

## Customer Behavior Analysis using Web Usage Mining

Shahrzad Jalaly<sup>1</sup>, Neda Abdolvand<sup>2</sup>, Saeedeh Rajae Harandi<sup>3</sup>

1- Master's Degree, Information Technology Management, Faculty of Social Science and Economics, Alzahra University, Tehran, Iran

shahrzad.jalaly@gmail.com

2- Assistant Professor and Faculty Member of Management Department, Alzahra University, Tehran, Iran

n.abdolvand@alzahra.ac.ir

3- Master's Degree, Information Technology Management, Faculty of Social Science and Economics, Alzahra University, Tehran, Iran

Saeedeh.rh@gmail.com

### Abstract

Recent marketing strategies consider customers as important sources of the organization. Therefore, acquiring knowledge about customers and understanding their needs is necessary for keeping customers in an e-commerce business. Online customer shopping behavior is difficult to predict because they rarely visit the stores for real shopping, which is a challenge for marketers and researchers. Therefore, online business needs to analyze customers' behavior in order to be successful. Hence, this study aims to provide a framework for increasing the accuracy of the analysis and recognition of customer groups as well as providing the model and rules for predicting customers' behavior. Therefore, CRISP-DM and K-means algorithm were used for clustering data. Then, by assigning three tags of purchase, waiting and not purchase to customers, customers were categorized by C5 decision tree. Finally, a model with a precision of 63.6% and a collection of 261 rules with a confidence of 70% was obtained.

**Keywords:** Customer Behavior Analysis; Data Mining; Web Usage Mining; Online Retailing

### تحلیل رفتار مشتری با استفاده از کاوش کاربری وب

شهرزاد جلالی<sup>۱</sup>، ندا عبدالوند<sup>۲</sup>، سعیده رجائی هرندی<sup>۳</sup>

۱- دانش‌آموخته کارشناسی ارشد مدیریت فناوری اطلاعات، دانشکده علوم اجتماعی و اقتصاد، دانشگاه الزهرا (س)

shahrzad.jalaly@gmail.com

۲- استادیار گروه مدیریت دانشکده علوم اجتماعی و اقتصاد، دانشگاه الزهرا (س)

n.abdolvand@alzahra.ac.ir

۳- دانش‌آموخته کارشناسی ارشد مدیریت فناوری اطلاعات، دانشکده علوم اجتماعی و اقتصاد، دانشگاه الزهرا (س)

saeedeh.rh@gmail.com

### چکیده

در راهبردهای اخیر بازاریابی، مشتریان از منابع مهم سازمان قلمداد می‌شوند. بر اساس این، کسب دانش درباره مشتریان و درک نیازهای آنها برای حفظ مشتریان در تجارت الکترونیک بسیار ضروری است. پیش‌بینی رفتار خرید مشتریان برخط دشوار است؛ زیرا به‌ندرت بازدید آنها از فروشگاه‌ها به خرید واقعی ختم می‌شود و این موضوع برای بازاریابان و پژوهشگران نوعی چالش شده است؛ از این‌رو، برای داشتن کسب‌وکار برخط موفق باید رفتار مشتریان را تحلیل کرد. بنابراین، این پژوهش با دو هدف الف) طرح چارچوبی برای افزایش دقت تحلیل و شناخت گروه‌های مشتریان و ب) ارائه مدل و قوانینی برای پیش‌بینی رفتار آنها، رفتار مشتریان را تحلیل می‌کند. در این پژوهش از روش کریسپ و الگوریتم کا-میانگین برای خوشه‌بندی مشتریان استفاده شده است؛ سپس با اختصاص سه نوع برجسب خرید، خرید نکردن و انتظار خرید به مشتریان و با استفاده از درخت تصمیم C5 مشتریان دسته‌بندی شدند. در نهایت، مدلی با دقت ۶۳.۶٪ و مجموعه‌ای از ۲۶۱ قانون مناسب با اطمینان ۷۰٪ برای کسب‌وکار به دست آمد.

**کلید واژه‌ها:** تحلیل رفتار مشتری، داده کاوی، کاوش کاربری وب، خرده‌فروشی برخط

## ۱- مقدمه و بیان مسئله

در عصر الکترونیک، شرکت‌های تجارت الکترونیک از دنیای قدیمی که در آن، محصولات استاندارد، بازارهای همگن و چرخه توسعه و عمر محصول طولانی یک قانون بود، به دنیای جدیدی انتقال می‌یابند که در آن انواع محصولات استاندارد جایگزین وجود دارند. مصرف‌کنندگان می‌توانند از میان میلیون‌ها کالا در یک فروشگاه برخط به جای ده‌ها هزار کالا در فروشگاه‌های بزرگ، کالاهای ضروری خود را انتخاب کنند (آرورا و چوپرا<sup>۱</sup>، ۲۰۱۶). از یک طرف رشد سریع و انتشار فناوری اینترنت، محبوبیت بازدید از پایگاه‌های وب را در میان بازدیدکنندگان افزایش داده است و از طرفی دیگر، در راهبردهای اخیر بازاریابی، مشتریان جزو منابع مهم یک سازمان قلمداد می‌شود (پارک و چونگ<sup>۲</sup>، ۲۰۰۹). بنابراین کسب دانش درباره مشتریان و درک نیازهای آنها برای حفظ مشتریان در تجارت الکترونیک بسیار ضروری است، زیرا رقبا به اندازه یک کلیک از ما فاصله دارند (ژنگ<sup>۳</sup> و همکاران، ۲۰۰۴). برای خرده‌فروشی برخط که در آن حفظ مشتری از عوامل کلیدی موفقیت است، توسعه مؤثر حضور در وب و عملیات بخش‌های مدیریتی لازم است. یک خرده‌فروشی برخط موفق باید سطحی ایده‌آل از سیستم، اطلاعات و کیفیت خدمات را ارائه کند و مشتریان را برای بازدید دوباره از پایگاه وب خود جلب کند. هنگامی که بازدیدکنندگان، بازدید لذت‌بخشی داشته باشند، احتمال بازدید دوباره آنها از پایگاه وب افزایش می‌یابد (آهن<sup>۴</sup> و همکاران، ۲۰۰۷). پیش‌بینی رفتار خرید مشتریان برخط دشوار است؛ زیرا به‌ندرت بازدید آنها از فروشگاه‌ها به خرید واقعی

می‌انجامد و این برای بازاریابان و پژوهشگران نوعی چالش شده است. میزان تبدیل بازدید به خرید برای کسب و کارهای برخط متوسط حداکثر ۳ درصد است (پارک و چونگ، ۲۰۰۹). در این صورت سازمان‌ها باید به جای هدف قراردادن تمام مشتریان به یک اندازه یا پیشنهاد مشوق‌های یکسان به همه آنها، تنها مشتریانی را هدف قرار دهند که براساس نیازهای فردی یا رفتارهای خریدشان به معیارهای سودبخش مشخصی دست یافته‌اند (لیو و تزنگ<sup>۵</sup>، ۲۰۱۰).

وب‌کاوی گام مؤثری برای رسیدن به این هدف است. وب‌کاوی، راهی تعیین‌کننده برای درک کاربران تجارت الکترونیک و تبدیل اطلاعات به مزیت رقابتی است و سازمان‌ها را قادر می‌سازد تا تصمیمات مبتنی بر داده بگیرند و راهبردهای تصمیم‌گیری خود را بهبود و توسعه دهند. همچنین در به دست آوردن مشتریان جدید و حفظ مشتریان موجود و بهبود رضایت مشتری نیز کمک می‌کند (ژنگ و همکاران، ۲۰۰۴؛ شانتی<sup>۶</sup>، ۲۰۱۷). کاوش کاربری وب شاخه‌ای از وب‌کاوی است که بر کاربرد تکنیک‌های داده‌کاوی برای یافتن الگوهای مفید تمرکز دارد و می‌تواند رفتار کاربر را هنگامی تعامل با وب پیش‌بینی کند (روا و آرورا<sup>۷</sup>، ۲۰۱۷). تجزیه و تحلیل این اطلاعات به سازمان‌ها در تعیین ارزش طول عمر مشتریان، طراحی راهبرد بازاریابی در محصولات و خدمات، ارزیابی اثربخشی کمپین‌های تبلیغاتی، بهینه‌سازی عملکرد برنامه‌های کاربردی مبتنی بر وب، به دست دادن محتوای شخصی‌سازی شده به بازدیدکنندگان و پیدا کردن مؤثرترین ساختار منطقی برای فضای وب خود، کمک می‌کند. این نوع تجزیه و تحلیل شامل شناسایی خودکار الگوها و روابط معنی‌دار از مجموعه بزرگی از

<sup>1</sup> Arora & Chopra

<sup>2</sup> Park & Chung

<sup>3</sup> Zhang

<sup>4</sup> Ahn

<sup>5</sup> Liou & Tzeng

<sup>6</sup> Shanthi

<sup>7</sup> Rao & Arora

این طریق شرکت برای مشتریان هر خوشه سیاست‌های متفاوتی را اعمال می‌کند و مبتنی بر دانش به‌دست‌آمده تصمیم می‌گیرد و محتوا و طراحی سایت را براساس نیازهای هر خوشه تغییر می‌دهد. در دسته‌بندی بررسی می‌شود که براساس ویژگی‌های رفتاری مشتریان، چه رفتارهایی باعث خرید شده و چه رفتارهایی به خرید ختم نمی‌شود. با ایجاد مدلی جامع و مقبول، می‌توان رفتار مشتریان جدید را با داشتن ویژگی‌های رفتاری آنان، در یکی از این دسته‌ها طبقه‌بندی و رفتار آینده آنها را پیش‌بینی کرد. این کار باعث می‌شود شرکت بتواند مشتریان باارزش خود را شناسایی کند و با سپردن مزایا و خدماتی به آنها به‌سوی جذب سود بیشتر حرکت کند.

