی‌ها بر هدف، یعنی خرید بیمه‌نامه‌ی کاروان، انتخاب می‌شوند. این رویه‌ی انتخاب ویژگی منجر به کاهش قابل توجه، معادل ۹۳٪ در داده‌های ورودی می‌شود. در مرحله‌ی بعد پیاده‌سازی سامانه‌ی هوشمند تقسیم‌بندی مشتریان طراحی شده، خریداران بیمه‌نامه‌ی کاروان براساس ویژگی‌های منتخب بهینه خوشه‌بندی می‌شوند.

#### ۴. تقسیم‌بندی مشتریان

در این مرحله پس از تعیین ویژگی‌های مؤثر بر رفتار خرید مشتری به تقسیم‌بندی مشتریان جهت شناسایی ویژگی‌های مشترک مشتریان هر خوشه برای صرف بهینه هزینه‌های بازاریابی خواهیم پرداخت. بنابراین، نیاز به ابزاری است که یاد بگیرد چگونه به ویژگی‌های مختلف وزن‌دهی کرده و همزمان یک نقشه تعمیم یافته که بیش از حد منطبق نباشد، تولید کند. بدین منظور از آنجا که دسته‌بندی اولیه‌ی مشخصی وجود ندارد (مستلایه تقسیم‌بندی از نوع خوشه‌بندی است نه طبقه‌بندی) و با توجه به روابط غیرخطی که در میان ویژگی‌ها دیده می‌شود، از شبکه‌ی عصبی کوهون



شکل ۱. نقطه‌ی بهینه برای هرس درخت رگرسیونی با کم‌ترین هزینه.

چنان که در تصویر مشاهده می‌شود، خط ترسیم شده نشان‌گر هزینه‌ی برآوردی برای هر اندازه درخت است؛ خط چین نشان‌گر یک خطای استاندارد بالاتر از کم‌ترین خطاست؛ و مربع کوچک بر روی خط نشان‌گر کم‌ترین هزینه‌ی درخت زیر خط چین است. درخت کامل شامل ۳۱۶ گره پایانی است که هزینه بر مبنای معیار میانگین مربعات خطا برای زیردرخت‌های مختلف محاسبه شده است. به طور کلی، درخت ۶۳۱ گره دارد که هزینه‌ی درخت به تدریج و همزمان با رشد اندازه درخت و رسیدن به مقدار ۵۲۸٪ با ۸ گره پایانی، کاهش می‌یابد. هزینه‌ی درخت در ادامه، در روندی صعودی شروع به نوسان می‌کند به طوری که، با افزایش تعداد گره‌های پایانی پس از رسیدن به نقطه‌ی بهینه، تعداد ویژگی‌های بیشتری برگزیده می‌شوند اما منجر به هزینه‌ی بالاتری می‌شود. بنابراین، بهترین انتخاب، بهترین سطح برآوردی هرس یا نقطه‌ی بهینه را نشان می‌دهد که شامل هشت گره پایانی می‌شود. در این حالت، می‌توان متوجه شد که پیچیدگی‌های درخت از طریق کاستن تعداد گره‌های پایانی کاهش می‌یابد. زمانی که درخت رگرسیونی هرس شده در مجموعه‌ی آموزشی ترسیم می‌شود، ویژگی‌هایی که در ترسیم آن به نظر می‌رسند انتخاب می‌شوند. به طوری که، ویژگی‌هایی که در مسیرهای منتهی به هر گره پایانی در درخت هرس شده به نظر

جدول ۱. قواعد درخت بعد از مرحله‌ی هرس بهینه.

اندازه	احتمال $p(t)$	خطا $e(t)$	ریسک $R(t)$	قاعده	گره $(t)$
۵۸۲۲	۱	۰٫۰۵۶۲	۰٫۰۵۶۲	اگر بیمه‌نامه‌ی خودرو کوچک‌تر از ۵/۵ است برو به گره ۲ در غیر این صورت گره ۳.	۱
۳۴۵۹	۰٫۵۹۴۱	۰٫۰۲۴۲	۰٫۰۱۴۴	تطابق = ۰٫۲۴۹	۲
۲۳۶۳	۰٫۴۰۵۸	۰٫۰۹۸۶	۰٫۰۴۰۰	اگر نوع مشتری کوچک‌تر از ۲/۵ است، برو به گره ۴ در غیر این صورت گره ۵.	۳
۴۶۴	۰٫۷۹۶	۰٫۱۶۱۵	۰٫۰۱۲۹	اگر بیمه‌نامه‌ی آتش‌سوزی کوچک‌تر از ۳/۵ است، برو به گره ۶ در غیر این صورت گره ۷.	۴
۱۸۹۹	۰٫۳۲۶۱	۰٫۰۸۰۶	۰٫۰۲۶۳	اگر بیمه‌نامه‌ی قاق کوچک‌تر از ۰/۵ است، برو به گره ۸ در غیر این صورت گره ۹.	۵
۲۴۱	۰٫۰۴۱۳	۰٫۱۱۲۱	۰٫۰۰۴۶	تطابق = ۰٫۱۲۸۶	۶
۲۲۳	۰٫۰۳۸۳	۰٫۲۰۲۷	۰٫۰۰۷۸	اگر کارگر بی‌تجربه کوچک‌تر از ۳/۵ است، برو به گره ۱۰ در غیر این صورت گره ۱۱.	۷
۱۸۸۲	۰٫۳۲۳۲	۰٫۰۷۷۳	۰٫۰۲۵۰	اگر بیمه‌نامه‌ی آتش‌سوزی کوچک‌تر از ۲/۵ است، برو به گره ۱۲ در غیر این صورت گره ۱۳.	۸
۱۷	۰٫۰۰۲۹	۰٫۲۴۹۱	۰٫۰۰۰۷	اگر بیمه‌ی شخص ثالث مشخصی کوچک‌تر از ۱ است، برو به گره ۱۴ در غیر این صورت گره ۱۵.	۹
۲۱۴	۰٫۰۳۶۷	۰٫۱۹۳۲	۰٫۰۰۷۱	تطابق = ۰٫۲۶۱۷	۱۰
۹	۰٫۰۰۱۵	۰٫۱۷۲۸	۰٫۰۰۰۳	تطابق = ۰٫۷۷۷۸	۱۱
۹۶۱	۰٫۱۶۵۰	۰٫۰۴۸۴	۰٫۰۰۸۰	تطابق = ۰٫۰۵۱۰	۱۲
۹۲۱	۰٫۱۵۸۱	۰٫۱۰۵۲	۰٫۰۱۶۶	تطابق = ۰٫۱۱۹۴	۱۳
۱۱	۰٫۰۰۱۸	۰٫۱۴۸۸	۰٫۰۰۰۳	تطابق = ۰٫۸۱۸۲	۱۴
۶	۰٫۰۰۱۰	۰	۰	تطابق = ۰	۱۵

یا نقشه‌ی خودسازمانده (SOM) به‌عنوان یکی از بهترین روش‌های خوشه‌بندی یا دسته‌بندی هدایت‌نشده استفاده می‌شود. محققین بسیاری SOM را الگوریتمی کارا برای خوشه‌بندی مشتریان ارزیابی کرده‌اند.

به‌علاوه، باید توجه داشت که یک مرحله‌ی مهم در طراحی هر سامانه، مرحله‌ی ارزیابی عملکرد است که در آن احتمال خطای تقسیم‌بندی سامانه‌ی طراحی‌شده برآورد می‌شود. اگر ویژگی‌هایی با قدرت تمایز کم انتخاب شوند، طراحی سامانه منجر به عملکردی ضعیف خواهد شد. از این رو، باید ویژگی‌هایی انتخاب شوند که در فضای بردار ویژگی منجر به فاصله بین دسته‌ی بزرگ و واریانس درون دسته‌ی کوچک می‌شوند.<sup>[۲۸]</sup> به‌علاوه، توجه به این نکته ضروری است که داده‌های مورد استفاده برای استخراج ویژگی‌ها باید کاملاً از داده‌های اعتبارسنجی مستقل باشند، در غیر این صورت خطر تطابق بیش از حد وجود خواهد داشت.

