



## Predicting the Purchase of Self-Employed Pension Schemes in the Iranian Social Security Organization Using Decision Tree and Random Forest Classification Algorithms

Yasin Ghasemi<sup>1</sup>  | Abbas Khandan<sup>2\*</sup>  | Narges Akbarpour Roshan<sup>3</sup> 

1. Master Student, Faculty of Economics, Kharazmi University, Tehran, Iran.  
Email: Yasin8769@gmail.com
2. Assistant Professor, Faculty of Economics, Kharazmi University, Tehran, Iran.  
(Corresponding Author), Email: khandan.abbas@khu.ac.ir
3. Assistant Professor, Faculty of Economics, Institute for Humanities and Cultural Studies, Tehran, Iran. Email: n.akbarpour@ihcs.ac.ir

Article Info	ABSTRACT
<b>Article type:</b> Research Article	The pension coverage of the Iranian Social Security Organization for self-employed workers is offered at three contribution rates of 12, 14 and 18 percent, but looking at the statistics shows that the demand for these types of insurances is low. This research investigates the characteristics of these insured groups by using data mining and applying two machine learning algorithms, decision tree and random forest, and predicts their behavior by providing a classification model. This will help the Social Security Organization to improve customer relationship management. For this purpose, the information of 1286174 insured persons of self-employed in 2020 was used, which includes the characteristics of age, gender, average monthly income, the years of service, and the type of self-employed pension scheme. The obtained results show that women mainly apply for the scheme with 12 percent contribution, while men tend to be covered by schemes with contribution rates of 14
<b>Article history:</b> Received: 2022/11/15	
<b>Received in revised form:</b> 2022/11/15	
<b>Accepted:</b> 2023/04/15	
<b>Keywords:</b> Pension Insurance of self-employed, Social Security Organization, Data mining, Machine learning, Classification	
<b>JEL:</b> H55 .C81 J26.	

---

and 18 percent due to the burden of supporting the family. Also, for men, the demand for schemes of 14 and 18 percent increases with the increase of age, income and years of service, but there are no such trends for women. According to the obtained results, years of service and then gender are decisive in choosing the type of pension scheme in such a way that according to the prediction of the model, people with less than 4.5 years of service are known as definite applicants for 12 percent self-employed pension scheme.

---

**Cite this article:** Ghasemi, Yasin; Khandan, Abbas; Akbarpour Roshan, Narges. (2023). Predicting the Purchase of Self-Employed Pension Schemes in the Iranian Social Security Organization Using Decision Tree and Random Forest Classification Algorithms. *Journal of Economic Modeling Research*, 13 (47), 115-165.  
DOI: 00000000000000000000



© The Author(s).

Publisher: Kharazmi University

---



## پیش بینی خرید بیمه نامه حرف و مشاغل آزاد سازمان تامین اجتماعی با استفاده از الگوریتم طبقه بندی درخت تصمیم و جنگل تصادفی

یاسین قاسمی<sup>۱</sup> | عباس خندان\*<sup>۲</sup> | نرگس اکبرپور روشن<sup>۳</sup>

۱. دانشجوی کارشناسی ارشد مهندسی صنایع، دانشکده اقتصاد، دانشگاه خوارزمی، تهران، ایران yasein8769@gmail.com

۲. استادیار اقتصاد، دانشکده اقتصاد، دانشگاه خوارزمی، تهران، ایران. (نویسنده مسئول) khandan.abbas@khu.ac.ir

۳. استادیار اقتصاد، پژوهشکده اقتصاد، پژوهشگاه علوم انسانی و مطالعات فرهنگی، تهران، ایران n.akbarpour@ihcs.ac.ir

اطلاعات مقاله	چکیده
<b>نوع مقاله:</b> مقاله پژوهشی	پوشش بیمه سازمان تامین اجتماعی برای حرف و مشاغل آزاد به صورت اختیاری در سه نرخ ۱۲، ۱۴ و ۱۸ درصد ارائه می شود اما نگاه به آمار نشان می دهد که تقاضای این بیمه نامه ها بسیار پایین است. این پژوهش با استفاده از داده کاوی و با به کارگیری دو الگوریتم یادگیری ماشین یعنی درخت تصمیم و جنگل تصادفی به بررسی مشخصه های خریداران این نوع بیمه نامه ها پرداخته و با ارائه یک مدل طبقه بندی، رفتار آن ها را پیش بینی می کند تا از این طریق به سازمان تامین اجتماعی در جهت بهبود مدیریت ارتباط با مشتری کمک کند.
<b>تاریخ دریافت:</b> ۱۴۰۱/۰۸/۲۴	برای این منظور، از اطلاعات ۱۲۸۶۱۷۴ نفر از خریداران انواع بیمه نامه های حرف و مشاغل آزاد سال ۱۳۹۹ استفاده شد که مشخصه های سن، جنسیت، متوسط درآمد ماهانه، میزان سابقه کار و نوع بیمه نامه خریداری شده را در بر می گیرد.
<b>تاریخ ویرایش:</b> ۱۴۰۱/۰۸/۲۴	نتایج به دست آمده نشان می دهند که زنان به طور عمده متقاضی بیمه نامه با نرخ ۱۲ درصد هستند در حالی که مردان به دلیل بر عهده داشتن بار تکفل خانواده عمدتاً تمایل به خرید بیمه نامه های با نرخ ۱۴ و ۱۸ درصدی دارند.
<b>تاریخ پذیرش:</b> ۱۴۰۲/۰۱/۲۶	همچنین، در مردان با افزایش سن، درآمد و سابقه، تقاضای بیمه های با نرخ ۱۴ و ۱۸ درصد افزایش می یابد، اما چنین روندهایی برای زنان وجود ندارد.
<b>واژه های کلیدی:</b> بیمه بازنشستگی حرف و مشاغل آزاد، سازمان تامین اجتماعی، داده کاوی، یادگیری ماشین، طبقه بندی.	طبق نتایج به دست آمده متغیرهای میزان سابقه کار و پس از آن جنسیت در انتخاب نوع بیمه نامه تعیین کننده هستند، به گونه ای که طبق پیش بینی مدل افراد با سابقه کار کمتر از ۴/۵ سال متقاضیان قطعی بیمه نامه ۱۲ درصدی شناخته
<b>طبقه بندی JEL:</b> H55, C81, J26.	

---

شده‌اند. با توجه به نتایج و انگیزه پایین زنان و جوانان برای انتخاب بیمه‌های با خدمات گسترده‌تر، سازمان تأمین اجتماعی می‌تواند از طریق ارائه مشوق‌ها یا خدمات کوتاه‌مدت، جذابیت این نوع بیمه‌ها با خدمات گسترده‌تر را در بین این گروه خاص افزایش دهد.

---

**استناد:** قاسمی، یاسین؛ خندان، عباس؛ اکبرپور روشن، نرگس. (۱۴۰۱). پیش‌بینی خرید بیمه‌نامه حرف و مشاغل آزاد سازمان تأمین اجتماعی با استفاده از الگوریتم طبقه‌بندی درخت تصمیم و جنگل تصادفی؛ تحقیقات مدل‌سازی اقتصادی، ۱۳ (۴۷)، ۱۱۵-۱۶۵.

DOI: 0000000000000000000000



© نویسنده‌گان.

ناشر: دانشگاه خوارزمی.