با توجه به هدف، نخست مبانی نظری پژوهش و سپس مدل و روش پژوهش بررسی می‌شود و درنهایت پژوهش با بحث و نتیجه‌گیری و ارائه پیشنهادها خاتمه می‌یابد.

## ۲- مبانی نظری پژوهش

خرده‌فروشی برخط<sup>۲</sup>، یک کسب‌وکار مبتنی بر اینترنت است که محصولات و خدمات را در وب ارائه می‌دهد. خرده‌فروشی برخط نه تنها یک سیستم اطلاعاتی است، ارائه‌ای کامل از یک فروشگاه به مشتری است (آهن و همکاران، ۲۰۰۷). با توجه به رشد روزافزون خرده‌فروشی برخط، خرده‌فروشان برخط به درک دلایل خاصی نیاز دارند که چرا مصرف‌کنندگان خرید برخط را انتخاب می‌کنند (دکا<sup>۳</sup>، ۲۰۱۷). رفتار مشتری، حاصل تعامل پیچیده میان تعدادی از عوامل است که این عوامل شامل سطح فعالیت بازاریابی، رقابت‌پذیری محیط، درک نام تجاری، تأثیر فناوری‌های جدید و

داده‌های نیمه‌ساخت یافته است که بیشتر در پایگاه وب و برنامه‌های کاربردی لاگ سرور و همچنین در منابع داده عملیاتی مرتبط ذخیره می‌شوند (لیو<sup>۱</sup> و همکاران، ۲۰۱۱). برای داشتن راه‌حل تجارت الکترونیک موفق، نیاز است که رفتارهای کلیک مشتری در پایگاه وب جمع‌آوری و بررسی شود؛ زیرا جذب مشتریان جدید و حفظ مشتریان باارزش از اهمیت بسیاری برخوردار است (ژنگ و همکاران، ۲۰۰۴). بنابراین این سؤالات مطرح می‌شود که چگونه می‌توان با ترکیب روش‌های داده‌کاوی و کاوش کاربری وب به درک و شناخت بهتری از گروه‌های مشتریان دست یافت و چگونه می‌توان با ترکیب روش‌های داده‌کاوی و کاوش کاربری وب و براساس ویژگی‌های رفتاری مشتریان پیش‌بینی کرد که آیا رفتار آنها به خرید ختم می‌شود یا خیر؟ از این‌رو، این پژوهش، دو هدف اصلی را دنبال می‌کند.

نخستین هدف، طرح چارچوبی برای افزایش دقت تحلیل و شناخت گروه‌های مشتریان و دیگری، ارائه مدل و قوانینی برای پیش‌بینی رفتار مشتریان بر مبنای ویژگی‌های رفتاری آنها با استفاده از ترکیب رویکردهای وب‌کاوی و داده‌کاوی است. در این پژوهش از رویکرد ترکیبی کمتر به‌کاررفته داده‌کاوی و کاوش کاربری وب استفاده شده است که به مزیت رقابتی برای شرکت می‌انجامد.

درواقع نوآوری پژوهش در اعمال روش‌های داده‌کاوی بر ابعاد مختلف داده‌های مشتری شامل داده‌های رفتاری مرور وب و داده‌های رفتار خرید است که به بهبود دانش از مشتری خواهد انجامید. در این پژوهش، داده‌ها و رفتار کاربران درون سایت با استفاده از دو تکنیک خوشه‌بندی و دسته‌بندی تحلیل می‌شود. خوشه‌بندی مشتریان سبب می‌شود شرکت بتواند مشتریان خود و مشتریان باارزش را شناسایی کند. از

<sup>2</sup> Online Retailing

<sup>3</sup> Deka

<sup>1</sup> Liu

چالش برانگیز برای داده کاوی است. در واقع وب کاوی، تحلیل رفتار الکترونیکی مشتری است و به نوعی استفاده از تکنیک‌های داده کاوی برای کشف و استخراج اطلاعات از مستندات و خدمات وب تعریف می‌شود (شیخ و مناریا<sup>۴</sup>، ۲۰۱۷). کاربرد تکنیک‌های داده کاوی روی داده‌های کاربردی وب اصطلاح علمی جدیدی با نام کاوش کاربردی وب ایجاد کرده است (دارمارانجان و دورایزنگاسوامی<sup>۵</sup>، ۲۰۱۶). با کاوش کاربردی وب، اطلاعات مفیدی مانند الگوهای پیمایش کاربران با استفاده از داده‌های لاگ وب استخراج و تحلیل می‌شود. هم طراحان و هم کاربران وب از کاوش کاربردی وب سود می‌برند. از طرفی، با تحلیل الگوهای پیمایش در لاگ وب سرور، طراحان وب سایت رفتارهای بازدید کاربران وب را تعیین می‌کنند؛ بنابراین آنها درمی‌یابند معروف‌ترین صفحات روی سایت کدامند و کدام صفحات با احتمال بیشتری با هم بازدید می‌شوند. از طرفی دیگر، کاربران وب همچنین می‌توانند از این موضوع برای دسترسی مؤثرتر به وب استفاده کنند (سان و ژنگ<sup>۶</sup>، ۲۰۰۴).

## ۲-۱- پیشینه پژوهش

در بسیاری از پژوهش‌های انجام شده درباره رفتار مشتری، مقادیر دموگرافیک مشتریان برای تحلیل رفتار آنها استفاده می‌شود (لیو و تزنگ<sup>۱</sup>، ۲۰۱۰)؛ برای نمونه، ژنگ و همکاران (۲۰۰۴) با استفاده از داده‌های جامع که شامل لاگ وب و اطلاعات مشتریان در وب سایت‌های تجارت الکترونیک هستند، الگوهای رفتاری کاربران و درک رفتارهای خرید آنها در پایگاه وب را شناسایی کرده‌اند. در این مقاله از سه تکنیک دسته‌بندی، خوشه‌بندی و قوانین انجمنی برای بررسی روش پیشنهادی استفاده کرده‌اند. داده‌های دموگرافیک

نیازهای فردی است (لیو و تزنگ<sup>۱</sup>، ۲۰۱۰). به منظور افزایش رضایت مشتری و جلوگیری از ترک سازمان از سوی مشتری، سازمان باید بر بخش‌بندی و تأمین نیازهای فردی مشتریان متمرکز شود (تسای<sup>۱</sup> و همکاران، ۲۰۱۵). به طور کلی «بخش‌بندی مشتری» فرایند تقسیم مشتریان سازمان به گروه‌های مختلف بر مبنای اطلاعات مختلف جغرافیایی، جمعیت‌شناختی، رفتار شناختی و اتخاذ راهبردهای مناسب هر گروه با توجه به مصرف کالا و خدمات و تاریخچه خرید مشتریان است (تسای و همکاران، ۲۰۱۵). بخش‌بندی مشتری با تکنیک‌ها و روش‌های تحلیل داده‌ای متفاوت انجام می‌شود که در این میان استفاده از تکنیک‌های داده کاوی روی داده‌هایی که از تراکشن‌های برخط تولید می‌شوند، متداول‌تر از سایر روش‌هاست. داده کاوی قابلیت پیچیده جست‌وجوی داده است که از الگوریتم‌های آماری برای کشف الگوها و همبستگی داده‌ها استفاده می‌کند. به بیان ساده، داده کاوی فرایند خودکار یافتن اطلاعات مفید از مخازن بزرگ داده است (بهشتیان اردکانی<sup>۲</sup> و همکاران، ۲۰۱۸). یکی از متداول‌ترین روش‌ها برای مرتب‌سازی و تحلیل مشتریان، امتیازدهی به آنها بر اساس دفعات خرید و مقدار پرداخت است که از شناخته شده‌ترین روش‌ها برای انجام این کار مدل RFM (تازگی خرید، تکرار خرید و ارزش پولی خریده‌ها) است که اساس بخش‌بندی برای بازاریابی مستقیم است. هنگامی که امتیازات RFM مشتریان تعیین می‌شود، می‌توان مشتریان را به صورت بخش‌هایی گروه‌بندی و متعاقباً سودآوری آنها را تحلیل کرد (مقدم<sup>۳</sup> و همکاران، ۲۰۱۷).