#### ۱.۴. نقشه‌ی خودسازمانده

گاهی ممکن است متغیرها به‌روشی کاملاً غیرخطی به هم مرتبط باشند که در این صورت، مدل‌های تحلیلی خطی قادر به شناسایی این روابط نخواهند بود. از آنجا که SOM یک روش نمایش غیرخطی است، غالباً می‌توان چنین حالت‌ها یا خوشه‌های مشخصه‌ی را بدون تشریح مدل‌سازی سامانه در نقشه‌ی خودسازمانده مشاهده کرد. این نوع شبکه‌ی عصبی بیشترین ساختارهای مشخصه‌ی تابع چگالی ورودی را در نمایش‌گری با ابعاد کم ارائه می‌کند. از این رو SOM ابزاری قوی برای کشف و تصویرسازی ساختارهای کلی فضای حالت و رفتار سامانه است که عملیات تصویرسازی و خوشه‌بندی را با هم ترکیب می‌کند. عمده‌ترین مزایای این روش، نسبت به سایر روش‌ها، عبارت‌اند از:<sup>[۳۰]</sup>

- مجموعه‌ی باز از مشاهدات چندمتغیره را توسط مجموعه‌ی نامحدود از مشاهدات مدل مانند خوشه‌بندی K-means نمایش می‌دهد.
- مشاهدات را در قالب یک شبکه‌ی منظم دوبعدی مرتب‌شده نمایش داده و سپس مشاهدات مدل با گره‌های شبکه همبسته می‌شوند.
- محاسبه‌ی مجدد کل نقشه برای هر نمونه‌ی جدید ضروری نیست، زیرا در صورت ثابت فرض شدن آماره‌ها، نمونه‌ی جدید می‌تواند مستقیماً در قالب نزدیک‌ترین مشاهده‌ی مدل قدیمی نقشه قرار گیرد.

خوشه‌بندی  $Q$ ، به‌معنای افراز مجموعه‌ی از داده‌ها به مجموعه‌ی از خوشه‌های  $Q, z = 1, \dots, C$  است. شبکه‌ی عصبی کوهون توصیف‌شده چنین عمل می‌کند:<sup>[۳۰]</sup>

مرحله‌ی ۱: واحد خروجی برنده را به‌عنوان بزرگ‌ترین سنجه‌ی مشابهت (یا کوچک‌ترین سنجه‌ی عدم تجانس) بین تمامی بردارهای وزنی  $w_i$  و بردار ورودی  $x$  انتخاب کنید. اگر فاصله‌ی اقلیدسی به‌عنوان سنجه‌ی عدم مشابهت انتخاب شود، آنگاه واحد برنده ( $c$ ) معادله‌ی ۶ را ارضاء می‌کند.

$$\|x - w_c\| = \min \|x - w_i\| \quad (۶)$$

که در آن شاخص  $c$  به واحد برنده نسبت داده می‌شود.

مرحله‌ی ۲:  $NB_c$  مجموعه‌ی از شاخص متناظر با یک همسایه اطراف  $c$  برنده را مشخص می‌کند. وزن‌های برنده و واحدهای همسایه‌اش از طریق رابطه‌ی ۷ بروز می‌یابد:

$$\Delta w_i = \eta \gamma(i)(x - w_i), \quad i \in NB_c \quad (۷)$$

که در آن  $\eta$  نرخ یادگیری مثبت کوچک است. به‌جای تعریف همسایگی یک واحد برنده، می‌توان از یک تابع همسایگی  $\gamma(i)$  در اطراف واحد برنده‌ی  $c$  استفاده کرد.

$$\gamma(i) = \exp\left(\frac{-\|p_i - p_c\|^2}{2\sigma^2}\right) \quad (۸)$$

که در آن  $p_i - p_c$  فاصله‌ی فیزیکی بین واحد  $i$  و واحد برنده‌ی  $c$  است.

الگوریتم با افزایش تعداد دوره‌های آموزشی و کاهش بکنواخت پارامتر آموزشی اندازه‌ی همسایگی (نشان‌دهنده‌ی فاصله‌ی تمامی نرون‌های قرارگرفته در شعاعی مشخص از واحد برنده‌ی  $c$ )، کاهش نرخ یادگیری (ضریب تطبیق) و کاهش مقدار  $\sigma^2$  نمایش‌گر پهنای تابع همسایگی تکرار می‌شود تا وزن‌ها به ثبات برسند. در طول فرایند، حوزه‌ی داده‌های مرتبط به تدریج در قالب یک گروه کوچک‌تر بر مبنای مفهوم همسایگی محدود خواهد شد.

یک تعریف پذیرفته‌شده از خوشه‌بندی بهینه افزایی است که فواصل بین نمونه‌های داخلی را کمینه، و فواصل بین خوشه‌ها را بیشینه کند.<sup>[۳۱]</sup> در این نوشتار نیز برای تعیین بهترین تعداد خوشه از شاخص دیویس-بولدین استفاده می‌شود که در آن  $C$  تعداد خوشه‌ها،  $S_c$  فاصله‌ی درون‌خوشه‌ی (مجموع فواصل بین کلیه‌ی بردارهای ورودی قرارگرفته در یک خوشه از مرکز همان خوشه)، و  $d_{ce}$  فاصله‌ی بین خوشه‌ی (مجموع فواصل بین مراکز کلیه خوشه‌ها) را نمایش می‌دهد. براساس شاخص اعتبارسنجی دیویس-بولدین، بهترین خوشه‌بندی رابطه‌ی ۹ را کمینه می‌کند:

$$\frac{1}{C} \sum_{j=1}^C \max_{i \neq j} \left\{ \frac{S_c(Q_j) + S_c(Q_i)}{d_{ce}(Q_j, Q_i)} \right\} \quad (۹)$$

این شاخص هم فاصله‌ی درون‌خوشه‌ی و هم فاصله‌ی بین خوشه‌ی را هنگام ارزیابی خوشه‌بندی حاصله مورد ارزیابی قرار می‌دهد و شاخص مناسبی برای سنجش تقسیم‌بندی SOM است، زیرا هر قدر مقدار رابطه‌ی ۹ کم‌تر باشد، نشان‌دهنده‌ی خوب بودن نتایج خوشه‌بندی به‌لحاظ کروی بودن خوشه‌ها<sup>[۳۱]</sup> است.

#### ۲.۴. نرمال‌سازی داده‌ها

غالباً طراح با ویژگی‌هایی مواجه می‌شود که مقادیر در محدوده‌های متفاوتی قرار می‌گیرند. بنابراین ویژگی‌های با مقادیر بزرگ ممکن است تأثیر بیشتری به نسبت ویژگی‌های با مقادیر کوچک در تابع هزینه داشته باشند؛ اگرچه لزوماً اهمیت نسبی آن‌ها در طراحی ابزار تقسیم‌بندی را منعکس نمی‌کند. موضوع غلبه بر این مشکل از طریق نرمال‌سازی ویژگی‌هاست، به‌طوری که مقادیر آن‌ها در محدوده‌های مشابهی قرار گیرد. در این تحقیق نیز با توجه به این که، شش ویژگی منتخب بهینه مقیاس‌های متفاوتی دارند، از نرمال‌سازی با میانگین صفر و واریانس ۱ استفاده می‌شود. میانگین و انحراف معیار داده‌های ورودی با استفاده از رابطه‌ی ۱۰ محاسبه می‌شوند:

$$y_{new} = \frac{y_{old} - mean}{std} \quad (۱۰)$$

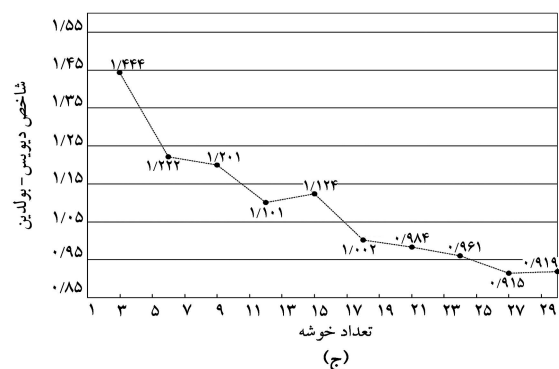
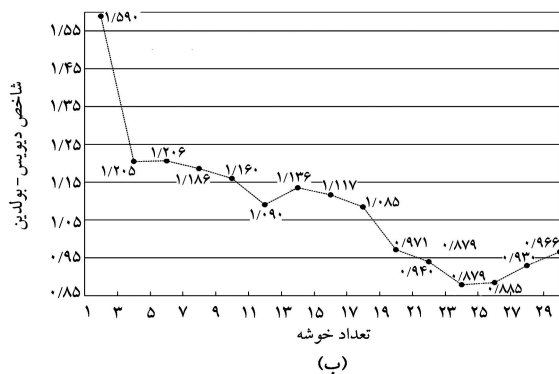
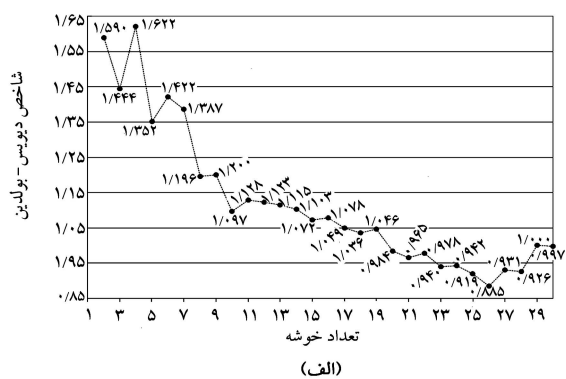
که در آن  $y_{old}$  مقدار اصلی،  $y_{new}$  مقدار جدید، و  $mean$  و  $std$  به‌ترتیب میانگین و انحراف معیار داده‌های اصلی هستند.