## ۱. مقدمه

پوشش ریسک کاهش درآمد در دوران سالمندی از اساسی‌ترین نیازهای انسان است که به لحاظ اقتصادی، سیاسی و اجتماعی بسیار اهمیت داشته و همواره مورد توجه دولت‌ها بوده است. سازمان تامین اجتماعی ایران بزرگ‌ترین سازمان فعال در حوزه بیمه‌های اجتماعی و بازنشستگی کشور محسوب می‌شود و طبق آخرین سالنامه منتشر شده در سال ۱۳۹۹ جمعیتی حدود ۴۴ میلیون و پانصد هزار نفر را تحت پوشش دارد. این سازمان یک بیمه‌گر اجتماعی است که مأموریت اصلی آن پوشش کارگران مزد و حقوق‌بگیر (به صورت اجباری) و صاحبان حرف و مشاغل آزاد (به صورت اختیاری) در برابر ریسک‌های مختلفی از جمله ریسک فوت، از کارافتادگی و کاهش درآمد در دوران سالمندی است. بازنشستگان و از کارافتادگان تحت پوشش تا زمان حیات خود از مستمری برخوردار می‌شوند، و در صورت فوت بیمه‌شده اصلی، چه در دوران اشتغال و چه بازنشستگی، مستمری به بازماندگان مضمول آن‌ها پرداخت می‌شود. این وجه در واقع ریسک فوت بیمه‌شده اصلی را پوشش می‌دهد.

بیمه حرف و مشاغل آزاد تامین اجتماعی که در گروه بیمه‌های خویش‌فرما قرار می‌گیرد، شامل افرادی می‌شود که حرفه و شغل آزاد دارند (به تنهایی) یا تحت عنوان کارفرما (با داشتن کارگر) در حال فعالیت هستند که می‌توانند پس از بازنشستگی از مزایای مستمری آن استفاده نمایند. طبق آخرین سالنامه منتشر شده، بیمه‌شدگان اصلی سازمان حدود ۱۴ میلیون و ۵۰۰ هزار نفر هستند که از این تعداد حدود ۱ میلیون و ۱۰۰ هزار نفر (۸ درصد) بیمه‌شده حرف و مشاغل آزاد هستند. این آمار نشان می‌دهد استقبال شاغلین حرف و مشاغل آزاد از این نوع بیمه‌نامه‌ها کم است و سازمان نیازمند ایجاد سازوکارهایی برای جذب حداکثری افراد به این نوع بیمه است. مطالعات بسیار کمی به بررسی بیمه حرف و مشاغل آزاد سازمان تامین اجتماعی و انگیزه‌های افراد از تقاضای آن پرداخته‌اند، اما همین مطالعات اندک هم به خوبی نشان می‌دهند که هنوز بسیاری از مردم متقاضی برخورداری از این بیمه‌نامه‌ها نیستند و در بسیاری از موارد قراردادهای این بیمه‌نامه‌ها پس از یک مدت بیمه‌پردازی، قطع می‌شود.

از این جهت، مطالعه و بررسی ویژگی‌های افرادی که متقاضی این نوع بیمه‌نامه‌ها هستند بسیار اهمیت دارد، چراکه می‌تواند در رفع مشکل ضریب نفوذ پایین این بیمه‌نامه‌های بازنشستگی اختیاری کمک‌کننده باشد. البته دستیابی به این هدف و بررسی ویژگی‌های افراد تحت پوشش این نوع بیمه‌ها با پیچیدگی‌هایی همراه است، اما ابزارهای جدید مانند یادگیری ماشین<sup>۱</sup> به راحتی و بدون مداخله مستقیم انسانی در روند کار، می‌توانند با حجم زیادی از داده‌های چندبعدی و چندمتغیره کار کنند و الگوها و روندهای موجود را شناسایی و در نتیجه باعث افزایش دقت و سرعت انجام کار شوند. این مقاله نخستین مطالعه‌ای است که خرید بیمه‌نامه‌های حرف و مشاغل آزاد تأمین اجتماعی را با استفاده از شیوه‌های یادگیری ماشینی و داده‌کاوی مورد بررسی قرار می‌دهد.

در مقاله حاضر به دو دسته از انواع بیمه‌نامه‌های حرف و مشاغل آزاد شامل پوشش بازنشستگی با حق بیمه ۱۲ درصد<sup>۲</sup> و پوشش بازنشستگی و سایر ریسک‌ها با حق بیمه ۱۴ و ۱۸ درصدی (پوشش از کارافتادگی علاوه بر بازنشستگی با حق بیمه ۱۴ درصد و پوشش ریسک از کارافتادگی و مستمری بازماندگان علاوه بر بازنشستگی با حق بیمه ۱۸ درصد) پرداخته خواهد شد و از آنجایی که تعداد افراد با نرخ ۱۴ درصد نسبت به سایر نرخ‌ها پایین تر است، آن را با ۱۸ درصد در یک دسته قرار دادیم. در ادامه تلاش می‌شود با تجزیه و تحلیل الگوهای خرید پوشش بیمه‌ای توسط افراد مختلف و بررسی ویژگی‌های افرادی که تحت پوشش این نوع بیمه‌نامه‌ها هستند به روشن شدن چرایی پایین بودن ضریب نفوذ این بیمه‌نامه‌ها کمک کرده و دانش جدیدی که از بررسی خصوصیات مشتریان به دست می‌آید به عنوان توصیه سیاستی در جهت بهبود فعالیت‌های سازمان برای جذب بیمه‌شده پیشنهاد شود. نتیجه این تجزیه و تحلیل می‌تواند باعث بهبود مدیریت روابط سازمان با مشتریان و همچنین افزایش جمعیت تحت پوشش سازمان شود. در این راستا، با پیش پردازش مناسب داده‌های

#### 1. Machine Learning

۲. علاوه بر خدمات بازنشستگی، ریسک فوت در زمان بازنشستگی را نیز پوشش می‌دهد.

۳. علاوه بر این نرخ، ۲ درصد هم سهم دولت است که باید به سازمان پرداخت کند.

بیمه‌شدگان حرف و مشاغل آزاد که شامل اطلاعات شخصی آن‌ها مانند سن، جنسیت، میزان درآمد و سابقه کار است مدلی ساخته خواهد شد که به پیش‌بینی تقاضای بیمه‌نامه مناسب هر مشتری با توجه به ویژگی‌های فردی آن‌ها کمک می‌کند (متقاضیان بیمه‌نامه با خدمات مستمری بازنشستگی یا افرادی که علاوه بر بازنشستگی، متقاضی خدمات از کارافتادگی و بازمانگان خواهند بود). طبقه‌بندی مشتریان در نهایت به مدیریت بهتر مستمری‌ها و هزینه‌ها منجر می‌شود.

این مقاله در پنج بخش سامان‌دهی شده است. در ادامه و در بخش دوم به ادبیات موضوع و پیشینه پژوهش پرداخته خواهد شد. بخش سوم مقاله به روش‌شناسی و توصیف داده‌ها اختصاص دارد. در این بخش وظایف و مراحل روش کریپس<sup>۱</sup> که یکی از معروف‌ترین روش‌های داده‌کاوی و مورد استفاده در این مقاله است بیان خواهد شد و مجموعه‌ای از آمار توصیفی از نحوه توزیع مشخصه‌های افراد و مدل‌سازی شامل دو الگوریتم درخت تصمیم و جنگل تصادفی ارائه می‌شود. در بخش چهارم یافته‌های پژوهش به همراه ارزیابی نتایج مدل‌سازی و پیش‌بینی‌ها مطرح می‌شود. در پایان و در بخش پنجم جمع‌بندی و نتیجه‌گیری مقاله به همراه پیشنهادهایی کاربردی برای سازمان تامین اجتماعی ارائه خواهد شد.