وب بزرگ‌ترین پایگاه داده در دسترس و موضوعی

<sup>4</sup> Sheikh & Menaria

<sup>5</sup> Dharmarajan & Dorairangaswamy

<sup>6</sup> Sun & Zhang

<sup>1</sup> Tsai

<sup>2</sup> Beheshtian-Ardakani

<sup>3</sup> Moghadam

پارک و چونگ (۲۰۰۹) استفاده از داده‌های جریان‌های کلیک برای پیش‌بینی رفتارهای خرید کاربران را ارائه کردند. آنها با استفاده از تحلیل رگرسیون سلسله‌مراتبی ثابت کردند که کاربران انتقال داده‌شده از پایگاه وب ارجاع‌دهنده، کمتر از کاربرانی که مستقیم وارد سایت می‌شوند، خرید می‌کنند و هرچه مدت زمان ماندن این کاربران در سایت بیشتر باشد و صفحات مشاهده کرده کمتر باشد، احتمال خریدشان بیشتر است. علاوه بر این، استفاده از تکنیک‌های داده‌کاوی شامل آمار توصیفی و قوانین انجمنی برای تحلیل رفتار هدایتی کاربران و شناسایی الگوهای گشت و گذار آنها را سیسودیا و ورما<sup>۵</sup> (۲۰۱۲) مطرح کردند. هونگ<sup>۶</sup> و همکاران (۲۰۱۳) در پژوهش مشابهی از روش کاوش کاربری وب روی سرویس مراقبت از خود برای افراد سالمند، به منظور بهبود درک و تحلیل رفتارهای آنها استفاده کرده‌اند. آنها این کار را با تکنیک‌های تحلیل انجمنی و مدل مارکوف همراه با الگوریتم بهبود یافته کامپانگین انجام داده‌اند. با استفاده از نمونه‌ای از لاگ وب سرور مرکز فضایی ناسا و کاوش کاربری وب، پامونتا<sup>۷</sup> و همکاران (۲۰۱۲) به اطلاعات آماری از نشست کاربر دست یافتند که می‌توان از آن برای شناسایی الگوهای دسترسی کاربران و تحلیل رفتار آنها استفاده کرد. جنامانی و همکاران (۲۰۰۳) نیز مدل فرایندی نیمه مارکوف (یک ابزار کاوش کاربری وب)، برای درک رفتار مشتریان الکترونیک را مطرح کردند که نتایج این مدل به بهبود طراحی سایت و تشخیص عملکرد آن کمک کرد. ها<sup>۸</sup> (۲۰۰۲) سیستم شخصی سازی شده مبتنی بر وبی را پیشنهاد داد که از کاوش کاربری وب

و خرید کاربران و همچنین اطلاعات بازدید آنها از سایت شامل تعداد صفحات بازدید شده در هر نشست و همچنین تعداد صفحات بازدید شده در بخش‌های مختلف سایت نیز برای خوشه‌بندی مشتریان استفاده شده است.

همچنین در تحلیل رفتار مشتریان، اهمیت متغیرهای ارزشمند رفتار مشتری به طور گسترده مطالعه شده است (لیو و تزنگ، ۲۰۱۰). پژوهشگران مشاهده کرده‌اند که متغیرهای RFM نه تنها برای تحلیل رفتار مشتریان مفید هستند، می‌توانند به صورت مؤثر در تجسس ارزش مشتری و بازارهای گوشه<sup>۱</sup> استفاده شوند (لیو و تزنگ، ۲۰۱۰). هسیه<sup>۲</sup> (۲۰۰۴) روشی پیشنهاد می‌دهد که داده‌کاوی و مدل‌های امتیازدهی رفتار یعنی RFM را برای مدیریت مشتریان یک بانک یکپارچه می‌کند. او توانست مشتریان بانک را به سه گروه بزرگ سودمند دسته‌بندی کند.

به تازگی کاوش کاربری وب، توجه پژوهشگران و متخصصان تجارت الکترونیک را در تحلیل رفتار مشتری به خود جلب کرده است. پژوهش‌ها در زمینه کاوش کاربری وب بیشتر بر توسعه تکنیک‌های کشف دانش، به خصوص آنهایی که برای تجزیه و تحلیل داده‌های کاربری وب طراحی شده، متمرکز شده است. بیشتر این تلاش‌ها، بیشتر بر سه پارادایم اصلی قوانین وابستگی، الگوهای ترتیبی و خوشه‌بندی توجه دارند (فاکا و لنزی<sup>۳</sup>، ۲۰۰۵)؛ برای مثال، در پژوهش ینگ و سو<sup>۴</sup> (۲۰۱۲) از الگوریتم SVM برای خوشه‌بندی رفتار مشتریان شبکه به منظور ارائه خدمات بهتر به آنها استفاده شده است که سازمان‌ها با استفاده از نتایج خدمات خود را بیشتر شخصی‌سازی می‌کنند. همچنین

<sup>5</sup> Sisodia & Verma

<sup>6</sup> Hung

<sup>7</sup> Pamutha

<sup>8</sup> Ha

<sup>1</sup> Niche Markets

<sup>2</sup> Hsieh

<sup>3</sup> Facca & Lanzi

<sup>4</sup> Yang & Su

آنها چارچوبی را ارائه می‌دهد. در واقع، از رویکرد ترکیبی و مغفول‌مانده داده کاوی و کاوش کاربری وب استفاده شده است که می‌تواند به مزیت رقابتی برای شرکت منجر شود. در این پژوهش، داده‌ها و رفتار کاربران سایت با استفاده از دو تکنیک خوشه‌بندی و دسته‌بندی تحلیل می‌شوند که به شناخت بهتر مشتریان با ارزش منجر شده و براساس ویژگی‌های رفتاری مشتریان، مشخص می‌شود چه رفتارهایی باعث خرید شده و چه رفتارهایی نمی‌شود؛ شرکت با تحلیل این موارد می‌تواند با اعطای مزایا و خدماتی به مشتریان، به سوی جذب سود بیشتر حرکت کند. مزیت اصلی این روش کسب نتایج بهتر از فرایند تحلیل با استفاده از اطلاعات بیشتر مشتریان است.

### ۳- روش پژوهش

این پژوهش دارای رویکردی سازنده است و با دو هدف الف) ارائه چارچوبی برای افزایش دقت تحلیل و شناخت گروه‌های مشتریان و ب) ارائه مدل و قوانینی برای پیش‌بینی رفتار آنها رفتار مشتریان را تحلیل می‌کند. داده‌ها و رفتار کاربران درون سایت با استفاده از دو تکنیک خوشه‌بندی و دسته‌بندی تحلیل می‌شود. این پژوهش بر مبنای متدولوژی کریسپ<sup>۴</sup> انجام شده که یکی از قوی‌ترین روش‌های تحلیلی برای اجرای پروژه‌های داده کاوی است و شامل شش مرحله فهم کسب و کار، فهم داده، آماده‌سازی، مدل‌سازی، ارزیابی و به‌کارگیری است (بهشتیان اردکانی و همکاران، ۲۰۱۸). ترتیب و توالی این شش مرحله انعطاف‌پذیر است. روش کریسپ بسیار کامل و مستند است. تمامی مراحل آن، به‌موقع سازماندهی، ساخته و تعریف می‌شود که اجازه می‌دهد یک پروژه بتواند به راحتی

برای دادن پیشنهادهایی شخصی به مشتریان برخط استفاده می‌کند. این سیستم برای ارائه اطلاعات به مشتریان طراحی شده و به آنها در خرید کالاها کمک می‌کند. بایی<sup>۱</sup> و همکاران (۲۰۰۳) سیستم انتخاب آگهی وبی طراحی کردند که کاربران پایگاه وب با ترجیحات مشابه را از راه کاوش کاربری وب به چند بخش تقسیم می‌کند. این سیستم با استنتاج فازی، تبلیغات مناسب را پیشنهاد می‌دهد. کیم<sup>۲</sup> و همکاران (۲۰۰۳) روشی برای پیش‌بینی رفتار خرید مشتریان تجارت الکترونیک با ترکیب چند دسته‌بندی بر مبنای الگوریتم ژنتیک پیشنهاد داد که عملکرد بهتری نسبت به دسته‌بندی‌های تکی داشت. جنامانی<sup>۳</sup> و همکاران (۲۰۰۳) با پیشنهاد مدل فرایندی نیمه‌مارکوف (یک ابزار کاوش کاربری وب)، به درک رفتار مشتریان الکترونیک پرداختند که نتایج این مدل به بهبود طراحی سایت و تشخیص عملکرد آن کمک می‌کند.