#### ۳.۴. الگوهای آرایه‌ی خوشه‌ها و نرم فاصله

هدف از خوشه‌بندی، تقسیم‌بندی بهینه‌ی مشتریان محصول مورد نظر، به‌منظور شناسایی ویژگی‌های مشترک محتمل‌ترین مشتریان برای تدوین استراتژی‌های بازاریابی

#### ۵.۴. نتایج خوشه‌بندی

چنان‌که پیش‌تر نیز بیان شد، SOM برای ۹۵ آرایه‌ی مختلف با به‌کارگیری الگوریتم آموزشی دسته‌بندی و در دو مرحله‌ی ترتیب و تنظیم آموزش داده شد. نتایج حاصل از آموزش هر یک از آرایه‌ها، مختصات مراکز هر یک از نرون‌ها و مکان قرارگیری هر بردار ورودی درون هر یک از نرون‌ها را ارائه می‌کند. سپس برای تعیین بهترین تعداد خوشه از میان آرایه‌های مختلف، شاخص اعتبارسنجی دیویس-بولدین از نرم فاصله‌ی اقلیدسی محاسبه می‌شود. با مقایسه‌ی نتایج حاصله برای آرایه‌های مختلف بهترین حالت خوشه‌بندی مشخص می‌شود. در شکل ۲ روند مقادیر شاخص دیویس-بولدین برای الگوی آرایه‌های دارای عنصر اول ۱ و عنصر دوم از ۲ تا ۳۰ (شکل الف)، عنصر اول ۲ و عنصر دوم از ۱ تا ۱۵ (شکل ب)، و عنصر اول ۳ و عنصر دوم از ۱ تا ۱۰ مشاهده می‌شود، کم‌ترین مقدار شاخص دیویس-بولدین از



شکل ۲. مقادیر شاخص دیویس-بولدین برای الگوی آرایه‌های دارای عنصر اول ۱، ۲ و ۳.

آنی و صرف بهینه‌ی منابع است. بنابراین افزایش تعداد خوشه‌ها تا جایی ادامه می‌یابد که بهبودی در کاهش شاخص دیویس-بولدین مشاهده نشود. به همین دلیل تقسیم‌بندی را با دو خوشه (نرون) آغاز، و تا دست‌یابی به بهترین حالت خوشه‌بندی -- با توجه به هدف تعریف شده -- ادامه داده‌ایم. در این راستا از ۲ تا ۳۰ خوشه مورد آموزش قرار گرفتند که شامل ۹۵ آرایه‌ی مختلف می‌شود: (۱ ۲)، (۱ ۳)، (۱ ۴)، (۱ ۵)، (۱ ۶)، (۱ ۷)، (۱ ۸)، (۱ ۹)، (۱ ۱۰)، (۱ ۱۱)، (۱ ۱۲)، (۱ ۱۳)، (۱ ۱۴)، (۱ ۱۵)، (۱ ۱۶)، (۱ ۱۷)، (۱ ۱۸)، (۱ ۱۹)، (۱ ۲۰)، (۱ ۲۱)، (۱ ۲۲)، (۱ ۲۳)، (۱ ۲۴)، (۱ ۲۵)، (۱ ۲۶)، (۱ ۲۷)، (۱ ۲۸)، (۱ ۲۹)، (۱ ۳۰)، (۱ ۳۱)، (۱ ۳۲)، (۱ ۳۳)، (۱ ۳۴)، (۱ ۳۵)، (۱ ۳۶)، (۱ ۳۷)، (۱ ۳۸)، (۱ ۳۹)، (۱ ۴۰)، (۱ ۴۱)، (۱ ۴۲)، (۱ ۴۳)، (۱ ۴۴)، (۱ ۴۵)، (۱ ۴۶)، (۱ ۴۷)، (۱ ۴۸)، (۱ ۴۹)، (۱ ۵۰)، (۱ ۵۱)، (۱ ۵۲)، (۱ ۵۳)، (۱ ۵۴)، (۱ ۵۵)، (۱ ۵۶)، (۱ ۵۷)، (۱ ۵۸)، (۱ ۵۹)، (۱ ۶۰)، (۱ ۶۱)، (۱ ۶۲)، (۱ ۶۳)، (۱ ۶۴)، (۱ ۶۵)، (۱ ۶۶)، (۱ ۶۷)، (۱ ۶۸)، (۱ ۶۹)، (۱ ۷۰)، (۱ ۷۱)، (۱ ۷۲)، (۱ ۷۳)، (۱ ۷۴)، (۱ ۷۵)، (۱ ۷۶)، (۱ ۷۷)، (۱ ۷۸)، (۱ ۷۹)، (۱ ۸۰)، (۱ ۸۱)، (۱ ۸۲)، (۱ ۸۳)، (۱ ۸۴)، (۱ ۸۵)، (۱ ۸۶)، (۱ ۸۷)، (۱ ۸۸)، (۱ ۸۹)، (۱ ۹۰)، (۱ ۹۱)، (۱ ۹۲)، (۱ ۹۳)، (۱ ۹۴)، (۱ ۹۵)، (۱ ۹۶)، (۱ ۹۷)، (۱ ۹۸)، (۱ ۹۹)، (۱ ۱۰۰)، (۱ ۱۰۱)، (۱ ۱۰۲)، (۱ ۱۰۳)، (۱ ۱۰۴)، (۱ ۱۰۵)، (۱ ۱۰۶)، (۱ ۱۰۷)، (۱ ۱۰۸)، (۱ ۱۰۹)، (۱ ۱۱۰)، (۱ ۱۱۱)، (۱ ۱۱۲)، (۱ ۱۱۳)، (۱ ۱۱۴)، (۱ ۱۱۵)، (۱ ۱۱۶)، (۱ ۱۱۷)، (۱ ۱۱۸)، (۱ ۱۱۹)، (۱ ۱۲۰)، (۱ ۱۲۱)، (۱ ۱۲۲)، (۱ ۱۲۳)، (۱ ۱۲۴)، (۱ ۱۲۵)، (۱ ۱۲۶)، (۱ ۱۲۷)، (۱ ۱۲۸)، (۱ ۱۲۹)، (۱ ۱۳۰)، (۱ ۱۳۱)، (۱ ۱۳۲)، (۱ ۱۳۳)، (۱ ۱۳۴)، (۱ ۱۳۵)، (۱ ۱۳۶)، (۱ ۱۳۷)، (۱ ۱۳۸)، (۱ ۱۳۹)، (۱ ۱۴۰)، (۱ ۱۴۱)، (۱ ۱۴۲)، (۱ ۱۴۳)، (۱ ۱۴۴)، (۱ ۱۴۵)، (۱ ۱۴۶)، (۱ ۱۴۷)، (۱ ۱۴۸)، (۱ ۱۴۹)، (۱ ۱۵۰)، (۱ ۱۵۱)، (۱ ۱۵۲)، (۱ ۱۵۳)، (۱ ۱۵۴)، (۱ ۱۵۵)، (۱ ۱۵۶)، (۱ ۱۵۷)، (۱ ۱۵۸)، (۱ ۱۵۹)، (۱ ۱۶۰).