## ۲. ادبیات موضوع

سازمان تامین اجتماعی یک سازمان بیمه‌گر اجتماعی است که مأموریت اصلی آن پوشش کارگران مزد و حقوق‌بگیر به صورت اجباری و صاحبان حرف و مشاغل آزاد به صورت اختیاری است. این سازمان یک سازمان عمومی غیردولتی است که عمده منابع مالی آن از محل حق بیمه‌ها (با مشارکت بیمه‌شده، کارفرما و دولت) تامین می‌شود و ارائه خدمات بلندمدت (مستمری‌های بازنشستگی، از کارافتادگی، بازمانگان)، کوتاه‌مدت (حمایت در برابر حوادث، بیماری‌ها و بارداری، غرامت دستمزد ایام بیماری، غرامت دستمزد ایام

بارداری، پرداخت هزینه وسایل کمک پزشکی، کمک هزینه ازدواج، کمک هزینه کفن و دفن) و همچنین خدمات درمان و سلامت و مقرری بیمه بیکاری را بر عهده دارد. بر این اساس می‌توان سازمان تامین اجتماعی را در دسته سیستم‌هایی قرار داد که اغلب ریسک‌های مهم ذکر شده در استانداردهای سازمان بین‌المللی کار (ILO) را پوشش می‌دهد (برای مطالعه بیشتر ر.ک. مقاله‌نامه ۱۰۲، با عنوان حداقل استانداردهای تامین اجتماعی<sup>۱</sup>).

بیمه حرف و مشاغل آزاد سازمان تامین اجتماعی اختیاری است و شامل افرادی می‌شود که غیر مزد و حقوق‌بگیر (تحت عنوان کارفرما یا بدون کارفرما) هستند. این افراد اختیار دارند که انواع بیمه‌نامه‌ها در سه سطح حق بیمه ۱۲ درصد (خدمات مستمری بازنشستگی)، حق بیمه ۱۴ درصد (خدمات مستمری بازنشستگی و ازکارافتادگی) و حق بیمه ۱۸ درصد (خدمات مستمری بازنشستگی، ازکارافتادگی و فوت یا بازمندگان) را انتخاب و در هر مقطعی قرارداد خود را فسخ کنند. طبق آخرین آمار منتشر شده در سال ۱۳۹۹ حدود ۱ میلیون و ۱۰۰ هزار نفر از بیمه‌شدگان سازمان را بیمه‌شدگان حرف و مشاغل آزاد تشکیل می‌دهند که تنها ۸ درصد از بیمه‌شدگان اصلی سازمان هستند. این در حالی است که در همین سال طبق گزارش مرکز آمار ایران، خوداشتغالی (کارکن مستقل و کارفرما) سهمی ۴۷ درصدی از بازار کار داشته‌اند.

برای رفع مسئله پایین بودن ضریب نفوذ و خریداری بیمه‌نامه‌های حرف و مشاغل آزاد لازم است تا سازمان تامین اجتماعی از طریق سازماندهی فرآیندهای بازاریابی و ارائه خدمات به مشتریان، باعث ایجاد یک سیستم تعاملی بین مشتریان و سازمان شود. از این طریق می‌توان ارتباط نزدیک‌تر و دقیق‌تری با خواسته‌های مشتریان ایجاد و آن‌ها را به طور کامل با خدمات سازمان آشنا کرد. این مفاهیم تحت عنوان مدیریت ارتباط با مشتری<sup>۲</sup> شناخته می‌شود.

## ۲-۱. مبانی نظری

1. Social Security (Minimum Standards) Convention  
2. Customer Relationship Management



مدیریت ارتباط با مشتری ترکیبی از افراد، فرآیندها و فناوری است که به دنبال درک مشتریان و ایجاد یک رویکرد یکپارچه جهت مدیریت روابط با تمرکز بر حفظ مشتری و توسعه روابط می‌باشد. مدیریت ارتباط با مشتری از پیشرفت‌های فناوری اطلاعات و تغییرات سازمانی در فرآیندهای مشتری محور تکامل یافته و شرکت‌هایی که مدیریت ارتباط با مشتری را با موفقیت پیاده‌سازی کنند می‌توانند از وفاداری مشتری و سودآوری بلندمدت بهره‌مند شوند (چن و پوپوویچ<sup>۱</sup>، ۲۰۰۳).

مدیریت ارتباط با مشتری در بازاریابی در چهار بعد اهمیت بیشتری یافته که شامل: شناسایی مشتری<sup>۲</sup>، جذب مشتری<sup>۳</sup>، نگهداری مشتری<sup>۴</sup> و توسعه مشتری<sup>۵</sup> می‌باشد (حسینی و همکاران، ۲۰۱۰). این چهار بعد را می‌توان به صورت یک چرخه بسته از سیستم مدیریت مشتری در نظر گرفت که هدف مشترک آن‌ها ایجاد درک عمیق‌تر از مشتریان برای به حداکثر رساندن ارزش مشتری برای سازمان در درازمدت می‌باشد (ان‌گای<sup>۶</sup> و همکاران، ۲۰۰۹). بُعد اول شناسایی مشتری است که شامل هدف قرار دادن جمعیتی است که به احتمال زیاد مشتری می‌شوند یا بیشترین سود را برای شرکت دارند. علاوه بر این شامل تجزیه و تحلیل مشتریانی است که در رقابت از دست می‌روند و اینکه چگونه می‌توان آن‌ها را به دست آورد (کراکلانور<sup>۷</sup> و همکاران، ۲۰۰۴). پس از شناسایی مشتریان، سازمان می‌تواند منابع خود را به سمت جذب مشتریان هدف هدایت کند، یکی از عناصر جذب مشتری بازاریابی مستقیم است. مرحله سوم، نگهداری مشتری نیز از اهمیت بالایی برخوردار است و شامل عناصر بازاریابی تک به تک<sup>۸</sup>، برنامه‌های وفاداری و مدیریت شکایت می‌باشد. مرحله چهارم توسعه مشتری نام دارد که به معنای افزایش با ثبات و مداوم تراکنش‌های مالی، ارزش

1. Chen & Popovich
2. Customer identification
3. Customer attraction
4. Customer retention
5. Customer development
6. Ngai
7. Kracklauer

۸. به کمپین‌های بازاریابی شخصی اطلاق می‌شود که با تجزیه و تحلیل، شناسایی و پیش‌بینی تغییرات در رفتارهای مشتری پشتیبانی می‌شود.

تراکنش‌ها و سودآوری مشتری است، و شامل تجزیه و تحلیل ارزش طول عمر مشتری<sup>۱</sup>، توسعه فروش و تجزیه و تحلیل سبد خرید می‌باشد (ان‌گای و همکاران، ۲۰۰۹). این مراحل چهارگانه مدیریت ارتباط با مشتری باعث ایجاد ارزش برای مشتریان و سازمان تأمین اجتماعی شده و به حفظ مشتریان فعلی و جذب مشتریان جدید برای این سازمان منتج خواهد شد.<sup>۲</sup> البته مدیریت ارتباط با مشتری وابسته به تحلیل آمار و اطلاعات مشتریان است و به این منظور داده‌کاوی<sup>۳</sup> اطلاعات مشتریان و خریداران بیمه‌نامه‌ها لازم است. داده‌کاوی با کشف روندهای عمومی، دانش مستقیم و غیر مستقیم زیادی برای مدیران و تصمیم‌گیرندگان به همراه دارد و نتیجه آن اتخاذ تصمیم‌های درستی است که می‌تواند باعث افزایش کارایی و عملکرد سازمان شود.