با توجه به پژوهش‌ها انجام شده در این زمینه، بیشتر پژوهش‌ها به پیش‌بینی رفتار مشتری با استفاده از رویکردهای وب کاوی یا داده کاوی تکیه کرده و پژوهش‌های اندکی به پیش‌بینی رفتار مشتریان با ترکیب این دو رویکرد توجه کرده‌اند و پژوهشی در این زمینه در ایران انجام نشده است. اگر اطلاعات بیشتر در دسترس باشد که با داده‌های لاگ وب ارتباط داشته و همچنین به مشکل هم مربوط باشند، می‌توانند به صورت چشمگیری نتایج را بهبود بخشند و به نتایج دقیق‌تر هم برسند. داده‌های مشتریان، اطلاعات بیشتری هستند که همراه با داده‌های لاگ استفاده می‌شوند. از این‌رو، در این پژوهش به منظور افزایش دقت تحلیل و شناخت گروه‌های مشتریان و ارائه مدل و قوانینی برای پیش‌بینی رفتار مشتریان بر مبنای ویژگی‌های رفتاری

<sup>1</sup> Bae

<sup>2</sup> Kim

<sup>3</sup> Jenamani

<sup>4</sup> CRISP-DM

دسته‌بندی مشتریان براساس ویژگی‌های رفتاری آنها و ایجاد مدلی برای پیش‌بینی نوع رفتار مشتریان آینده استفاده می‌شود. درخت تصمیم در بین الگوریتم‌های طبقه‌بندی، روش قدرتمندی است که محبوبیت آن با رشد داده‌کاوی به‌طور فزاینده‌ای در حال افزایش است. درخت تصمیم روشی برای نمایش دسته‌ای از قوانین است که به یک رده یا مقدار منتهی می‌شود. الگوریتم درخت تصمیم C5 روشی افزایشی از هرس کردن درخت را به کار می‌گیرد تا خطای طبقه‌بندی ناشی از نویز یا جزئیات خیلی زیاد را در داده‌های آموزشی کاهش دهد. این درخت می‌تواند دسته‌بندی‌هایی هم برای درخت تصمیم و هم برای مجموعه قوانین تولید کند. در پایان با استفاده از معیارهای ارزیابی الگوریتم‌های خوشه‌بندی و دسته‌بندی نتایج این دو تکنیک ارزیابی می‌شود.

#### ۴- پیاده‌سازی

این پژوهش بر مبنای فرایند کریسپ انجام شده و شامل مراحل زیر است:

##### فاز ۱. درک کسب و کار

این فاز بر درک اهداف و نیازمندی‌های پروژه از دیدگاه کسب و کار تمرکز می‌کند. این پژوهش در شرکت XYZ انجام شده است که مرجع تخصصی نقد و بررسی و فروش اینترنتی کالا در ایران است. این شرکت در زمینه فروش برخط محصولات مختلف فعالیت می‌کند و در حال حاضر در حدود ۶۷ هزار کالای مختلف را در سایت خود به فروش می‌رساند. سایت این شرکت به‌طور میانگین روزانه بیش از ۴۰۰ هزار بازدیدکننده دارد و از پربازدیدترین سایت‌های کشور است. کاربران و مشتریان XYZ می‌توانند با حق انتخاب متنوع و با دریافت اطلاعاتی کامل برای انتخاب درست کالای مدنظر خود، وب‌سایت این شرکت را

درک یا تجدیدنظر شود (آزودو<sup>۱</sup>، ۲۰۰۸). برای بخش‌بندی و ارزش‌گذاری و همچنین الگویی رفتاری خرید مشتری از الگوریتم‌ها و تکنیک‌های مختلفی استفاده می‌شود که در دو گروه کلی خوشه‌بندی و کشف قواعد انجمنی دسته‌بندی می‌شوند. بر همین اساس، در این پژوهش سعی بر آن است تا از تکنیک خوشه‌بندی و الگوریتم کا-میانگین استفاده شود. الگوریتم خوشه‌بندی کا-میانگین یکی از ساده‌ترین و البته مشهورترین الگوریتم‌های یادگیری بدون نظارت است. این الگوریتم، کاربردی‌ترین روش خوشه‌بندی داده‌هاست که از مزایایی همچون سرعت و ساده‌بودن در اجرا برخوردار است و در مسائل بزرگ بسیار کارایی دارد. این الگوریتم، روشی بسیار مناسب برای یافتن گروه‌های مشتریان با رفتارهای مشابه است و در تلاش است تا عدم تشابه میان گروه‌های مشتریان را به حداقل برساند (تاسینگر و هوبر<sup>۲</sup>، ۲۰۰۰)؛ سپس به منظور سنجش صحت<sup>۳</sup> نتایج، از معیار ارزیابی دیویس بولدین با استفاده از نرم‌افزار رپیدمایز روی ۷۴.۳۴۴ رکورد نرمال شده انجام می‌شود (اطلاعات مربوط به سازمان و داده‌ها در بخش بعد آورده شده است). برای بررسی تمایز خوشه‌ها آزمون آنووا با استفاده از نرم‌افزار اسپاس‌اس ۱۹ روی خوشه‌های به‌دست‌آمده اجرا می‌شود. در مرحله بعد با توجه به ضرورت استخراج قوانین، به دلیل اینکه شبکه عصبی مانند یک جعبه سیاه عمل می‌کند، با وجود دقت بیشتر دسته‌بندی شبکه عصبی، از الگوریتم درخت تصمیم C5، با اختصاص سه برچسب به نام‌های خرید، انتظار خرید و خریدنکردن با استفاده از نرم‌افزار کلمنتاین برای بررسی رفتارها با توجه به ویژگی‌های رفتاری کاربران و

<sup>1</sup> Azevedo

<sup>2</sup> Theusinger & Huber

<sup>3</sup> Goodness

بررسی و با حداکثر اطمینان کالای خود را انتخاب و خرید کنند. XYZ با ارائه طیف گسترده‌ای از معتبرترین برندها در گروه‌های مختلف و با همکاری نزدیک با واردکنندگان و توزیع‌کنندگان اصلی این کالاها در ایران، تلاش می‌کند نیازهای متفاوت مشتریان با کاربری‌های متفاوت آنان را برآورده سازد. کیفیت و سهولت استفاده از پایگاه وب و خدمات ارائه‌شده در آن، همواره یکی از مهم‌ترین و بااولویت‌ترین موضوعات در XYZ بوده است و این شرکت تلاش می‌کند در پایگاه وب XYZ، خدماتی شایسته و تجربه‌ای خوشایند را برای مخاطبان خود به ارمغان بیاورند.

### فاز ۲ و ۳. شناخت و آماده‌سازی داده

این فاز با جمع‌آوری داده‌های اولیه و تصمیم اینکه کدام داده‌ها و در چه فرمت و اندازه‌ای لازم خواهند بود، آغاز می‌شود و فعالیت‌ها را به‌منظور آشنایی با داده‌ها، شناسایی مشکلات کیفیت داده‌ها، دریافت بینش‌های مقدماتی درباره داده‌ها یا کشف زیرمجموعه‌های جالب برای شکل‌دادن فرضیه‌ها برای اطلاعات پنهان ادامه می‌دهد. براساس رویکرد پیشنهادی، دو سری داده (پایگاه داده لاگ و پایگاه داده مشتریان) ضروری است. بنابراین، داده‌های ۱۰۰.۰۰۰ کاربر تصادفی و یکتا در بازه شش‌ماه دوم سال ۱۳۹۳ از پایگاه داده لاگ جمع‌آوری شده و داده‌های مربوط به سبد این کاربران هم از پایگاه داده مشتریان گرفته شده است. داده‌های دموگرافیک این کاربران به دلیل وجود برخی مسائل امنیتی و فنی در دسترس نبودند؛ از این‌رو، از استفاده از آنها صرف‌نظر می‌شود. تعداد نشست‌های این کاربران در حدود ۱.۲ میلیون و داده‌های لاگ در حدود ۷.۵ میلیون رکورد و داده‌های مربوط به سبد خرید این کاربران در حدود ۱۴۲ هزار رکورد بود. در پایگاه داده مشتریان کلیه

اطلاعات مربوط به پروفایل مشتریان ذخیره می‌شوند؛ اما در این پژوهش، داده‌های سبد خرید مشتریان شامل شناسه کاربر، شماره سبد، تاریخ و زمان، مبلغ، شناسه کالای موجود در سبد خرید، نام کالای موجود در سبد خرید و دسته هر کالای موجود در سبد خرید و داده‌های لاگ سرور شامل تعداد نشست‌ها، زمان کل، صفحات مختلف، صفحات کالا، صفحات جست‌وجو، تأخر<sup>۱</sup>، مدت مراجعه، مدت زمان بین نشست‌ها، کالاهای سبد، ارزش مالی، دسته‌های کالاهای سبد و برچسب دسته است. برای ایجاد مسیرهای طی‌شده کاربر در هر نشست و همچنین تکمیل مسیرهای آنها برنامه‌ای به زبان برنامه‌نویسی پایتون<sup>۲</sup> نوشته شد و برای هر کاربر مسیرهای هر نشست و زمان صرف‌شده در هر نشست آن به دست آمد. پس از به دست آوردن مسیر دسترسی هر کاربر در هر نشست، با ترکیب ویژگی‌های رفتاری استخراج‌شده از سبد کالا و لاگ سرور بر اساس مدل RFM، ویژگی‌های مشترک آنها شامل ۱۱ ویژگی تعداد نشست‌ها، زمان کل، صفحات مختلف، صفحات کالا، صفحات جست‌وجو، تأخر، مدت زمان بین نشست‌ها، کالاهای سبد، ارزش مالی، دسته‌های کالاهای سبد، برچسب دسته از پایگاه داده لاگ و مشتریان استخراج شد.