یادآور می‌شود تعداد خوشه‌ها در هر آرایه از ضرب دو مقدار آرایه در یکدیگر حاصل می‌شود. همچنین برای سنجش فواصل درون‌خوشه‌ی و بین‌خوشه‌ی از معیار فاصله‌ی اقلیدسی استفاده می‌شود، بدین ترتیب که با فرض در اختیار داشتن یک ماتریس داده‌ی  $y$  ( $k * n$ ) که شامل  $k$  بردار ورودی  $y_1, y_2, \dots, y_n$  باشد، فواصل مختلف بین بردار  $y_r$  و  $y_s$  به‌عنوان فاصله‌ی اقلیدسی و با استفاده از رابطه‌ی ۱۱ محاسبه می‌شود:

$$d_{rs} = \sqrt{(y_r - y_s)^2} \quad (11)$$

#### ۴.۴. آموزش SOM

نقشه‌ی خودسازمانده با استفاده از الگوریتم آموزشی دسته‌بندی ۱۴ برای تقسیم‌بندی مجموعه داده‌های ۳۴۸ مشتری محصول مورد نظر آموزش داده شده است. در این الگوریتم کل مجموعه داده‌ها یک باره، و پیش از این که هیچ وزنی به‌روز شود، به شبکه ارائه می‌شود. سپس الگوریتم یک نرون برنده را برای هر بردار ورودی تعیین می‌کند. در ادامه، هر بردار وزنی به موقعیت میانگین کلیه بردارهای ورودی، که یک برنده است یا در همسایگی یک برنده است حرکت می‌کند.

یادآور می‌شود نرخ یادگیری و فاصله‌ی همسایگی در قالب دو مرحله تغییر داده می‌شوند: ۱. ترتیب؛ ۲. تنظیم. در این تحقیق، مرحله‌ی «ترتیب» با انجام ۱۰۰۰ گام پایان می‌پذیرد. در طول این مرحله، نرخ یادگیری از ۰/۹ تا ۰/۰۲۰ تعدیل می‌شود و فاصله‌ی همسایگی از بیشترین فاصله‌ی نرونی تا ۱ تعدیل می‌شود. در طول این مرحله انتظار می‌رود که وزن‌های نرونی خودشان را در فضای ورودی سازگار با موقعیت‌های نرونی همبسته مرتب کنند. در طول مرحله‌ی «تنظیم» نرخ یادگیری از ۰/۰۲ به کندی کاهش می‌یابد و فاصله‌ی همسایگی همیشه برابر با ۱ است. در طول این مرحله انتظار می‌رود وزن‌ها تقریباً به صورت تصادفی در کل فضای ورودی پراکنده شوند در حالی که نظم توپولوژیکی‌شان در مرحله‌ی ترتیب را حفظ می‌کنند. به این ترتیب، نقشه‌های ویژگی ۱۵ در حین یادگیری گروه‌بندی ورودی‌شان، توپولوژی و توزیع ورودی‌شان را نیز یاد می‌گیرند.<sup>[۳۳]</sup> همچنین، هر یک از شبکه‌های عصبی SOM با ۲۰۰۰ دوره آموزشی مورد آموزش قرار می‌گیرد.

نیز مقدارش کم‌تر است -- اگرچه هر دو آرایه ۲۴ خوشه را به‌عنوان حالت بهینه برای خوشه‌بندی مشتریان نمایش می‌دهند. بنابراین، الگوی آرایه‌ی (۶ ۴) به‌عنوان بهینه‌ترین حالت برای تقسیم‌بندی مشتریان تعیین می‌شود. در شکل ۳، چیدمان بهینه‌ی قرارگیری خوشه‌های مشتریان مبتنی بر شش ویژگی بهینه و تعداد مشتریان هر خوشه در مجموعه‌ی آموزشی نمایش داده شده است.

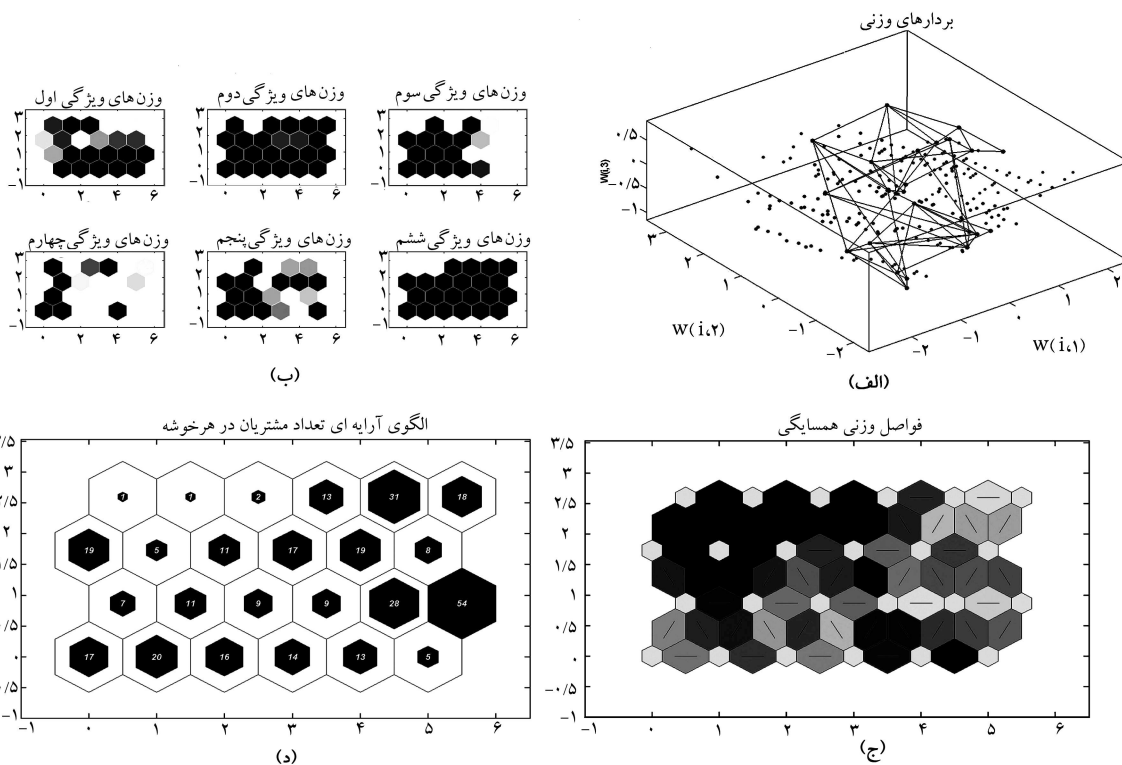
در شکل ۳ الف تعداد ۳۴۸ بردار ورودی (مشتری) و چگونگی خوشه‌بندی فضای ورودی با چیدمان ۲۴ بردار وزنی هر نرون (مرکز یک خوشه از بردارهای ورودی) متصل شده به نرون‌های همسایه نشان داده شده است. با توجه به وجود شش

میان ۲۹ آرایه‌ی اجراشده با عنصر اول ۱ مربوط به آرایه‌ی (۱ ۲۶) با مقدار ۰/۸۸۵، میان ۱۵ آرایه با عنصر اول ۲ مربوط به آرایه‌ی (۲ ۱۲) با مقدار ۰/۸۷۹ و میان ۱۰ آرایه با عنصر اول ۳ مربوط به آرایه‌ی (۳ ۹) با مقدار ۰/۹۱۵ است که میان این سه آرایه، الگوی آرایه‌ی (۲ ۱۲) مقدار کم‌تری را ارائه می‌کند. همچنین محاسبات شاخص دیویس - بولدین برای ۴۱ آرایه‌ی دیگر نامبرده نیز انجام شد. جدول ۲ مقادیر شاخص دیویس - بولدین برای سایر الگوهای آرایه‌ی را نشان می‌دهد.

براساس جدول ۲، محاسبات نشان می‌دهد که میان این آرایه‌ها، الگوی آرایه‌ی (۶ ۴) با مقدار ۰/۸۷۷ کم‌ترین مقدار را دارد که در مقایسه با آرایه‌ی (۲ ۱۲)

جدول ۲. مقادیر شاخص دیویس - بولدین برای سایر الگوهای آرایه‌ی.