داده‌کاوی در واقع فرآیندی با استفاده از ابزارهای توصیفی<sup>۴</sup> و تحلیلی برای جستجوی مدل‌ها، الگوها و رابطه بین داده‌های تاریخی و ایجاد مدل‌هایی است که کارکرد آن‌ها پیش‌بینی<sup>۵</sup> است. بنابراین داده‌کاوی را می‌توان این‌گونه توصیف کرد: با توجه به هدف تعیین شده یک سازمان، روشی موثر و پیشرفته در کاوش و تجزیه و تحلیل و مدیریت حجم زیادی از داده‌های سازمانی و آشکارسازی قوانین ناشناخته و مدل‌سازی آن‌ها برای پیش‌بینی موارد آینده است (چن و هو، ۲۰۰۵). داده‌کاوی که مبتنی بر اصول آمار است، فرآیند کاوش، تجزیه و تحلیل مقادیر زیادی از داده‌ها را برعهده دارد تا الگوهای موجود در آن‌ها را کشف کند. از الگوریتم‌ها برای یافتن روابط و الگوها در داده‌ها استفاده می‌شود و سپس این اطلاعات در مورد الگوها برای پیش‌بینی<sup>۶</sup> و تخمین<sup>۷</sup> به کار می‌رود. به طور کلی، هدف

۱. تحلیل ارزش طول عمر مشتری به عنوان پیش‌بینی کل درآمد خالصی است که یک شرکت می‌تواند از مشتری انتظار داشته باشد.

۲. البته باید به این نکته توجه کرد که منظور از مشتری در اینجا بیمه‌شده است که ممکن است الزاماً ویژگی‌ها و مؤلفه‌های مشتری به معنای مرسوم را نداشته باشد.

3. Data mining  
4. Descriptive  
5. Predictive  
6. Chen & Hu  
7. Predictive  
8. Forecasting

داده‌کاوی استخراج داده‌ها از یک مجموعه داده بزرگ‌تر برای اهداف طبقه‌بندی<sup>۱</sup> یا پیش‌بینی است. از نظر جهت‌گیری فرآیند، فعالیت‌های داده‌کاوی به سه دسته کلی تقسیم می‌شوند (ریجسکی<sup>۲</sup> و همکاران، ۲۰۰۲): کشف<sup>۳</sup> به معنی فرآیند جستجو در پایگاه داده برای یافتن الگوهای پنهان بدون یک ایده یا فرضیه از پیش تعیین شده در مورد اینکه الگوها ممکن است چه باشند؛ مدلسازی پیش‌بینی کننده<sup>۴</sup> که به فرآیند گرفتن الگوهای کشف شده از پایگاه داده و استفاده از آن‌ها برای پیش‌بینی آینده گفته می‌شود؛ و بررسی صحت مدل<sup>۵</sup> که فرآیند بکارگیری الگوهای استخراج شده برای یافتن الگوهای داده غیرمعمول و بی‌شباهت است. داده‌کاوی کاربرد زیادی در تجزیه و تحلیل اطلاعات مربوط به مشتریان دارد و می‌تواند در تمامی مراحل ارتباط با مشتری از جمله شناسایی مشتری، جذب مشتری، نگهداری و توسعه مشتری استفاده شود.

## ۲-۲. پیشینه پژوهش

مطالعات بسیار کمی به بررسی بیمه حرف و مشاغل آزاد سازمان تأمین اجتماعی و انگیزه‌های افراد از پوشش اختیاری تحت این بیمه‌نامه‌ها پرداخته‌اند. عبدی (۱۳۸۵) به بررسی مشکلات بیمه‌شدگان حرف و مشاغل آزاد و اختیاری سازمان تأمین اجتماعی پرداخته است. در این مطالعه یک نمونه ۴۰۰ نفری از بیمه‌شدگان حرف و مشاغل آزاد و اختیاری مراجعه کننده به ۱۱ شعبه سازمان تأمین اجتماعی تهران و شهری (از مجموع ۷۱۲۵۶ نفر بیمه‌شده اختیاری و ۴۷۸۴۴ نفر حرف و مشاغل آزاد تحت پوشش سازمان در ۳۰ شعبه) در پایان سال ۱۳۸۳ به صورت پرسشنامه‌ای مورد پرسش قرار گرفته‌اند و متغیرهای فردی (سن، جنسیت، تحصیلات، وضعیت تاهل، وضع درآمد و...) و متغیرهای بیمه‌ای آنان (سابقه بیمه، سن در زمان تقاضا، انواع پوشش انتخابی، مبلغ بیمه و...) بررسی شده است. نتایج به دست آمده از

---

1. Classification  
2. Rygielski  
3. Discovery  
4. Predictive Modeling  
5. Forensic Analysis

این مطالعه میدانی نشان می‌دهد که حدود ۹۰ درصد بیمه‌شدگان حرف و مشاغل آزاد را مردان تشکیل می‌دهند که فقط حدود ۶ درصد مجرد هستند. متوسط سن بیمه‌شدگان در مقطع پژوهش حدود ۴۵ سال بوده و افراد دارای تحصیلات بالاتر در مقایسه با جامعه آماری سهم بیشتری در نمونه داشته‌اند. سطوح تحصیلی زنان بیمه‌شده در مقایسه با سطوح تحصیلی مردان بیمه‌شده بیشتر بوده است. همچنین متوسط سطح درآمد افراد بیمه‌شده برابر با ۳۸۴ هزار تومان در ماه بوده است. اکثر بیمه‌شدگان اختیاری بوده که بر حسب آمار موجود از بیمه‌های اختیاری تهران استنباط شده (حدود ۶۰ درصد) و متوسط سابقه بیمه‌ای آن‌ها قبل از اولین تقاضا به طور متوسط حدود ۷/۲ سال بوده است. یافته‌های این پژوهش همچنین نشان می‌دهند که اولاً از مجموعه این بیمه‌شدگان حدود ۱۷ درصد از آنان سابقه قطع قرارداد خود را داشته‌اند. البته بیشتر فقط یک بار قطع قرارداد داشته‌اند اما پیش‌بینی می‌شود که حدود ۴۰ درصد به احتمال زیاد در دو سال آینده بیکار شوند و ۳۵ درصد از آن‌ها هم به احتمال زیاد قراردادشان در آینده به دلیل مشکلات مالی به طور دائم یا موقتی قطع خواهد شد. دوم این که آگاهی آنان از مسائل بیمه‌ای به شدت پایین بوده به گونه‌ای که حتی عده‌ای نوع پوشش بیمه‌نامه خود را نمی‌دانسته‌اند. و سوم این که انگیزه‌ها و اهداف اصلی آنان از بیمه به ترتیب بازنشستگی، درمان، از کارافتادگی و فوت بوده است.

ملا محمدی و مستوفی (۱۳۹۳) به بررسی عوامل موثر بر موفقیت سازمان تامین اجتماعی در پوشش و برقراری مستمری بازنشستگی پرداخته‌اند. این مطالعه با استفاده از داده‌های یک نمونه ۱۹۱ نفری که بر اساس فرمول کوکران به طور تصادفی به روش پیمایشی و پرسشنامه از مجموعه ۸۱۲۳ نفری در بازه زمانی آبان ۱۳۹۱ تا شهریور ۱۳۹۳ از یک شعبه سازمان تامین اجتماعی در قم جمع‌آوری شده نشان می‌دهد که ۸۵/۷ درصد از نمونه بررسی شده بیش از ۴۱ سال سابقه کار داشته و ۹۵/۳ درصد از نمونه‌های تحقیق را مردان متأهل تشکیل می‌دهند. همین‌طور ۸۰/۲ درصد از نمونه‌ها بیش از دو نفر را تحت کفالت خود دارند. ۷۲/۳ درصد سن بین ۶۱ تا ۶۵ سال دارند و بر همین اساس ۶۷ درصد از آنان سواد ابتدایی به پایین را

دارند. همچنین ۶۷ درصد نمونه سابقه بیمه کمتر از ۱۵ سال دارند. طبق نتایج به دست آمده، عواملی از جمله برنامه‌های حمایت از کارگران، رفتار بیمه‌ای کارفرما، آموزش و اطلاع رسانی قوانین سازمان تأمین اجتماعی و در نهایت توان اقتصادی بیمه‌شده و کارفرما از مهم ترین عوامل مؤثر در برقراری مستمری بازنشستگی بوده‌اند.