وظیفه اصلی در کاوش کاربری وب، پیش‌پردازش داده است که شامل پاک‌سازی داده، شناسایی کاربر، شناسایی نشست، کامل کردن مسیر و شکل‌دهی داده است (پامونتا و همکاران، ۲۰۱۲). سه مرحله نخست از مراحل پیش‌پردازش داده، در هنگام ذخیره‌سازی لاگ‌ها در جداول پایگاه داده انجام شد. تنها یک مرحله پاک‌سازی دیگر روی داده‌های لاگ انجام شد. به این صورت که درخواست‌های تکراری پشت سر هم در هر نشست و نشست‌هایی که تنها یک درخواست

<sup>۱</sup> . Recency

<sup>۲</sup> . Python



پیش‌بینی نوع رفتار مشتریان آینده از تکنیک درخت تصمیم استفاده می‌شود.

### • شناخت گروه‌های مشتریان

اولین هدف این پژوهش، شناخت الگوها و گروه‌های مختلف مشتریان بر مبنای رفتار آنهاست. بدین منظور از الگوریتم خوشه‌بندی کا-میانگین استفاده می‌شود که باید تعداد خوشه‌های آن از پیش تعیین شوند و در ادامه بر مبنای تعداد خوشه‌های اولیه داده‌ها در خوشه‌های مختلف قرار می‌گیرند (مقدم و همکاران، ۲۰۱۷). یکی از معایب الگوریتم کا-میانگین این است که تعداد خوشه‌ها را همچون پارامتر ورودی الگوریتم می‌گیرد و نمی‌تواند تعداد خوشه‌های بهینه را بیابد. بنابراین، با استفاده از شاخص دیویس بولدین الگوریتم برای تعداد خوشه‌های مختلف امتحان شده و بهترین نتیجه برای تعداد خوشه‌ها انتخاب می‌شود. براساس سه معیار شباهت مختلف (فاصله اقلیدسی، فاصله چیشف و فاصله منهن) و برای تعداد خوشه‌های مختلف، الگوریتم کا-میانگین اجرا می‌شود تا زمانی که معیار دیویس بولدین آن به بالای عدد یک برسد. نتایج اعتبارسنجی الگوریتم کا-میانگین با استفاده از شاخص دیویس بولدین در جدول ۱ نشان داده شده است. با توجه به نتایج به دست آمده بهترین جواب برای شاخص دیویس بولدین مربوط به معیار فاصله اقلیدسی و با چهار خوشه است. در واقع هرچه مقدار این شاخص کمتر باشد، خوشه‌ها در بیشترین فاصله از هم قرار دارند.

جدول ۱. نتایج اعتبارسنجی الگوریتم کا-میانگین با استفاده از شاخص دیویس بولدین

تعداد خوشه	فاصله چیشف	فاصله اقلیدسی	فاصله منهن
۳	۰.۸۰۸	۰.۸۱۹	۰.۸۲۹
۴	۰.۷۲۵	۰.۷۱۶	۰.۷۱۸

گرفته‌اند. نقاط مرکزی هر چهار خوشه که در واقع میانگین فواصل نقاط موجود در آن خوشه به ازای تمامی ویژگی‌های به کاررفته، در جدول ۲ آورده شده است.

برای آنها ثبت شده بود، از داده‌ها حذف شد. همچنین تمامی داده‌ها از لحاظ داشتن مقادیر ازدست‌رفته و داده‌های غیرطبیعی بررسی شدند. پس از پاک‌سازی، تعداد کاربران به ۷۴.۳۴۴ رسید. پس از استخراج ویژگی‌ها و حذف داده‌های پرت نوبت به نرمال‌سازی آنها می‌رسد. به دلیل وجود محدوده‌ها و مقیاس‌های مختلف اندازه‌گیری برای هر ویژگی، لازم است برای یکسان کردن این محدوده‌ها، تمامی ویژگی‌ها نرمال شوند. در این پژوهش، برای نرمال‌سازی داده‌ها از روش نرمال‌سازی مینیم-ماکزیم استفاده شده است. در این روش نرمال‌سازی، تبدیل خطی روی داده‌های اصلی انجام می‌دهد. فرض کنید که  $\max_A$  و  $\min_A$  کمترین و بیشترین مقدار ویژگی  $A$  باشند؛ سپس نرمال‌سازی مینیم-ماکزیم، مقداری مانند  $v$  از ویژگی  $A$  را به  $v$  در محدوده  $[\text{newmin}_A, \text{newmax}_A]$  می‌نگارد. رابطه (۱) چگونگی این نگاشت را نشان می‌دهد:

رابطه (۱)

### فاز ۴. مدل‌سازی

در

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{newmax}_A - \text{newmin}_A) + \text{newmin}_A$$

مدل‌سا

زی، برای شناخت گروه‌های مختلف مشتریان از تکنیک خوشه‌بندی و به منظور بررسی رفتارها با توجه به ویژگی‌های رفتاری کاربران و ایجاد مدلی برای

سپس خوشه‌بندی با پارامترهای چهار خوشه و معیار شباهت برابر با فاصله اقلیدسی انجام می‌شود (جدول ۲). با توجه به جدول ۲، بیشترین تعداد مشتریان در خوشه اول و کمترین تعداد مشتریان در خوشه دوم قرار

## جدول ۲. نتایج خوشه‌بندی با چهار خوشه و نقاط مرکزی

نقاط مرکزی											تعداد خوشه	تعداد
دسته کالاهای سید	ارزش مالی	کالاهای سید	زمان بین نشست‌ها	مدت مراجعه	تأخر	صفحات جستجو	صفحات کالا	صفحات مختلف	زمان کل	تعداد نشست		
۰.۰۳	۰.۰۱	۰.۰۱	۱	۰	۰.۴۶	۰.۰۲	۰.۰۲	۰.۰۲	۰.۰۵	۰	۲۶۴۴۱	۱
۰.۰۲	۰.۰۱	۰.۰۰	۰.۰۸	۰.۶۴	۰.۱۳	۰.۰۲۰	۰.۰۲	۰.۰۲	۰.۰۵	۰.۰۴	۱۲۲۸۰	۲
۰.۰۳	۰.۰۱	۰.۰۱	۰.۰۳	۰.۱۱	۰.۶۸	۰.۰۲	۰.۰۲	۰.۰۲	۰.۰۵	۰.۰۱	۱۷۷۳۶	۳
۰.۰۳	۰.۰۱	۰.۰۱	۰.۰۳	۰.۱۲	۰.۱۵	۰.۰۲	۰.۰۲	۰.۰۲	۰.۰۵	۰.۰۱	۱۷۸۸۷	۴
-	-	-	-	-	-	-	-	-	-	-	۷۴۳۴۴	کل

به سه گروه آموزشی، آزمایشی و اعتبارسنجی تقسیم می‌شوند. در گره C5، برای افزایش قابلیت اطمینان اعتبارسنجی مدل، گزینه اعتبارسنجی متقابل با ده تکرار انتخاب شده است. این روش بر پایه تقسیم مجموعه داده به ۱۰ قسمت مساوی است که در آن ۹ قسمت از مجموعه داده آموزش مدل و بقیه آزمایش مدل را انجام می‌دهند. همچنین ساخت درخت تصمیم با استفاده از حالت ساده و مطلوب تعمیم انجام شده است.

این حالت باعث می‌شود بیشتر پارامترهای درخت به صورت خودکار تنظیم شود. به دلیل کم بودن تعداد نمونه‌های خرید مشتریان و ارزشمندی این دسته در پژوهش، در ماتریس هزینه درخت، هزینه تشخیص نادرست مشتریان واقعی با عنوان دسته‌های دیگر افزایش یافته است. پس از اجرای الگوریتم، پایین‌ترین سطح شکست‌ها دوباره بررسی می‌شوند و آنهایی که کمک شایانی به ارزش مدل نمی‌کنند، از درخت حذف یا به عبارتی دیگر هرس می‌شوند. پس از اجرای این روند، درختی با میانگین دقت ۶۳.۶٪ و خطای استاندارد ۰.۳ به دست آمد. در نهایت تعداد ۲۶۱ قانون با اطمینان ۷۰٪ حاصل شد.