عنصر دوم	عنصر اول	۱	۲	۳	۴	۵	۶	۷
۴	۱/۶۲۲	۱/۱۳۲	۱/۱۶۴	۱/۱۳۲	۱/۰۵۷	۱/۰۶۱	۰/۹۶۲	۰/۹۱۶
۵	۱/۳۵۲	۱/۰۸۳	۱/۰۹۲	۱/۰۸۳	۱/۰۲۷	۰/۹۱۸	۰/۹۹۹	
۶	۱/۴۲۲	۱/۱۳۳	۱/۰۹۸	۱/۰۹۸	۰/۸۷۷	۰/۸۸۱		
۷	۱/۳۸۷	۱/۰۷۹	۱/۰۷۹	۱/۰۱۵	۰/۸۸۸			
۸	۱/۱۹۶	۱/۰۱۹	۰/۹۲۷					
۹	۱/۲۰۰	۱/۰۶۴	۰/۹۳۱					
۱۰	۱/۰۹۷	۱/۰۵۰	۰/۹۸۰					
۱۱	۱/۱۲۸	۰/۹۳۷						
۱۲	۱/۱۲۳	۰/۹۶۲						
۱۳	۱/۱۱۵	۰/۸۸۸						
۱۴	۱/۱۰۳	۰/۹۲۴						
۱۵	۱/۰۷۲	۰/۸۸۲						



شکل ۳. چیدمان بهینه‌ی قرارگیری خوشه‌های مشتریان مبتنی بر شش ویژگی بهینه و تعداد مشتریان هر خوشه در مجموعه‌ی آموزشی.



جدول ۳. ویژگی‌های مشتریان هر خوشه.

شماره‌ی خوشه	ویژگی‌های مشتریان
۱	نوع مشتری (۶-۱۰)، کارگر بی‌تجربه (۲-۰)، بیمه شخص ثالث (۲،۰)، بیمه‌نامه‌ی خودرو (۰)، بیمه‌نامه‌ی آتش‌سوزی (۴،۳،۰)، بیمه‌نامه‌ی قایق (۰)
۲	نوع مشتری (۳-۱)، کارگر بی‌تجربه (۳-۰)، بیمه شخص ثالث (۲،۰)، بیمه‌نامه‌ی خودرو (۰)، بیمه‌نامه‌ی آتش‌سوزی (۳،۰)، بیمه‌نامه‌ی قایق (۰)
۳	نوع مشتری (۳-۱)، کارگر بی‌تجربه (۱،۰)، بیمه شخص ثالث (۰)، بیمه‌نامه‌ی خودرو (۶)، بیمه‌نامه‌ی آتش‌سوزی (۰)، بیمه‌نامه‌ی قایق (۰)
۴	نوع مشتری (۳-۱)، کارگر بی‌تجربه (۱،۰)، بیمه شخص ثالث (۱،۰)، بیمه‌نامه‌ی خودرو (۶،۵)، بیمه‌نامه‌ی آتش‌سوزی (۴،۳)، بیمه‌نامه‌ی قایق (۰)
۵	نوع مشتری (۳-۱)، کارگر بی‌تجربه (۳-۰)، بیمه شخص ثالث (۲-۰)، بیمه‌نامه‌ی خودرو (۰)، بیمه‌نامه‌ی آتش‌سوزی (۴-۲)، بیمه‌نامه‌ی قایق (۰)
۶	نوع مشتری (۳-۱)، کارگر بی‌تجربه (۲-۰)، بیمه شخص ثالث (۲)، بیمه‌نامه‌ی خودرو (۵،۶)، بیمه‌نامه‌ی آتش‌سوزی (۰،۱)، بیمه‌نامه‌ی قایق (۰)
۷	نوع مشتری (۹-۷)، کارگر بی‌تجربه (۴،۵)، بیمه شخص ثالث (۰)، بیمه‌نامه‌ی خودرو (۰)، بیمه‌نامه‌ی آتش‌سوزی (۰)، بیمه‌نامه‌ی قایق (۰)
۸	نوع مشتری (۱۳،۵)، کارگر بی‌تجربه (۴-۲)، بیمه شخص ثالث (۰)، بیمه‌نامه‌ی خودرو (۶)، بیمه‌نامه‌ی آتش‌سوزی (۱،۰)، بیمه‌نامه‌ی قایق (۱،۰)
۹	نوع مشتری (۳،۵-۱)، کارگر بی‌تجربه (۲،۳،۵)، بیمه شخص ثالث (۰)، بیمه‌نامه‌ی خودرو (۶)، بیمه‌نامه‌ی آتش‌سوزی (۳،۴،۶)، بیمه‌نامه‌ی قایق (۰)
۱۰	نوع مشتری (۳-۱)، کارگر بی‌تجربه (۴،۵)، بیمه شخص ثالث (۲)، بیمه‌نامه‌ی خودرو (۶)، بیمه‌نامه‌ی آتش‌سوزی (۵-۳)، بیمه‌نامه‌ی قایق (۰)
۱۱	نوع مشتری (۳،۵-۱)، کارگر بی‌تجربه (۲،۳)، بیمه شخص ثالث (۲)، بیمه‌نامه‌ی خودرو (۶)، بیمه‌نامه‌ی آتش‌سوزی (۴-۲)، بیمه‌نامه‌ی قایق (۰)
۱۲	نوع مشتری (۳،۵-۱)، کارگر بی‌تجربه (۰،۱)، بیمه شخص ثالث (۲،۳)، بیمه‌نامه‌ی خودرو (۶)، بیمه‌نامه‌ی آتش‌سوزی (۵-۳)، بیمه‌نامه‌ی قایق (۰)
۱۳	نوع مشتری (۹،۰۶،۸)، کارگر بی‌تجربه (۳-۰)، بیمه شخص ثالث (۰)، بیمه‌نامه‌ی خودرو (۶،۵)، بیمه‌نامه‌ی آتش‌سوزی (۰)، بیمه‌نامه‌ی قایق (۰)
۱۴	نوع مشتری (۹،۰۳،۸)، کارگر بی‌تجربه (۲-۰)، بیمه شخص ثالث (۰)، بیمه‌نامه‌ی خودرو (۶،۰)، بیمه‌نامه‌ی آتش‌سوزی (۰)، بیمه‌نامه‌ی قایق (۴-۲)
۱۵	نوع مشتری (۶-۱۰)، کارگر بی‌تجربه (۳-۰)، بیمه شخص ثالث (۲)، بیمه‌نامه‌ی خودرو (۶،۵)، بیمه‌نامه‌ی آتش‌سوزی (۶-۳)، بیمه‌نامه‌ی قایق (۲-۰)
۱۶	نوع مشتری (۹-۵،۷)، کارگر بی‌تجربه (۸،۵،۷،۴)، بیمه شخص ثالث (۰)، بیمه‌نامه‌ی خودرو (۶)، بیمه‌نامه‌ی آتش‌سوزی (۰،۳)، بیمه‌نامه‌ی قایق (۰)
۱۷	نوع مشتری (۹-۳،۵،۷)، کارگر بی‌تجربه (۶-۴)، بیمه شخص ثالث (۲،۱)، بیمه‌نامه‌ی خودرو (۶،۵)، بیمه‌نامه‌ی آتش‌سوزی (۴-۲،۰)، بیمه‌نامه‌ی قایق (۰)
۱۸	نوع مشتری (۹-۵)، کارگر بی‌تجربه (۲-۰)، بیمه شخص ثالث (۲)، بیمه‌نامه‌ی خودرو (۶،۵)، بیمه‌نامه‌ی آتش‌سوزی (۲-۰)، بیمه‌نامه‌ی قایق (۰)
۱۹	نوع مشتری (۷)، کارگر بی‌تجربه (۲)، بیمه شخص ثالث (۲)، بیمه‌نامه‌ی خودرو (۰)، بیمه‌نامه‌ی آتش‌سوزی (۴)، بیمه‌نامه‌ی قایق (۶)
۲۰	نوع مشتری (۵)، کارگر بی‌تجربه (۶)، بیمه شخص ثالث (۰)، بیمه‌نامه‌ی خودرو (۶)، بیمه‌نامه‌ی آتش‌سوزی (۳)، بیمه‌نامه‌ی قایق (۶)
۲۱	نوع مشتری (۱،۹)، کارگر بی‌تجربه (۱،۳)، بیمه شخص ثالث (۲)، بیمه‌نامه‌ی خودرو (۵)، بیمه‌نامه‌ی آتش‌سوزی (۴)، بیمه‌نامه‌ی قایق (۴)
۲۲	نوع مشتری (۹-۷)، کارگر بی‌تجربه (۴،۶،۷-۲،۰)، بیمه شخص ثالث (۲-۰)، بیمه‌نامه‌ی خودرو (۰)، بیمه‌نامه‌ی آتش‌سوزی (۳،۴)، بیمه‌نامه‌ی قایق (۰)
۲۳	نوع مشتری (۹-۷)، کارگر بی‌تجربه (۳،۲)، بیمه شخص ثالث (۳-۱)، بیمه‌نامه‌ی خودرو (۵،۶)، بیمه‌نامه‌ی آتش‌سوزی (۵-۳)، بیمه‌نامه‌ی قایق (۰)
۲۴	نوع مشتری (۱۰-۸)، کارگر بی‌تجربه (۱،۰)، بیمه شخص ثالث (۲)، بیمه‌نامه‌ی خودرو (۵،۶)، بیمه‌نامه‌ی آتش‌سوزی (۴،۶،۳)، بیمه‌نامه‌ی قایق (۰)

بر اساس اطلاعات در دسترس، اعداد داخل پرانتز درمورد کارگر بی‌تجربه از ۰ تا ۹؛ درمورد بیمه‌نامه‌ی خودرو شامل ۰ تا ۴؛ درمورد بیمه‌نامه‌ی آتش‌سوزی از ۰ تا ۸؛ و درمورد بیمه‌نامه‌ی قایق از ۰ تا ۶ هستند. به علاوه، دو جدول ۴ و ۵ دسته‌بندی نوع مشتری و مقادیر بیمه شخص ثالث را توصیف می‌کنند.