نجفی (۱۳۹۸) به بررسی و پیش‌بینی وفاداری و رویگردانی بیمه‌شدگان خویش‌فرمای سازمان تأمین اجتماعی پرداخته و هدف خود را شناسایی میزان رویگردانی (قطع رابطه) بیمه‌شدگان و همچنین یافتن مدلی بر مبنای الگوریتم‌های هوش مصنوعی<sup>۱</sup> قرار می‌دهد که میزان رویگردانی بیمه‌شدگان در سال‌های آینده را به درستی تخمین بزند. در این پژوهش اطلاعات و شاخص‌های مهم حدود ۲۱۴۰۷ نفر در قالب ۲۷ ویژگی از بانک‌های اطلاعاتی سازمان تأمین اجتماعی استخراج گردیده است. سپس با استفاده از الگوریتم ژنتیک رتبه‌بندی نامغلوب<sup>۲</sup>، تعداد ۷ ویژگی مهم که حداقل خطای طبقه‌بندی را داشته‌اند انتخاب شدند (ویژگی‌های جنسیت، نوع بیمه، تعداد روز سابقه ناشی از قرارداد، سال تولد بیمه‌شده، تاریخ قرارداد، تاریخ تولد بزرگترین فرزند، سال شروع قرارداد). در این مطالعه در مرحله انتخاب بهترین طبقه‌بندی کننده برای انجام پیش‌بینی‌های مورد نیاز نیز از سه الگوریتم شبکه عصبی چندلایه<sup>۳</sup>، الگوریتم ماشین بردار پشتیبان<sup>۴</sup> و الگوریتم کی-میانگین<sup>۵</sup> استفاده شده که نهایتاً شبکه عصبی چندلایه، بهترین دقت طبقه‌بندی را با مقدار ۹۶/۸ درصد به دست آورد. سپس برای شبکه عصبی مذکور با استفاده از ۷ ویژگی مربوط به داده‌های سال‌های ۱۳۶۷ تا ۱۳۹۵ آموزش انجام شد. شبکه عصبی آموزش دیده برای پیش‌بینی وفاداری و رویگردانی مشتریان سال‌های ۱۳۹۶ و ۱۳۹۷ که به تعداد ۸۳۶۴ رکورد می‌باشد، مورد استفاده قرار گرفت. با

---

1. Artificial Intelligence  
2. Nsga-li  
3. Multilayer Perceptron  
4. Support Vector Machines  
5. K-Nearest Neighbor Algorithm

توجه به نتایج به دست آمده، حدود ۲۷ درصد از بیمه‌شدگان سال‌های ۱۳۹۶ و ۱۳۹۷ در کلاس رویگردان دسته‌بندی می‌شوند.

برووفر و همکاران (۱۳۹۵) در مقاله‌ای با عنوان «شناسایی الگوی رفتاری مشتریان در بیمه عمر و تشکیل سرمایه با استفاده از داده کاوی» به ارائه یک مدل مناسب و کارآمد جهت بخش‌بندی مشتریان بر اساس برخی از مهم‌ترین ویژگی‌های مالی و جمعیت شناختی در قالب عوامل موثر بر شاخص‌های ارزش دوره عمر مشتری (آ.اف.ام)<sup>۱</sup> پرداختند. پس از تعیین مقادیر شاخص‌های مدل آ.اف.ام شامل تازگی مبادله<sup>۲</sup>، تعداد دفعات مبادله<sup>۳</sup> و ارزش پولی<sup>۴</sup> مبادله در ۱۸۰ هزار مشتری و وزن‌دهی آن‌ها با استفاده از فرآیند تحلیل سلسله مراتبی، تعداد خوشه بهینه بر اساس شاخص سیلوئت<sup>۵</sup> و نرخ تاثیر شاخص‌های آ.اف.ام با استفاده از الگوریتم دو مرحله‌ای<sup>۶</sup> تعیین شد و در مرحله بعد به خوشه‌بندی مشتریان با استفاده از روش کی-میانگین پرداخته شده است. آمار مورد مطالعه شامل اطلاعات خرید بیمه‌نامه عمر و تشکیل سرمایه مشتریان شرکت بیمه سامان در سال‌های ۱۳۹۰ تا ۱۳۹۳ می‌باشد که فیلدهای مختلفی از جمله مشخصات بیمه‌شده (تاریخ تولد، جنسیت، وضعیت تاهل، شغل، درآمد، میزان تحصیلات)، مشخصات بیمه‌نامه (حق بیمه سالیانه، سرمایه فوت، پوشش‌های تکمیلی، میزان خطرپذیری، روش پرداخت) و پوشش‌های تکمیلی (فوت بر اثر سانحه، از کار افتادگی، امراض صعب‌العلاج) را شامل می‌شود. در این تحقیق از روش‌شناسی کریپس استفاده شده است. پس از بکارگیری تکنیک‌های مختلف جهت آماده‌سازی داده‌ها (حذف داده‌های پرت و گم‌شده) در نهایت ۱۷۱۱۹۲ رکورد جهت خوشه‌بندی مورد استفاده قرار گرفت. پس از جلسات متعدد با خبرگان ۵ ویژگی موثر که اهمیت زیادی در تحقیق دارند

۱. ارزش دوره عمر مشتری (Recency Frequency Monetary) ارزشی است که مشتری در طول عمر خود برای سازمان ایجاد می‌کند. این

مفهوم علاوه بر ارزش فعلی مشتری به ارزش بالقوه و ارزش آتی مشتری نیز اشاره دارد.

2. Recency
3. Frequency
4. Monetary
5. Silhouette
6. Two-step

انتخاب شدند که شامل: نسبت سرمایه اولیه به حق بیمه اولیه، نسبت سرمایه فوت حادثه به حق بیمه اولیه، نسبت سرمایه از کار افتادگی حادثه به حق بیمه اولیه، نسبت سرمایه امراض به حق بیمه اولیه و ماه شروع بیمه‌نامه می‌باشد. این مطالعه سپس با مقایسه دو الگوریتم خوشه‌بندی کی-میانگین و دو مرحله‌ای نتیجه می‌گیرد که الگوریتم کی-میانگین از نظر معیار شاخص سیلوئت با مقدار  $0/70$  برتری دارد که در نهایت این الگوریتم به عنوان بهترین الگوریتم برای خوشه‌بندی مشتریان انتخاب شد.