## فاز ۵. ارزیابی

## • ارزیابی خوشه‌ها

پس از مشخص شدن نقاط مرکزی خوشه‌ها، به منظور بررسی این که به ازای تمامی ویژگی‌ها آیا

## • دسته‌بندی مشتریان براساس ویژگی‌های

## رفتاری

برای دسته‌بندی رفتار مشتریان شرکت XYZ مشکل توزیع نامتوازن کلاس‌ها وجود دارد. به بیان دیگر، دسته خرید در میان دو دسته دیگر، کلاس اقلیت است و نمونه‌های بسیار کمتری در مقایسه با دو دسته دیگر دارد؛ پس باید یک دسته‌بند ایجاد شود که هزینه کلی دسته‌بندی نادرست را کمینه کند. برای مجموعه داده‌های بزرگ با بیش از ۱۰۰۰۰ نمونه، الگوریتم یادگیری حساس به هزینه بهترین نتایج را نسبت به روش‌های نمونه‌برداری ایجاد می‌کند (ویس و همکاران، ۲۰۰۷). بر این اساس، در این پژوهش نیز با داشتن مجموعه داده بزرگ از روش یادگیری حساس به هزینه استفاده می‌شود. بهترین نسبت هزینه در این پژوهش هم به صورت تجربی و براساس مجموعه اعتبارسنجی تعیین شده است. هدف اصلی این پژوهش از دسته‌بندی، علاوه بر مدل‌سازی برای پیش‌بینی رفتار مشتریان، استخراج قوانین به صورت واضح و دقیق به منظور استفاده در کسب و کار است. انتخاب الگوریتم درخت تصمیم C5.0 هم بر این مبنا صورت گرفته است؛ زیرا قوانین درخت تصمیم بسیار ساده و قابل تفسیر است. در این الگوریتم فیله‌های ورودی ویژگی‌های به دست آمده و فیله خروجی برچسب‌های کاربران است. سپس داده‌ها با استفاده از گره پارتیشن

به دست آمده مشخص است، با بررسی ستون سطح معناداری مشاهده می شود که این مقدار برای تمامی متغیرها برابر با ۰.۰۰۰ و کمتر از ۰.۰۵ است؛ از این رو فرض همگن بودن میانگین های جامعه رد می شود و نشان می دهد گروه ها میانگین های مختلفی دارند.

میانگین های به دست آمده، اختلاف معناداری با هم دارند یا خیر، آزمون آنووا در سطح معناداری کمتر از ۰.۰۵ اجرا می شود. برای این کار، چهار میانگین به دست آمده به ازای هر ویژگی با اجرای این آزمون با هم مقایسه می شوند. همان طور که از جدول ۳ و نتایج

جدول ۳. نتایج آزمون آنووا برای الگوریتم های کا-میانگین

متغیر	منبع تغییر	مجموع مجذورات	درجه آزادی	میانگین مجذورات	آزمون F	معناداری
تعداد نشست ها	بین گروهی	۱۲.۸۸۹	۳	۴.۲۹۶	۶۰۰۵.۲۵۳	۰.۰۰۰
	درون گروهی	۵۳.۱۸۴	۷۴۳۴۰	۰.۰۰۱		
	مجموع	۶۶.۰۷۳	۷۴۳۴۳	-		
زمان کل	بین گروهی	۰.۴۲۰	۳	۰.۱۴۰	۶۵.۳۶۹	۰.۰۰۰
	درون گروهی	۱۵۹.۰۲۸۱	۷۴۳۴۰	۰.۰۰۲		
	مجموع	۱۵۹.۷۰۱	۷۴۳۴۳	-		
صفحات مختلف	بین گروهی	۰.۳۱۲	۳	۰.۱۰۴	۱۴۹.۵۹۴	۰.۰۰۰
	درون گروهی	۵۱.۶۳۰	۷۴۳۴۰	۰.۰۰۱		
	مجموع	۵۱.۹۴۲	۷۴۳۴۳	-		
صفحات کالا	بین گروهی	۰.۱۹۷	۳	۰.۰۶۶	۱۳۲.۴۱۵	۰.۰۰۰
	درون گروهی	۳۶.۹۲۸	۷۴۳۴۰	۰.۰۰۰		
	مجموع	۳۷.۱۲۶	۷۴۳۴۳	-		
صفحات جستجو	بین گروهی	۰.۰۲۲	۳	۰.۰۰۷	۱۳.۴۷۷	۰.۰۰۰
	درون گروهی	۴۱.۰۰۵	۷۴۳۴۰	۰.۰۰۱		
	مجموع	۴۱.۰۲۸	۷۴۳۴۳	-		
تأخر	بین گروهی	۳۴۱۹.۹۶۰	۳	۱۱۳۹.۹۸۷	۲۴۵۱۹.۳۲۴	۰.۰۰۰
	درون گروهی	۳۴۵۶.۳۱۹	۷۴۳۴۰	۰.۰۴۶		
	مجموع	۶۸۷۶.۲۷۸	۷۴۳۴۳	-		
مدت مراجعه	بین گروهی	۳۶۳۹.۱۸۱	۳	۱۲۱۳.۰۶۰	۱۰۰۳۰۲.۵۷۶	۰.۰۰۰
	درون گروهی	۸۹۹.۰۶۹	۷۴۳۴۰	۰.۰۱۲		
	مجموع	۴۵۳۸.۲۵۰	۷۴۳۴۳	-		
زمان بین	بین گروهی	۱۵۵۶۶.۱۹۳	۳	۵۱۸۸.۷۳۱	۱۵۳۳۶۶۷.۵۱۸	۰.۰۰۰
	درون	۲۵۱.۵۰۸	۷۴۳۴۰	۰.۰۰۳		

متغیر	منبع تغییر	مجموع مجدورات	درجه آزادی	میانگین مجدورات	آزمون F	معناداری
نشست‌ها	گروهی					
	مجموع	۱۵۸۱۷.۷۰۲	۷۴۳۴۳	-		
کالاهای سبد	بین گروهی	۰.۱۱۲	۳	۰.۰۳۷	۱۳۳.۵۳۹	۰.۰۰۰
	درون گروهی	۲۰.۸۴۸	۷۴۳۴۰	۰.۰۰۰		
	مجموع	۲۰.۹۶۱	۷۴۳۴۳	-		
ارزش مالی	بین گروهی	۰.۴۳۴	۳	۰.۱۴۵	۱۰۱.۵۴۹	۰.۰۰۰
	درون گروهی	۱۰۵.۹۱۰	۷۴۳۴۰	۰.۰۰۱		
	مجموع	۱۰۶.۳۴۴	۷۴۳۴۳	-		
دسته‌های کالای سبدها	بین گروهی	۱.۵۹۴	۳	۰.۵۳۱	۳۳۰.۴۱۶	۰.۰۰۰
	درون گروهی	۱۱۹.۵۳۸	۷۴۳۴۰	۰.۰۰۲		
	مجموع	۱۲۱.۱۳۲	۷۴۳۴۳	-		

### • تحلیل خوشه‌ها

اختصاص داده‌اند و همچنین چند درصد از کل هر گروه در هر خوشه است. نتایج بررسی تعداد هر برچسب در هر خوشه در جدول ۴ آورده شده است.

برای تحلیل بهتر خوشه‌ها، نخست بررسی می‌شود در هر خوشه هر مشتری چه برچسبی دارد، هر کدام از این گروه‌ها چه درصدی از کل خوشه را به خود

### جدول ۴. نتایج حاصل از بررسی تعداد برچسب در هر خوشه

خوشه	برچسب	تعداد نمونه	درصد در خوشه	درصد در برچسب	تعداد کل
خوشه اول	خرید	۴	۰.۰۱	۱.۳۴	۲۶۴۴۱
	انتظار	۱۰۸۷۸	۴۱.۱۴	۳۶	
	خریدنکردن	۱۵۵۵۹	۵۸.۸۵	۳۵.۵	
خوشه دوم	خرید	۱۷۱	۱.۴	۵۸.۳۳	۱۲۲۸۰
	انتظار	۳۲۰۵	۲۶.۱	۱۰.۶	
	نخریدن	۸۹۰۴	۷۲.۵	۲۰.۳	
خوشه سوم	خرید	۷۳	۰.۴۱	۲۵	۱۷۷۳۶
	انتظار	۸۱۰۴	۴۹.۵۹	۲۹.۱	
	خریدنکردن	۸۸۵۹	۵۰	۲۰.۲	
خوشه چهارم	خرید	۴۵	۰.۲۵	۱۵.۳۳	۱۷۸۸۷
	انتظار	۷۳۶۳	۴۱.۱۵	۲۴.۳	
	خریدنکردن	۱۰۴۷۹	۵۸.۶	۲۴	

خوشه اول: در خوشه نخست که بیشترین تعداد

باتوجه به نتایج جدول ۴ خوشه‌ها تحلیل می‌شود:

از کاربران دارای بیشترین تعداد کالا است که از متنوع‌ترین سبدها هم محسوب می‌شود. همچنین مبالغ سبدهای تشکیل شده از سایر خوشه‌ها بیشتر است.

**خوشه چهارم:** کاربرانی که در خوشه چهارم‌اند، همانند خوشه سوم تنها در بازه زمانی کوتاهی از سایت بازدید کرده‌اند؛ با این تفاوت که این بازدیدها اغلب در اواخر بازه مطالعه شده انجام گرفته است. همچنین فاصله زمانی بین بازدیدها در این گروه کم است. بعد از خوشه اول، بیشترین تعداد کاربرانی که تمایلی به تشکیل سبد نداشته‌اند، در این خوشه است؛ با این حال، باز هم سبدهای این گروه از کاربران دارای بیشترین تعداد و متنوع‌ترین کالاها بعد از خوشه سوم است.