همچنین به منظور صرف بهینه‌ی منابع شرکت در قالب هزینه‌های بازاریابی، باید از شاخصی برای سنجش آن بهره جست. شاخص PPC که از آن برای محاسبه‌ی

ویژگی (بعد)، از تحلیل مؤلفه‌ی اصلی برای نمایش موقعیت وزنی خوشه‌ها و هریک از مشتریان در فضای سه بعدی استفاده شده است. در شکل ۳، یک صفحه‌ی وزنی برای هریک از شش عنصر بردار ورودی (نوع مشتری، کارگر بی‌تجربه، بیمه‌ی شخص ثالث شخصی، بیمه‌نامه‌ی خودرو، بیمه‌نامه‌ی آتش‌سوزی و بیمه‌نامه‌ی قایق به ترتیب ورودی) به صورت جداگانه در یک فضای دوبعدی نشان داده شده است. هریک از عناصر تصویری از وزن‌ها هستند که هر ورودی را به یکی از نرون‌ها متصل می‌کند. هر قدر رنگ خوشه تیره‌تر باشد، نشان‌دهنده‌ی وزن بیشتر است. اگر الگوهای اتصال دو ورودی شباهت بسیاری به یکدیگر داشته باشند، می‌توان نتیجه گرفت که آن دو ورودی همبستگی شدیدی دارند. اما در خصوص این شش ورودی اتصالات بسیار متفاوت‌اند. در شکل ۳ نیز فاصله‌ی بردارهای وزنی (مرکز خوشه‌ها) هر نرون از نرون‌های همسایه با لحاظ هر شش ویژگی به صورت همزمان نشان داده شده است. میزان تیرگی رنگ فاصله‌ی بین خوشه‌ها، نشان‌دهنده‌ی فاصله بیشتر مرکز خوشه از خوشه‌های مجاور است. به طور کلی، فاصله‌ی همسایگی میان خوشه‌های قرار گرفته در سمت چپ بالای تصویر بسیار بیشتر از سایر مناطق شکل فاصله همسایگی است. بین دو خوشه‌ی ۱۳ و ۱۹ و دو خوشه‌ی ۱۵ و ۲۰ که فاصله‌ی همسایگی در آن‌ها تیره‌تر نشان داده شده، بیشترین فاصله بین مراکز خوشه‌ها دیده می‌شود. در شکل ۳، تعداد بردارهای ورودی در برگرفته توسط هر خوشه نشان داده شده است. در میان خوشه‌ها، خوشه‌ی ۱۲ با ۵۴ بردار ورودی بیشترین تعداد مشتری را شامل می‌شود. جدول ۳ ویژگی‌های مشتریان هر خوشه را نشان می‌دهد.

جدول ۴. نوع مشتری.

ردیف	نوع مشتری
۱	لذت‌جوهای موفق
۲	پرورش دهندگان گیاهان
۳	خانواده‌ی متوسط
۴	افراد تنها در کل زندگی
۵	افراد بدون نگرانی مالی
۶	افراد مهم در حال سفر
۷	بازنشسته و دیندار
۸	خانواده‌یی با افراد بزرگسال
۹	خانواده‌های محافظه‌کار
۱۰	کشاورزان

جدول ۵. مقادیر بیمه شخص ثالث.

مقادیر بیمه	ردیف
۰	۰
۴۹-۱	۱
۹۹-۵۰	۲
۱۹۹-۱۰۰	۳

جدول ۶. بهینه‌ترین خوشه‌ها برای صرف هزینه‌های بازاریابی و خوشه‌بندی مشتریان مجموعه داده‌های ارزیابی.

شماره‌ی خوشه	مجموعه آموزشی			مجموعه اعتبارسنجی		
	PC	TC	PPC (%)	PC	TC	PPC (%)
۱	۱۷	۵۹۷	۲٫۸	۱۲	۴۱۸	۲٫۹
۲	۲۰	۶۱۷	۳٫۲	۱۲	۴۳۴	۲٫۸
۳	۱۶	۲۵۷	۶٫۲	۸	۱۷۳	۴٫۶
۴	۱۴	۷۷	۱۸٫۲	۵	۵۶	۸٫۹
۵	۱۳	۴۴۴	۲٫۹	۲۱	۲۹۵	۷٫۱
۶	۵	۳۳	۱۵٫۲	۰	۲۴	۰٫۰
۷	۷	۵۶۴	۱٫۲	۶	۳۹۳	۱٫۵
۸	۱۱	۲۷۲	۴٫۰	۵	۱۸۴	۲٫۷
۹	۹	۶۹	۱۳٫۰	۱۰	۵۰	۲۰٫۰
۱۰	۹	۳۸	۲۳٫۷	۷	۳۲	۲۱٫۹
۱۱	۲۸	۲۰۰	۱۴٫۰	۱۱	۱۴۰	۷٫۹
۱۲	۵۴	۲۵۷	۲۱٫۰	۳۷	۱۶۸	۲۲٫۰
۱۳	۱۹	۴۷۵	۴٫۰	۱۷	۳۴۸	۴٫۹
۱۴	۵	۱۱	۴۵٫۵	۱	۲	۵۰٫۰
۱۵	۱۱	۱۶۹	۶٫۵	۸	۱۱۶	۶٫۹
۱۶	۱۷	۲۸۸	۵٫۹	۱۰	۱۶۸	۶٫۰
۱۷	۱۹	۲۵۴	۷٫۵	۱۲	۱۵۰	۸٫۰
۱۸	۸	۱۲۶	۶٫۳	۶	۷۵	۸٫۰
۱۹	۱	۴	۲۵٫۰	۰	۰	-
۲۰	۱	۱	۱۰۰	۰	۰	-
۲۱	۲	۱۰	۲۰٫۰	۲	۷	۲۸٫۶
۲۲	۱۳	۶۰۷	۲٫۱	۱۵	۴۳۸	۳٫۴
۲۳	۳۱	۲۹۸	۱۰٫۴	۱۹	۱۹۸	۹٫۶
۲۴	۱۸	۱۵۴	۱۱٫۷	۱۴	۱۳۱	۱۰٫۷
مجموع	۳۴۸	۵۸۲۲	۵٫۹۸	۲۳۸	۴۰۰۰	۵٫۹۵
۲۰٪ برتر	۱۷۷	۱۱۵۲	۱۵٫۳۶	۱۰۶	۸۰۸	۱۳٫۱۲

درصد مشتریان محصول بیمه‌نامه‌ی کاروان در هر خوشه استفاده شده، در قالب رابطه‌ی ۱۲ تعریف می‌شود:

$$PPC = \frac{PC}{TC} \quad (12)$$

که در آن  $PC$  نشان‌گر مشتریان محصول در هر خوشه، و  $TC$  نشان‌گر کل مشتریان شرکت در هر خوشه است. هرچه مقدار این شاخص بزرگ‌تر باشد، صرف هزینه‌های بازاریابی برای این قبیل از مشتریان با احتمال برگشت سرمایه‌ی بیشتری همراه است. به‌علاوه، چنان‌که پیش‌تر نیز بیان شد، برای اطمینان از عملکرد مناسب سامانه‌ی هوشمند طراحی‌شده، لازم است عملکرد سامانه مورد سنجش قرار گیرد. بدین منظور

از مجموعه داده‌های ارزیابی برای اعتبارسنجی سامانه استفاده می‌شود. در جدول ۶ بهینه‌ترین خوشه‌ها برای صرف هزینه‌های بازاریابی و خوشه‌بندی مشتریان مجموعه داده‌های ارزیابی مبتنی بر بردارهای وزنی تعیین شده برای هر خوشه در مجموعه آموزشی نمایش داده شده است.