پرمه و همکاران (۱۳۹۹) در مقاله‌ای با عنوان؛ متغیرهای کلان اقتصادی و تقاضای بیمه‌های خویش‌فرمایی در سازمان تامین اجتماعی، به تاثیر متغیرهای کلان اقتصادی بر تقاضای بیمه‌های خویش فرما در ۳۱ استان کشور پرداخته‌اند. هدف خاص این مطالعه بررسی تاثیر تولید ناخالص داخلی، نرخ تورم، نرخ بیکاری و تحصیلات بر میزان تقاضای بیمه‌های اختیاری و حرف و مشاغل آزاد تامین اجتماعی بوده است. بدین منظور داده‌های مربوط به متغیرهای تحقیق در ۳۱ استان کشور در دوره زمانی ۱۳۹۵-۱۳۸۰ جمع‌آوری و از الگوی پویای پانلی به روش گشتاورهای تعمیم یافته (GMM) برای بررسی موضوع استفاده شده است. نتایج تحقیق بیانگر آن است که تاثیر نرخ تورم بر تقاضای بیمه‌های خویش‌فرمایی به صورت منفی و معنادار و متغیرهای نرخ بیکاری، تحصیلات و تولید ناخالص داخلی مثبت و معنادار بوده است. عبدی و همکاران (۲۰۱۷) یک رویکرد داده‌کاوی سه مرحله‌ای را برای شناسایی مشتریان وفادار و برنامه‌ریزی فروش بیمه بر اساس ویژگی مشتریان در ایران را توسعه داده‌اند. مجموعه داده‌ای که در این مطالعه استفاده شده شامل ۴۷۸ مشتری بیمه می‌باشد. در مجموع ۷۰ مشتری یک پوشش بیمه خاص را خریداری و ۴۰۸ مشتری این نوع پوشش بیمه‌ای را خریداری نمی‌کنند که متغیرهای مورد استفاده در این مدل شامل سن، تحصیلات، شغل، سابقه کار، درآمد، تعداد اعضای خانواده و ... می‌باشد. این مطالعه از الگوریتم کی-نزدیک‌ترین همسایه<sup>۱</sup> برای خوشه‌بندی<sup>۲</sup> مشتریان بهره برده است. هدف این مطالعه پیش‌بینی

1. K-Nearest Neighbor  
2. Clustering

این بوده که آیا یک پوشش بیمه‌ای خاص توسط مشتریان خریداری می‌شود یا خیر و همچنین در این مقاله تلاش شده تا از تکنیک‌های داده‌کاوی برای بهبود مدیریت ارتباط با مشتری و شناسایی مشتریان ارزشمند استفاده شود.

دامغانی و همکاران (۲۰۱۹) با استفاده از تکنیک‌های داده‌کاوی مدلی ترکیبی از طبقه‌بندی و خوشه‌بندی ساخته‌اند که بیمه‌نامه مناسب هر مشتری را بر اساس سابقه پوشش بیمه‌ای آن فرد پیشنهاد می‌دهد. در این مقاله یک مطالعه موردی واقعی شامل ۴۷۶ مشتری یک شرکت بیمه در ایران که شامل ۲۵ متغیر می‌باشد، مورد تجزیه و تحلیل قرار گرفته است که از این ۲۵ متغیر ۱۵ متغیر مربوط به ویژگی‌های مشتری و ۱۰ متغیر مرتبط با پوشش بیمه‌ای مشتریان می‌باشد. متغیرها شامل سن، تحصیلات، میزان درآمد، تعداد اعضای خانواده و ... می‌باشد. در این مطالعه در نهایت الگوریتم خوشه‌بندی کی-میانگین انتخاب شد که طبق آن مشتریان به چهار خوشه تقسیم شده‌اند. برای طبقه‌بندی نیز در این مطالعه از الگوریتم کی-نزدیک-ترین همسایه استفاده شد که در آن معیار دقت برابر با ۸۴/۶۵، پوشش برابر با ۷۷/۶۰ و صحت نیز برابر با ۸۰/۸۲ رسید.

رحمان<sup>۱</sup> و همکاران (۲۰۱۷) با تجزیه و تحلیل داده‌های شرکت بیمه زندگی پرایم اسلامی بنگلادش<sup>۲</sup> در بازه زمانی ۲۰۱۱ تا ۲۰۱۴ و با استفاده از اطلاعات ۲۸۲۲۸۲ بیمه‌گذار به نحوه واکنش مشتریان به بیمه‌نامه‌های ارائه شده پرداخته‌اند. در این مطالعه از ۱۰ ویژگی (شرایط خط‌مشی<sup>۳</sup>، سن، جنسیت، شغل، شهری یا روستایی، وضعیت تاهل، مبلغ تضمین شده<sup>۴</sup>، بخش محل سکونت، حالت پرداخت حق بیمه سالانه یا ماهانه و منظم بودن در پرداخت) استفاده شد. هدف این مطالعه طبقه‌بندی مشتریان بر اساس ویژگی‌های آن‌ها بوده تا بتوانند برچسب کلاس را برای مشتریان آینده پیش‌بینی کنند. این هدف بر مبنای منظم بودن یا نبودن مشتری

1. Rahman

2. Prime Islamic Life Insurance Company Ltd., of Bangladesh

۳. خط مشی (Policy term) به بازه زمانی اشاره دارد که دارنده بیمه نامه از طریق آن حق بیمه خود را پرداخت می‌کند.

۴. مبلغ تضمین شده (Sum-Assured) مبلغ از پیش تعیین شده‌ای که بیمه‌گر متعهد به پرداخت آن در صورت فوت بیمه‌گذار می‌شود.



با استفاده از یک ویژگی به نام «نظم و قاعده»<sup>۱</sup> مشخص می‌شود. از تکنیک‌های مختلف ارزیابی ویژگی برای تعیین موثرترین ویژگی‌های مورد نیاز طبقه‌بندی استفاده شده که اولین تکنیک «ارزیابی ویژگی کسب اطلاعات»<sup>۲</sup> و دیگری «نسبت افزایش تکنیک‌های ارزیابی ویژگی»<sup>۳</sup> می‌باشد. در نهایت از میان تمام ویژگی‌های انتخاب شده از دو تکنیک، ویژگی حالت پرداخت حق‌بیمه به عنوان موثرترین ویژگی شناخته شده است.

ماتو<sup>۴</sup> و همکاران (۲۰۱۸) به بررسی ارزش داده‌های غنی شده مشتریان برای مدیریت تحلیلی ارتباط با مشتری در بخش بیمه خودرو و خانه پرداخته‌اند. برای این منظور اعلام‌های آنلاین یک بیمه‌گر سوئسی که دارای محصولات بیمه زندگی و غیرزندگی است را از سال ۲۰۱۲ تا ۲۰۱۵ مورد بررسی قرار داده‌اند. این داده‌ها شامل متغیرهای کمی شخصی (روز تولد، جنسیت و غیره) و متغیرهای عمومی بیمه‌نامه (شناسه بیمه‌نامه، نسخه بیمه‌نامه، تاریخ آغاز و تاریخ انقضا) و قیمت‌گذاری متغیرهای کمکی مربوط به محصول بیمه خاص (ارزش خانوار بیمه‌شده برای بیمه‌نامه خانگی) بوده است. از مجموع ۲/۵ تا ۳ میلیون مشاهده تقریباً ۲۷۵ هزار مورد قیمت بیمه اتومبیل و ۹۰ هزار مورد قیمت بیمه خانه به عنوان نمونه انتخاب شده است. برای محصول بیمه اتومبیل متغیرهای کمی مانند تاریخ تولد، کدپستی محل اقامت، جنسیت، تاریخ صدور گواهی‌نامه رانندگی و سیله نقلیه نیز ثبت شده است. در ادامه برای پیش‌بینی خرید و طبقه‌بندی مشتریان از مدل پیش‌بینی جنگل تصادفی استفاده شده که بر اساس دو معیار دقت<sup>۵</sup> و امتیاز-F<sup>۶</sup> دقت پیش‌بینی آزمون مربوطه بر روی داده‌های تست به ۰/۸۴۴ رسیده است.