#### • ارزیابی دسته‌بندی

برای ارزیابی درخت تصمیم از ماتریس درهم‌ریختگی استفاده می‌شود. این ماتریس ابزار مفیدی برای تحلیل چگونگی عملکرد روش دسته‌بندی در تشخیص داده‌ها یا مشاهدات دسته‌های مختلف است. مهم‌ترین معیار برای تعیین کارایی تکنیک دسته‌بندی، معیار دقت است. این پژوهش به دقت دسته‌بندی ۶۳.۶٪ دست یافت که با توجه به نوع دسته‌ها و تعداد نمونه‌های موجود در هر دسته، دقت قابل قبولی است. معیارهایی که به صورت جداگانه عملکرد یک دسته‌بند<sup>۱</sup> را روی دسته‌های مختلف برآورد می‌کنند حساسیت، شفافیت و صحت‌اند. این شاخص‌ها برای هر سه دسته با استفاده از ماتریس درهم‌ریختگی محاسبه شده‌اند. نتایج ارزیابی درخت تصمیم در جدول ۵ ارائه شده است.

کاربران را داشت، کاربرانی‌اند که در بازه مطالعه شده، تنها یک بار از سایت بازدید داشته و نسبت به خوشه‌های دیگر از صفحات مختلف بیشتری بازدید کرده‌اند. حدود نیمی از کاربران، برخی از این کالاها را به سبد خرید خود اضافه کرده‌اند؛ اما خرید خود را نهایی نکرده و سبد خرید آنها در حالت باز قرار دارد. نیمی دیگر از کاربران هم تمایلی به تشکیل سبد نداشته و تنها به بازدید از صفحات اکتفا کرده‌اند. در این خوشه، تعداد کاربرانی که بازدیدشان به خرید ختم شده بسیار محدود است.

**خوشه دوم:** در خوشه دوم که کمترین تعداد کاربران را داشت، کاربرانی هستند که در بازه مطالعه شده، بیشترین دفعات بازدید از سایت را داشته‌اند. نسبت به دیگر خوشه‌ها، کمترین مدت زمان بازدید از سایت و کمترین تعداد صفحات مختلف بازدیدشده در هر مراجعه را دارند و بیشتر به جست‌وجوی کالاها می‌پردازند. فاصله زمانی بین بازدیدهای این گروه از خوشه‌های دیگر بیشتر است. بیشترین تعداد سبدهایی که به خرید ختم شده است، در این گروه قرار دارد که می‌توان کاربران این خوشه را دارای احتمال خرید دانست.

**خوشه سوم:** کاربرانی که در خوشه سوم‌اند، در بازه زمانی کوتاهی از سایت بازدید داشته‌اند. همچنین کمترین فاصله زمانی بین بازدیدها هم متعلق به این گروه است. بعد از خوشه اول، بیشترین صفحات بازدیدشده متعلق به این خوشه است که می‌توان نتیجه گرفت بیشتر به جست‌وجوی کالاها در سایت و بازدید از آنها پرداخته‌اند. بعد از خوشه اول، بیشترین تعداد سبدهای باز در این خوشه قرار دارد. سبدهای این گروه

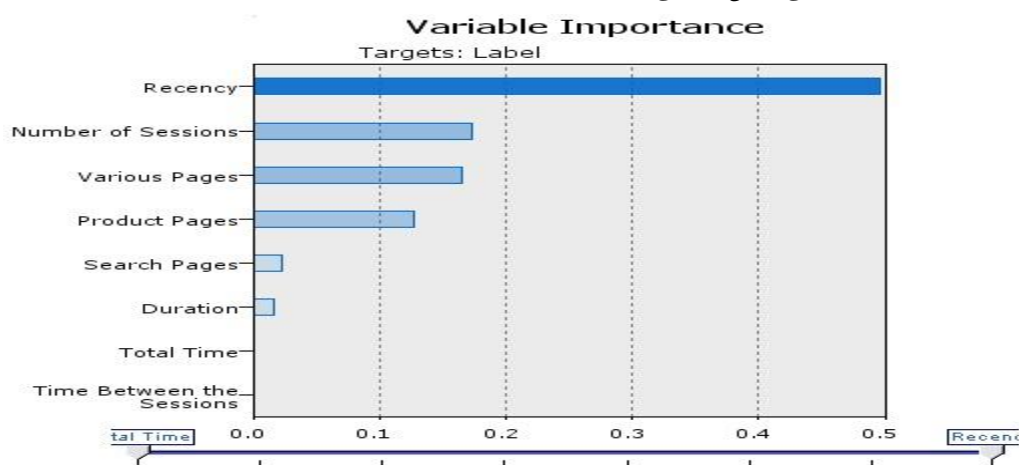
<sup>1</sup> . Classifier

## جدول ۵. نتایج ارزیابی درخت تصمیم

Precision-	Precision+	Recall-	Recall+	پرچسب
٪۹۹	٪۱.۶	٪۷۸	٪۸۰	خرید
٪۷۴	٪۷۶	٪۸۳	٪۶۶	انتظار خرید
٪۵۶	٪۸۲	٪۷۷	٪۶۵	خریدنکردن

بر اساس مدل (شکل ۱)، تأخر مهم ترین و مدت زمان و زمان بین نشست‌ها کم اهمیت ترین ویژگی‌ها در دسته بندی بوده اند.

با توجه به نامتعادل بودن دسته‌ها و دقت به دست آمده از اعتبارسنجی متقابل با ده تکرار یعنی ٪۶۳.۶، ثابت می‌شود که مدل دارای کارایی قابل قبولی است.



شکل ۱. ویژگی‌های مهم در دسته بندی

- ۲- کاربرانی با مدت مراجعه یک روزه و داشتن تعداد نشست‌های کمتر از ۱۲ و تعداد صفحات مختلف بزرگ تر از ۶، از سایت خرید نمی‌کنند.
- ۳- کاربرانی که تعداد نشست‌هایشان کمتر از ۱۲ و تعداد صفحات مختلف کمتر از ۶ و تعداد صفحات جست‌وجوی آنها بیشتر از ۳ است، خرید نمی‌کنند.
- ۴- کاربرانی که ۹ صفحه کالا را دیده‌اند و تعداد نشست‌هایشان بیشتر از ۱۲ است، زمانی را که درون سایت می‌گذرانند، بیشتر از ۲۸ دقیقه است و همچنین دارای مدت مراجعه بزرگ تر از ۶۷ هستند، از سایت خرید می‌کنند.
- ۵- کاربرانی که تأخرشان بین ۸۲ تا ۸۵ روز است،

## فاز ۶. به کارگیری

هدف از دسته بندی داده‌ها در این پژوهش، ارائه مدلی برای پیش بینی رفتار مشتریان و کشف دانشی مفید از پایگاه داده‌هاست. این مرحله شامل برنامه ریزی برای دانش کشف شده است که کجا و چگونه به کار رود. برنامه ای هم برای نظارت بر پیاده سازی دانش کشف شده باید ایجاد و تمام پروژه مستندسازی شود. کشف دانش از راه بررسی مدل و قوانین ایجاد شده انجام می‌گیرد. در این پژوهش ۲۶۱ قانون با اطمینان ٪۷۰ حاصل شد که تعدادی از آنها به شرح زیر است:

- ۱- اگر کاربر هیچ صفحه کالایی را مشاهده نکند، اصلاً خرید نخواهد کرد.

دارای تعداد نشست بالای ۱۲ و تعداد صفحات مختلف بیشتر از ۶ هستند، فقط تمایل به تشکیل سبد دارند.

## ۵- نتیجه گیری

شرکت‌های تجارت الکترونیک از دنیای قدیمی که در آن، محصولات استاندارد، بازارهای همگن و چرخه توسعه و عمر محصول طولانی یک قانون بود، به دنیای جدیدی انتقال می‌یابند که در آن انواع محصولات استاندارد جایگزین وجود دارند. مصرف‌کنندگان می‌توانند از بین میلیون‌ها کالا در یک فروشگاه برخط به جای ده‌ها هزار کالا در فروشگاه‌های بزرگ، کالاهای خود را انتخاب کنند. بنابراین، پیش‌بینی رفتار مشتریان دشوار است؛ چون به‌ندرت بازدید آنها از سایت‌ها به خرید ختم می‌شود؛ از این رو، برای داشتن کسب و کار برخط موفق باید به تحلیل رفتار مشتریان پرداخت. استفاده از داده‌کاوی و وب‌کاوی به ما در تحلیل رفتار مشتریان و کشف دانش در این زمینه کمک می‌کند. این دانش، سازمان‌ها را قادر می‌سازد تا تصمیمات مبتنی بر داده گرفته و راهبردهای تصمیم‌گیری خود را بهبود و توسعه دهند. همچنین برای به دست آوردن مشتریان جدید و حفظ مشتریان موجود و بهبود رضایت مشتری نیز کمک می‌کند. از آنجا که به نظر می‌رسد در ایران پژوهشی در این زمینه از تجارت الکترونیک وجود ندارد، این پژوهش با هدف طرح چارچوبی برای افزایش دقت تحلیل و شناخت گروه‌های مشتریان و همچنین ارائه مدل و قوانینی برای پیش‌بینی رفتار مشتریان بر مبنای ویژگی‌های رفتاری آنها انجام شده است. در این پژوهش، با استفاده از رویکرد ترکیبی داده‌کاوی و کاوش کاربری وب، رفتار مشتریان تحلیل شده است؛ این روش مزیت رقابتی برای شرکت را باعث می‌شود. در واقع نوآوری این پژوهش اعمال روش‌های