چنان‌که در جدول ۶ مشاهده می‌شود، براساس بردارهای وزنی بدست آمده برای ۲۴ خوشه، اقدام به خوشه‌بندی ۵۸۲۲ مشتری کلیه محصولات شرکت در مجموعه آموزشی کرده‌ایم. بدین ترتیب هر مشتری با توجه به میزان شباهت ویژگی‌هایش به مشتریان هر یک از خوشه‌ها در آن خوشه قرار می‌گیرد. ستون سوم جدول تعداد کل مشتریان در هر خوشه را نشان می‌دهد. سپس، از تقسیم مشتریان محصول بیمه‌نامه‌ی کاروان در هر خوشه بر تعداد کل مشتریان آن خوشه، درصد مشتریان محصول بیمه‌نامه‌ی کاروان در آن خوشه حاصل می‌شود. با مشاهده‌ی ستون چهارم جدول می‌توان متوجه شد که خوشه‌های ۴، ۶، ۹، ۱۰، ۱۱، ۱۲، ۱۴، ۱۹، ۲۰، ۲۱، ۲۳ و ۲۴ بالاترین درصد مشتریان بیمه‌نامه‌ی کاروان را در بر می‌گیرند. به‌طوری‌که، این ۱۲ خوشه با داشتن نزدیک به ۲۰٪ (۱۱۵۲/۵۸۲۲) از کل مشتریان، ۹۰٪ (۱۷۷/۳۴۸) از مشتریان محصول بیمه‌نامه‌ی کاروان را پوشش می‌دهند. بدین ترتیب، درصد مجموع مشتریان محصول بیمه‌نامه‌ی کاروان به کل مشتریان شرکت با بهبودی قابل توجه از ۵٫۹۸٪ به ۱۵٫۳۶٪ افزایش می‌یابد. در نتیجه، شرکت می‌تواند با تمرکز بر مشتریانی با ویژگی‌های مشابه با این ۱۲ خوشه منابع بازاریابی خود را به‌صورت بهینه‌تری هزینه کند. در مقابل، تمرکز بر مشتریان شش خوشه‌ی ۱، ۲، ۷، ۸، ۱۳ و ۲۲ تنها هزینه‌بر خواهد بود، زیرا مقدار شاخص در این خوشه‌ها کم‌تر از مقدار متوسط (یعنی ۵٫۹۸٪) است.

عموماً افرادی که بیمه‌نامه‌ی خودرو خریداری می‌کنند، محتمل‌ترین افراد برای خرید بیمه‌نامه‌ی کاروان هستند و افرادی که مقادیر ۴، ۷ یا ۸ به بیمه‌نامه‌ی خودرو آن‌ها تخصیص یافته، قطعاً مشتریان این محصول نیستند. به‌علاوه، احتمال خرید گروه «افراد تنها در کل زندگی» وجود نداشته است. یکی از بزرگ‌ترین شگفتی‌ها در تحلیل ما این حقیقت بود که گروه «افراد مهم در حال سفر» بعید به نظر می‌رسد که اقدام به خرید بیمه‌نامه‌ی کاروان کنند. اگر ویژگی کارگر بی‌تجربه مقدار ۹ داشته باشد، احتمال خرید وجود ندارد. اگر بیمه‌نامه‌ی آتش‌سوزی مقدار ۷ یا ۸ بگیرد، احتمال خرید برابر صفر است همچنان‌که اگر بیمه‌نامه‌ی قایق مقدار بیش از ۲ بگیرد، احتمال خرید بسیار زیاد است. البته این تحلیل‌ها اطلاعاتی کلی را ارائه می‌کنند و اطلاعات دقیق‌تر از داخل هر خوشه استخراج می‌شود.

در ادامه، خوشه‌بندی مشتریان مبتنی بر بردارهای وزنی تعیین شده برای هر خوشه در مجموعه آموزشی بر روی ۴۰۰۰ مشتری شرکت در مجموعه داده ارزیابی نیز انجام شد. باید توجه داشت که همچون مجموعه‌ی آموزشی خوشه‌های ۴، ۹، ۱۰، ۱۱، ۱۲، ۱۴، ۲۱، ۲۳ و ۲۴ بالاترین مقدار شاخص را به خود اختصاص می‌دهند. با توجه به کم‌بودن مراجعه‌ی مشتریانی با ویژگی‌های مشابه با مشتریان خوشه‌های ۱۹ و ۲۰، در مجموعه ارزیابی چنین مشتریانی وجود نداشته‌اند. اما نمی‌توان این دو خوشه را از میان خوشه‌های مناسب حذف کرد، زیرا به‌احتمال فراوان، این مشتریان در صورت مراجعه به شرکت، مشتری محصول بیمه‌نامه‌ی کاروان نیز خواهند بود. لیکن، تنها درمورد مشتریان خوشه‌ی ۶ باید تأمل بیشتری کرد. بهتر است در تخصیص منابع بازاریابی مشتریان این خوشه پس از ۱۱ خوشه‌ی مناسب‌تر قرار گیرد. در مجموع، ۱۲ خوشه‌ی نامبرده با پوشش نزدیک به ۲۰٪ (۸۰۸/۴۰۰۰) از کل مشتریان ۴۴٫۵٪ (۱۰۶/۲۳۸) از مشتریان محصول بیمه‌نامه‌ی کاروان را دارند. با این اوصاف، درصد مجموع مشتریان محصول بیمه‌نامه‌ی کاروان به کل مشتریان شرکت با بهبودی چشمگیر از ۵٫۹۵٪ به ۱۳٫۱۲٪ افزایش می‌یابد. درخصوص

ویژگی‌ها شد. این رویه همچنین سرعت و دقت محاسباتی بسیار زیادی در انتخاب ویژگی دارد. سپس تقسیم‌بندی مشتریان با توسعه‌ی SOM بر اساس ویژگی‌های منتخب انجام شد، به طوری که به کارگیری SOM مبتنی بر یک شاخص مناسب نتایج قابل توجهی در بر دارد. نتایج نشان داد امکان صرف منابع بازاریابی کم‌تر برای دستیابی به مشتریانی با احتمال جذب بیشتر وجود دارد. با انتخاب نیمی از خوشه‌ها شامل ۲۰٪ از کل مشتریان می‌توان حدود ۵۰٪ از مشتریان بیمه‌نامه‌ی کاروان را مورد پوشش قرار داد. همچنین به منظور ارزیابی اثربخشی سامانه‌ی هوشمند از یک مجموعه داده‌ی ارزیابی استفاده شد که انتخاب صحیح بهینه‌ترین خوشه‌ها در مجموعه‌ی آموزشی برای صرف کارا تر منابع را مورد تأیید قرار داد. بدین ترتیب، شرکت بیمه می‌تواند با تدوین راهکارهای بازاریابی خود برای خوشه‌های بهینه، انتظار کسب بیشترین سودآوری را داشته باشد.

نگارندگان این مقاله مدل خود را روی داده‌های یک شرکت بیمه اجرا کردند. اما یکی دیگر از حوزه‌های مهم کاربردی این مدل داده‌های بانک است که می‌تواند کمک شایانی به بانک در شناسایی سودآورترین مشتریان کند. استفاده از شبکه‌های عصبی به عنوان یک تکنیک غیرخطی برای انتخاب ویژگی بسیار زمان‌بر است، اما ممکن است منجر به انتخاب دقیق‌تر ویژگی‌ها شود. همچنین، نویسندگان بر این باورند که روش خوشه‌بندی فازی  $c - means$  به عنوان یک تکنیک خوشه‌بندی به‌تنهایی نمی‌تواند نتیجه‌ی خوبی در بر داشته باشد، اما ترکیب آن با SOM ممکن است منجر به نتایج بهتری شود که نیاز به تحقیقات بیشتری دارد.

## ۶. تقدیر و تشکر

بدینوسیله از پشتیبانی مالی مرکز تحقیقات مخابرات ایران (ITRC) در پیشبرد اهداف این تحقیق تقدیر و تشکر به عمل می‌آید.

خوشه‌های هزینه‌بر نیز مجموعه‌ی اعتبارسنجی بر دقت تشخیص شش خوشه‌ی نامبرده صحه می‌گذارد. بنابراین، مجموعه داده‌ی ارزیابی مؤید کارآمدی و دقت قابل توجه خوشه‌بندی حاصله است. بدین ترتیب شرکت می‌تواند بنا بر ویژگی‌های مشتریان بهینه‌ترین خوشه‌ها اقدام به تدوین استراتژی‌های بازاریابی خود برای ترغیب و جذب مشتریان مشابه کند. در صورتی که منابع بازاریابی شرکت محدود نباشد، شش خوشه‌ی (با احتمال جذب کم‌تر) که نام آن‌ها آورده نشد می‌تواند گزینه‌های بعدی آن باشند.