---

1. Regularity  
2. Information Gain Attribute Evaluation  
3. Gain Ratio Attribute Evaluation Techniques  
4. Mau  
5. Accuracy  
6. F-Score

عبدالرحمن<sup>۱</sup> و همکاران (۲۰۲۱) در مطالعه‌ای با استفاده از خوشه‌بندی k-modes و طبقه‌بندی درخت تصمیم به تقسیم‌بندی<sup>۲</sup> مشتری و پروفایل سازی<sup>۳</sup> برای بیمه عمر پرداخته‌اند. الگوریتم درخت تصمیم برای طبقه‌بندی پروفایل مشتری با مقایسه دو شاخص آنتروپی و جینی استفاده شد و همچنین خوشه‌بندی k-modes مشتریان را به سه گروه تقسیم‌بندی کرد؛ «مشتریان بالقوه با ارزش بالا»، «مشتریان کم ارزش» و «مشتریان بی‌علاقه». در نهایت درخت تصمیم با شاخص جینی و با اعتبار سنجی متقابل ۱۰ برابری<sup>۴</sup> به عنوان بهترین مدل برازش با دقت ۸۱/۳ درصد تعیین شد. آتزون و همکاران (۲۰۲۲) با استفاده از الگوریتم جنگل تصادفی و رگرسیون لجستیک به پیش‌بینی مدت زمان سپری شده بیمه‌نامه‌های عمر می‌پردازند. نتایج بدست آمده نشان می‌دهند که عوامل غیراقتصادی (مدت زمان مانده تا سررسید، شرکت بیمه و رویکرد بازاریابی آن) در تعیین مدت زمان سپری شده نقش مهمی دارند و در مقابل عوامل اقتصادی و مالی بی‌تأثیر بوده‌اند (البته به غیر از نرخ رشد درآمد که تأثیر زیادی داشت). این مطالعه همچنین نشان می‌دهد که مدل‌های خطی مانند رگرسیون لجستیک در مدل‌سازی تصمیمات مالی پیچیده مانند تقاضای بیمه‌نامه‌ها عملکرد مناسبی ندارند.

کونگ و همکاران (۲۰۲۲) با طراحی یک فیلتر همکارانه (CF) با استفاده از الگوریتم‌های یادگیری ماشین به پیش‌بینی رفتار خریداران بیمه‌نامه‌های زندگی در کره جنوبی می‌پردازد. این مقاله سپس عملکرد این الگوریتم‌ها را با مدل رگرسیون لجستیک مقایسه می‌کند و نشان می‌دهد مدل‌های مبتنی بر یادگیری ماشین و فیلترهای همکارانه (CF) مبتنی بر آن‌ها در پیش‌بینی و توصیه‌های سیاستی برای طبقه‌بندی مشتریان بسیار مناسب هستند به ویژه زمانی که مشخصه‌های قراردادی و اطلاعات خرید بیمه‌نامه نیز به مشخصه‌های خریداران افزوده شود. سیورینو و پنگ<sup>۵</sup>

---

1. Abdul-Rahman  
2. Segmentation  
3. Profiling  
4. 10-Fold Cross Validation  
5. Severino & Peng

(۲۰۲۱) در مطالعه‌ای تحت عنوان الگوریتم‌های یادگیری ماشین برای پیش‌بینی تقلب در بیمه اموال با استفاده از شواهد تجربی و داده‌های دنیای واقعی، به پیش‌بینی تقلب در ادعاهای بیمه دارایی در یک شرکت بزرگ بیمه‌ای در برزیل پرداخته‌اند، که مجموعه داده‌های آن مربوط به خسارت‌های ثبت شده برای بیمه‌های مسکونی و تجاری از سال ۲۰۰۹ تا ۲۰۱۸ می‌باشد و متغیرهای آن شامل روزهای بین پایان قرارداد و ادعا<sup>۱</sup>، روزهای بین شروع قرارداد و ادعا<sup>۲</sup>، زمان تایید<sup>۳</sup>، مبلغ بیمه شده، حق بیمه، تعداد اقساط، سن و ادعای قبلی می‌باشد. در این مطالعه ۹ مدل به صورت بازگشتی آزمایش شدند و میانگین نتایج پیش‌بینی کننده با کنترل مثبت کاذب<sup>۴</sup> و منفی کاذب<sup>۵</sup> مقایسه شدند. نتایج با استفاده از معیارهای مختلف ارزیابی شدند و مشخص شد که روش‌های مبتنی بر مجموعه<sup>۶</sup> (جنگل تصادفی<sup>۷</sup> و تقویت گرادیان<sup>۸</sup>) و شبکه‌های عصبی عمیق<sup>۹</sup> بهترین نتایج را به همراه داشته‌اند و عملکرد متوسط بالاتری را در مقایسه با سایر طبقه‌بندی کننده‌ها از جمله رگرسیون لجستیک<sup>۱۰</sup> رایج نشان می‌دهند.

به‌طور کلی، مرور مطالعات مربوطه در ایران نشان می‌دهد که عمده پژوهش‌ها در خصوص شناسایی ویژگی‌ها و انگیزه‌های مشتری در صنعت بیمه مربوط به بیمه‌های بازرگانی است و تعداد مطالعاتی که به شناسایی ویژگی‌های بیمه‌شدگان (به‌نوعی مشتریان) اختیاری بیمه‌های اجتماعی از طریق الگوریتم‌های یادگیری ماشینی پرداخته‌اند بسیار کم است. از طرفی، معدود مطالعات انجام شده نشان می‌دهند که انعقاد انواع بیمه‌نامه بازنشستگی یا قطع قراردادهای بیمه می‌تواند به دلایل مختلفی اتفاق بیفتد که بررسی آن‌ها برای درک چرایی پایین بودن ضریب نفوذ بیمه‌های حرف و مشاغل آزاد در ایران ضروری است. اما درک

1. Days Between Contract End And Claim
2. Days Between Contract Start And Claim
3. Time Of Approval
4. False Positive
5. False Negative
6. Ensemble
7. Random Forest
8. Gradient Boosting
9. Deep Neural Networks
10. Logistic Regression

نیازها و تقاضای مشتریان و شناخت ویژگی های آن ها به سادگی امکان پذیر نیست و ناگزیر بسیاری از شرکت های بازنشستگی و بیمه عمر به استفاده از ابزارهای جدید مانند یادگیری ماشین روی آورده اند. الگوریتم های یادگیری ماشینی به راحتی و بدون مداخله مستقیم انسانی در روند کار می توانند با حجم زیادی از داده های چند بعدی و چند متغیر کار کرده و الگوها و روندهای موجود را شناسایی کنند که باعث افزایش دقت و سرعت انجام کار می شود. مطالعات زیادی از این روش های جدید در حوزه بیمه و شناخت ویژگی های مشتریان آن ها استفاده کرده اند و توانسته اند به نتایجی دست یابند که در ارائه محصولات جدید بیمه ای متناسب با نیاز مشتریان مفید واقع گردند.

### ۳. روش شناسی و توصیف داده ها

داده کاوی یکی از مهم ترین ابزارها در جهت شناسایی الگوهای رفتاری مشتری می باشد. بر همین اساس، در این مقاله از روش کریپس که یکی از معروف ترین روش های داده کاوی است استفاده شده است. این روش اولین بار توسط چمپان<sup>۱</sup> و همکاران (۲۰۰۰) مطرح شد که نقطه عطفی در روش های داده کاوی و مدل های فرآیند بود چراکه یک نمای کلی از چرخه عمر یک پروژه داده کاوی شامل شش مرحله را در اختیار قرار می دهد (ویرث و هیپ<sup>۲</sup>، ۲۰۰۰).