داده‌کاوی بر ابعاد مختلف داده‌های مشتری شامل داده‌های رفتاری مرور وب و داده‌های رفتار خرید است که به بهبود دانش از مشتری می‌انجامد. این پژوهش مبتنی بر داده‌های واقعی است و داده‌های واقعی مشتریان شرکت XYZ را بررسی می‌کند که محدوده فعالیت آن، تجارت الکترونیک و از نوع خرده‌فروشی اینترنتی است. در این پژوهش اطلاعات حساب‌های کاربری مشتریان به همراه اطلاعات موجود در لاگ‌های وب سرور که در پایگاه داده شرکت ذخیره می‌شوند، استفاده شده است. شرکت XYZ به منظور درک سازمانی بهتر نیازها، خواسته‌ها و تمایلات مشتریان و همچنین وجود انسجام و دیدی مشخص در سازمان برای پیام‌های بازاریابی و تعریف مشتریان بالقوه، می‌تواند از خوشه‌های به‌دست آمده در این پژوهش، برای این منظور بهره‌برد و به شناخت دقیقی از مشتریان خود دست یابد و می‌تواند راهبردهای مقتضی را در هر گروه پیاده‌سازی کند و احتمال موفقیت خود را افزایش دهد. در کل خرده‌فروشان برخط باید فعالیت‌های بازاریابی خود در رسانه‌های اجتماعی را گسترش داده تا بتوانند متناسب با نیازها و فعالیت‌های مشتریان خود (کاربران این رسانه‌ها) عمل کنند و فعالیت‌های جذب مشتری جدید و حفظ مشتریان موجود را ادامه دهند.

محدودیت‌های زیادی در زمینه استفاده از داده‌های دنیای واقعی برای تحلیل رفتار مشتریان وجود دارد. به دلیل وجود بازار رقابتی، شرکت‌ها به‌سادگی داده‌های مشتریان خود را در اختیار قرار نمی‌دهند یا در صورت انجام این کار، داده‌های دست‌کاری شده و ناقص در اختیار آنها قرار می‌دهند. در واقع مدیران شرکت برای سپردن اطلاعات به افراد خارج از سازمان تمایلی ندارند و چون نتایج داده‌کاوی کاملاً به صحت داده‌های اولیه بستگی دارد، این امر ممکن است سبب به دست آمدن

- (2008). KDD, SEMMA and CRISP-DM: a parallel overview. *IADS-DM*.
- 4- Bae, S. M., Park, S. C., & Ha, S. H. (2003). Fuzzy web ad selector based on web usage mining. *IEEE intelligent Systems*, 18(6), 62-69.
- 5- Beheshtian-Ardakani, A., Fathian, M., & Gholamian, M. (2018). A novel model for product bundling and direct marketing in e-commerce based on market segmentation. *Decision Science Letters*, 7(1), 39-54.
- 6- Deka, P. K. (2017). A Conceptual Model for Determining Factors Influencing Online Purchasing Behavior. *Journal of Management in Practice*, 2(1).
- 7- Dharmarajan, K., & Dorairangaswamy, D. M. (2016). Web Usage Mining: Improve The User Navigation Pattern Using Fp-Growth Algorithm. *Elysium journal of engineering research and management (EJERM)*, 3(4).
- 8- Facca, F. M., & Lanzi, P. L. (2005). Mining interesting knowledge from weblogs: a survey. *Data & Knowledge Engineering*, 53(3), 225-241.
- 9- Ha, S. H. (2002). Helping online customers decide through web personalization. *IEEE Intelligent systems*, 17(6), 34-43.
- 10- Hsieh, N. C. (2004). An integrated data mining and behavioral scoring model for analyzing bank customers. *Expert systems with applications*, 27(4), 623-633.
- 11- Hung, Y. S., Chen, K. L. B., Yang, C. T., & Deng, G. F. (2013). Web usage mining for analysing elder self-care behavior patterns. *Expert Systems with Applications*, 40(2), 775-783.
- 12- Jenamani, M., Mohapatra, P. K., & Ghose, S. (2003). A stochastic model of e-customer behavior. *Electronic Commerce Research and Applications*, 2(1), 81-94.
- 13- Kim, E., Kim, W., & Lee, Y. (2003). Combination of multiple classifiers for the customer's purchase behavior prediction. *Decision Support*

نتایج نامعتبر شود. همچنین در برخی از شرکت‌ها، ممکن است کارمندان به دلیل ترس از دست دادن موقعیت شغلی خود، حاضر به همکاری با پژوهشگران نباشند و در شناخت کسب و کار و داده‌ها کمک چندانی نکنند.

در این پژوهش به دلیل زیاد بودن حجم داده‌ها و کم بودن مقدار حافظه، امکان اجرای تکنیک‌های دیگر داده کاوی مانند الگوریتم SOM وجود نداشت؛ از این رو پیشنهاد می‌شود پژوهش‌های آینده از سایر الگوریتم‌های داده کاوی که به مشخص کردن تعداد خوشه‌های بهینه نیاز ندارند و خودشان این تعداد را به دست می‌آورند، برای خوشه‌بندی استفاده کنند.

در این پژوهش داده‌های دموگرافیک مشتریان شامل جنس، سن، مکان سکونت، تحصیلات و غیره به دلیل برخی مسائل فنی و امنیتی در دسترس نبودند؛ از این رو پیشنهاد می‌شود این ویژگی تحلیل رفتار مشتریان بررسی شود. علاوه بر آن، پیشنهاد می‌شود ویژگی‌های دیگری مانند نوع خرید، روش خرید، نوع پرداخت و غیره استخراج و تأثیر آنها در این رویکرد ترکیبی بررسی شود. همچنین بررسی تأثیر وب کاوی بر ارزش مشتری برای پاسخ به این سؤال که آیا وب کاوی می‌تواند ارزش مشتری را پیش‌بینی کند یا خیر، نیز می‌تواند در پژوهش‌های آینده بررسی شود.

## منابع

- 1- Ahn, T., Ryu, S., & Han, I. (2007). The impact of Web quality and playfulness on user acceptance of online retailing. *Information & management*, 44(3), 263-275.
- 2- Arora, M., & Chopra, A. B. (2016). Impact of online selling on physical retail in India. *International Journal of Research in IT and Management*, 6(10), 57-68.
- 3- Azevedo, A. I. R. L., & Santos, M. F.



- Berlin, Heidelberg.
- 23- Theusinger, C., & Huber, K. P. (2000, August). Analyzing the footsteps of your customers. In *Proc. of the Sixth ACM SIGKDD Internat. Conf. on Web KDD 2000* (pp. 44-52).
  - 24- Tsai, C. F., Hu, Y. H., & Lu, Y. H. (2015). Customer segmentation issues and strategies for an automobile dealership with two clustering techniques. *Expert Systems*, 32(1), 65-76.
  - 25- Weiss, G. M., McCarthy, K., & Zabar, B. (2007). Cost-sensitive learning vs. sampling: Which is best for handling unbalanced classes with unequal error costs?. *DMIN*, 7, 35-41.
  - 26- Yang, Z., & Su, X. (2012). Customer behavior clustering using SVM. *Physics Procedia*, 33, 1489-1496.
  - 27- Zhang, X., Gong, W., & Kawamura, Y. (2004, January). Customer behavior pattern discovering with web mining. In *APWeb* (pp. 844-853). *Systems*, 34(2), 167-175.
  - 14- Liou, J. J., & Tzeng, G. H. (2010). A dominance-based rough set approach to customer behavior in the airline market. *Information Sciences*, 180(11), 2230-2238.
  - 15- Liu, B., Mobasher, B., & Nasraoui, O. (2011). Web usage mining. In *Web Data Mining* (pp. 527-603). Springer Berlin Heidelberg.
  - 16- Moghaddam, Q., S. Abdolvand, N., & Harandi, R. S. (2017). A RFMV Model and Customer Segmentation Based on Variety of Products. *Information Systems & Telecommunication*, 5(3), 155- 161
  - 17- Pamutha, T., Chimphee, S., Kimpan, C., & Sanguansat, P. (2012). Data preprocessing on web server log files for mining users access patterns. *International Journal of Research and Reviews in Wireless Communications (IJRRWC) Vol, 2*.
  - 18- Park, J., & Chung, H. (2009). Consumers' travel website transferring behaviour: analysis using clickstream data-time, frequency, and spending. *The Service Industries Journal*, 29(10), 1451-1463.
  - 19- Shanthi, S. (2017). Survey on Web Usage Mining using Association Rule Mining. *International Journal of Innovative Computer Science & Engineering*, (4) 3; 65-67
  - 20- Sheikh, A. M., & Menaria, S. (2017). An Approach of Security in E-Commerce with Web Mining Framework. *International Education and Research Journal*, 3(5).
  - 21- Sisodia, D. S., & Verma, S. (2012, May). Web usage pattern analysis through web logs: A review. In *Computer Science and Software Engineering (JCSSE), 2012 International Joint Conference on* (pp. 49-53). IEEE.
  - 22- Sun, L., & Zhang, X. (2004, April). Efficient frequent pattern mining on web logs. In *Asia-Pacific Web Conference* (pp. 533-542). Springer,