## ۵. نتیجه‌گیری

پیشرفت‌های اخیر در انتخاب ویژگی مسئله‌ی بهبود عملکرد تقسیم‌بندی مشتریان را از نقطه‌نظر عملی مورد توجه قرار می‌دهد. موضوعی که به‌ندرت در ادبیات به آن پرداخته‌اند. در این نوشتار، تلاش کردیم یک سامانه‌ی هوشمند تلفیقی برای تقسیم‌بندی بهینه‌ی مشتریان یک محصول مشخص، به نام بیمه‌نامه‌ی کاروان، توسعه دهیم. رویکرد جدیدی برای تقسیم‌بندی مشتریان با طراحی یک شبکه عصبی SOM بهینه‌شده مبتنی بر شش ویژگی منتخب بهینه با به کارگیری درخت رگرسیونی هرس شده پیشنهاد شده، به طوری که کم‌ترین نرخ میانگین مربعات خطا را برای تعیین بهترین درخت رگرسیونی هرس شده و شاخص دیویس-بولدین را برای تعیین بهترین الگوی آرایه‌ی لحاظ کردیم. این سامانه‌ی هوشمند تلفیقی به یک بهینه‌سازی ترکیبی نامحدود منجر شد که در آن نرخ MSE و شاخص نامبرده معیار جست‌وجو هستند. معیار کم‌ترین میانگین مربعات خطا برای تعیین نقطه‌ی بهینه برای هرس درخت رگرسیونی به کار برده شده بود که نهایتاً منجر به انتخاب شش ویژگی شد. به کارگیری این رویه‌ی انتخاب ویژگی پیشنهادی منجر به بهبودی قابل توجه در کاهش

## پانویس‌ها

1. customer relationship management
2. self organizing map
3. nested partition
4. simulated annealing
5. rough set
6. support vector machine
7. Davies-Bouldin index
8. classification and regression tree
9. pruning algorithm
10. cross-validation
11. goodness-of-fit
12. leave-one-out
13. spherical clusters
14. batch training algorithm
15. feature maps

## منابع (References)

1. Anderson, E.T. "Sharing the wealth: When should firms treat customers as partners?", *Management Science*, **48**(8), pp. 955-71 (2002).
2. Ha, S.H. "Applying knowledge engineering techniques to customer analysis in the service industry", *Advanced Engineering Informatics*, **21**, pp. 293-301 (2007).
3. Hsieh, N.C. "An integrated data mining and behavioral scoring model for analyzing bank customers", *Expert Systems with Applications*, **27**, pp. 623-633 (2004).
4. Jonker, J.J.; Piersma, N. and Poel, D.V. "Joint optimization of customer segmentation and marketing policy to maximize long-term profitability", *Expert Systems with Applications*, **27**, pp. 159-168 (2004).

5. Hwang, H.; Jung, T. and Suh, E. "An LTV model and customer segmentation based on customer value: A case study on the wireless telecommunication industry", *Expert Systems with Applications*, **26**, pp. 181-188 (2004).
6. Kim, S.Y.; Jung, T.S.; Suh, E.H. and Hwang, H.S. "Customer segmentation and strategy development based on customer lifetime value: A case study", *Expert Systems with Applications*, **31**, pp. 101-107 (2006).
7. Chan, C.C.H. "Intelligent value-based customer segmentation method for campaign management: A case study of automobile retailer", *Expert Systems with Applications*, **34**, pp. 2754-2762 (2008).
8. Lee, J.H. and Park, S.Ch., "Intelligent profitable customers segmentation system based on business intelligence tools", *Expert Systems with Applications*, **29**, pp. 145-152 (2005).
9. Chen C.H.; Khoo, L.P. and Yan, W. "A strategy for acquiring customer requirement patterns using laddering technique and ART2 neural network", *Advanced Engineering Informatics*, **16**, pp. 229-240 (2002).
10. Shieh, M.D.; Yan, W. and Chen, C.H., "Soliciting customer requirements for product redesign based on picture sorts and ART2 neural network", *Expert Systems with Applications*, **34**, pp. 194-204 (2008).
11. Espinoza, M.; Joye, C.; Belmans, R. and Moor, B.D. "Short-term load forecasting, profile identification, and customer segmentation: A methodology based on periodic time series", *IEEE Transactions on Power Systems*, **20**(3), pp. 1622-1630 (2005).
12. Hu, T.L. and Sheu, J.B. "A fuzzy-based customer classification method for demand-responsive logistical distribution operations", *Fuzzy Sets and Systems*, **139**, pp. 431-450 (2003).
13. Kohavi, R. and John, G.H. "Wrappers for feature subset selection", *Artificial Intelligence*, (1-2), pp. 273-324 (1997).
14. Blum, A.L. and Rivest, R.L. "Training a 3-node neural network is NP-complete", *Neural Networks*, **5**, pp. 117-127 (1992).
15. Ng, K.S. and Liu, H. "Customer retention via data mining", *Artificial Intelligence Review*, **14**(6), pp. 569-590 (2000).
16. Kim, Y.S. "Toward a successful CRM: Variable selection, sampling, and ensemble", *Decision Support Systems*, **41**, pp. 542-553 (2006).
17. Hadden, J.; Tiwari, A.; Roy, R. and Ruta, D. "Computer assisted customer churn management: State-of-the-art and future trends", *Computers & Operations Research*, **34**, pp. 2902-2917 (2005).
18. Kim, Y.S. and Street, W.N. "An intelligent system for customer targeting: A data mining approach", *Decision Support Systems*, **37**, pp. 215-228 (2004).
19. Yu, E. and Cho, S. "Constructing response model using ensemble based on feature subset selection", *Expert Systems with Applications*, **30**, pp. 352-360 (2006).
20. Ahn, H.; Kim, K. and Han, I. "A case-based reasoning system with the two-dimensional reduction technique for customer classification", *Expert Systems with Applications*, **32**, pp. 1011-1019 (2007).
21. Yan, L. and Changrui, Y. "A new hybrid algorithm for feature selection and its application to customer recognition", *LNCS*, **4616**, pp. 102-111 (2007).
22. Tseng, T.L. and Huang, C.C. "Rough set-based approach to feature selection in customer relationship management", *Omega*, **35**, pp. 365-383 (2007).
23. Lessmann, S. and Voß, S. "A reference model for customer-centric data mining with support vector machines", *European Journal of Operational Research*, **199**, pp. 520-530 (2009).
24. Han, J. and Kamber, M., *Data Mining: Concepts and Techniques*, Morgan Kaufmann, Second Edition (2006).
25. Wu, X.; Kumar, V.; Quinlan, J.R.; Ghosh, J.; Yang, Q.; Motoda, H.; McLachlan, G.J.; Ng, A.; Liu, B.; Yu, Ph.S.; Zhou, Zh.H.; Steinbach, M.; Hand, D.J. and Steinberg, D. "Top 10 algorithms in data mining", *Knowledge Information System*, **14**, pp. 1-37 (2008).
26. Liu, H.H. and Ong, Ch.Sh. "Variable selection in clustering for marketing segmentation using genetic algorithms", *Expert Systems with Applications*, **34**, pp. 502-510 (2008).
27. Breiman, L.; Friedman, J.; Olshen, R. and Stone, C. "Classification and regression trees", Boca Raton, FL: CRC Press (1984).
28. Theodoridis, S. and Koutroumbas, K. "Pattern recognition", *Academic Press*, Third Edition (2006).
29. Webb, A.R., *Statistical Pattern Recognition*, John Wiley & Sons, Ltd., Second Edition (2002).
30. Kohonen, T. "Self-organizing maps", *Springer Series*, Third Edition (2001).
31. Vesanto, J. and Alhoniemi, E. "Clustering of the self-organizing map", *IEEE Transactions on Neural Networks*, **11**(3), pp. 586-600 (2000).
32. Demuth, H.; Beale, M. and Hagan, M., *Neural Network Toolbox<sup>TM</sup> User's Guide*, Version 6.0.3, The MathWorks, Inc. (2009).
33. Hagan, M.; Demuth, H. and Beale, M., *Neural Network Design*, USA, PWS Publishing Company, First Edition (1996).