- مرحله اول، درک کسب و کار<sup>۳</sup>، متمرکز است بر فهم اهداف و الزامات پروژه از دیدگاه تجاری، تبدیل این دانش به تعریف مسئله داده کاوی و سپس طرح اولیه پروژه که برای دستیابی به اهداف طراحی شده است. در واقع در این مرحله تلاش بر این است که مشکل و مسئله مورد نظر در فضای کسب و کار تامین اجتماعی تعیین و مشخص شود.
- مرحله دوم، درک داده ها<sup>۴</sup>، با جمع آوری داده های اولیه شروع می شود و با فعالیتهایی به منظور آشنایی با داده ها، شناسایی مشکلات مربوط به کیفیت داده ها، کشف اولین بینش در داده ها یا شناسایی زیرمجموعه های جالب برای تشکیل فرضیه هایی به جهت پی بردن به

1. Champan et al

2. Wirth & Hipp

3. Business Understanding

4. Data Understanding

اطلاعات پنهان ادامه می‌یابد. بین درک کسب و کار و درک داده ارتباط نزدیکی وجود دارد، چون درک مسئله داده کاوی حداقل نیاز به درک کمی از داده‌های موجود دارد.

- مرحله سوم، آماده‌سازی داده‌ها<sup>۱</sup> است که همه فعالیت‌ها برای تبدیل داده‌های خام اولیه به مجموعه داده‌های نهایی که در مدل وارد می‌شوند را در بر می‌گیرد. برای این منظور متغیرهایی که قرار است تحلیل شوند و یا برای تحلیل مناسب هستند انتخاب می‌شوند (انتخاب ویژگی<sup>۲</sup>). این مرحله کمی‌سازی متغیرها مثلاً تبدیل متغیر ترتیبی<sup>۳</sup> به عددی<sup>۴</sup> و همچنین حذف داده‌های پرت و نامرتب از مجموعه داده را در بر می‌گیرد.
- مرحله چهارم مدل‌سازی<sup>۵</sup> است. در این مرحله تکنیک‌های مدل‌سازی مختلف انتخاب و اعمال می‌شوند و پارامترهای آن‌ها تا مقادیر بهینه کالیبره<sup>۶</sup> می‌شوند، به عنوان مثال با تغییر نسبت تعداد داده‌های مجموعه تست<sup>۷</sup> و داده‌های مجموعه آموزش<sup>۸</sup> و یا تغییر پارامترهای اصلی<sup>۹</sup> الگوریتم که تعیین آن‌ها با توجه به شرایط به دست خیره یا کاربر صورت می‌گیرد. به طور معمول، برای یک نوع مسئله داده کاوی چندین تکنیک وجود دارد. برخی از تکنیک‌ها به فرمت‌های داده خاصی نیاز دارند. همچنین ارتباط نزدیکی بین آماده‌سازی داده و مدل‌سازی وجود دارد. در اغلب موارد هنگام مدل‌سازی است که شخص متوجه مشکلات داده‌ها می‌شود یا ایده‌هایی برای ساخت داده‌های جدید به دست می‌آورد.
- مرحله پنجم اعتبارسنجی<sup>۱۰</sup> است. در واقع، مدلی که ساخته شده ممکن است از دیدگاه تحلیل داده‌ها کیفیت بالایی داشته یا نداشته باشد و لذا پیش از اقدام به استقرار نهایی

---

1. Data Preparation  
2. Feature Selection  
3. Ordinal  
4. Nominal  
5. Modeling  
6. Calibrated  
7. Test Set  
8. Train Set  
9. Hyper Parameter  
10. Evaluation

مدل، ارزیابی دقیق مدل و بررسی مراحل اجرا شده در ساخت مدل برای اطمینان از دستیابی مناسب به اهداف تجاری ضروری است.

- مرحله ششم، پیاده‌سازی مدل است. معمولاً دانش به دست آمده باید به گونه‌ای سازماندهی و ارائه شود که سازمان بتواند از آن استفاده کند. در بسیاری از موارد این کاربر، و نه تحلیلگر داده خواهد بود، که مراحل استقرار را انجام خواهد داد.

مرحله اول یعنی درک کسب و کار و شناخت مسئله در بخش‌های پیشین مقاله با جزئیات بررسی شد. در این بخش و در ادامه به درک داده‌ها، آماده‌سازی داده‌ها و معرفی روش مدلسازی پرداخته می‌شود. یافته‌های حاصل از مدلسازی و اعتبارسنجی آن در بخش چهارم مقاله ارائه خواهد شد.

### ۳-۱. توصیف و درک داده‌ها

برای پیش‌بینی خرید بیمه متناسب با خصوصیات مشتری و بررسی علت پایین بودن ضریب نفوذ بیمه‌نامه حرف و مشاغل آزاد می‌بایست مشخصات جمعیت شناختی را مورد بررسی قرار داد. بررسی پیشینه پژوهش نیز نشان داد که متغیرهای جمعیت شناختی از جمله عوامل موثر بر تقاضا یا پیش‌بینی خرید بیمه‌نامه‌ها هستند. در این پژوهش از مجموعه داده‌های مربوط به بیمه‌شدگان حرف و مشاغل آزاد سازمان تأمین اجتماعی در سال ۱۳۹۹ استفاده شده که شامل ۱۲۸۶۱۷۴ نفر بیمه شده است. جدول ۱ آمار توصیفی مربوط به متغیرها و مشخصه‌های اصلی مورد مطالعه را نمایش می‌دهد.

جدول ۱. مشخصات جمعیت شناختی بیمه شده‌های حرف و مشاغل آزاد سازمان تأمین اجتماعی

مشخصات جمعیت شناختی	جنسیت	
	زن	۵۷۳۹۴۷ (۴۴/۶ درصد)
سن	مرد	۷۱۲۲۲۷ (۵۵/۴ درصد)
	کمترین	۱۴

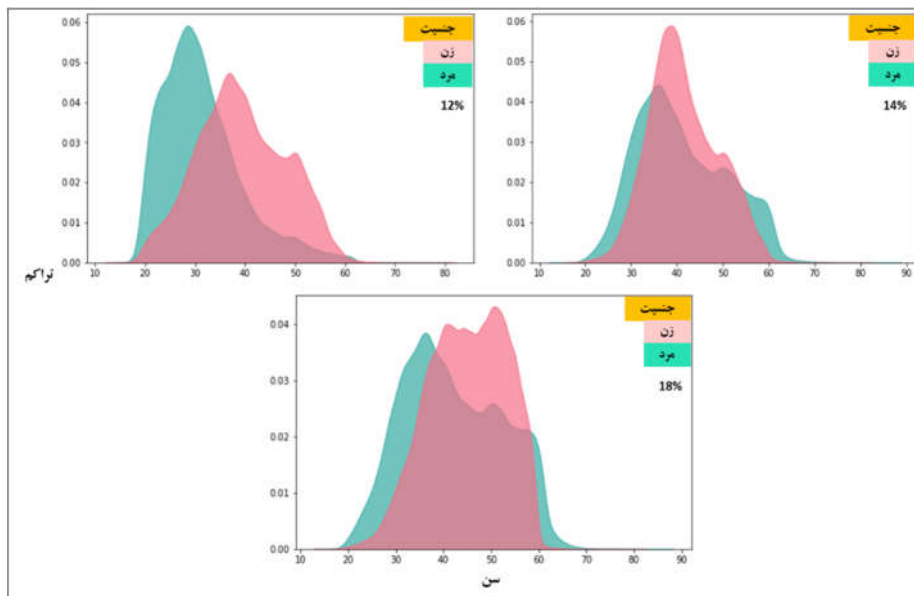
۴۰	متوسط	سابقه کار (روز)	
۸۶	بیشترین		
۱	کمترین		
۲۷۱۰	متوسط		
۱۶۱۲۵	بیشترین		
۱/۹۱	کمترین		
۱/۹۷	متوسط	میانگین درآمد ماهیانه (میلیون تومان)	
۱۸/۶	بیشترین		
۴۱۹۹۶۴	مستمری بازنشستگی (۱۲ درصد)		
۱۸۲۲۵۵	مستمری بازنشستگی و از کارافتادگی (۱۴ درصد)	نوع بیمه‌نامه	
۶۸۳۹۵۵	مستمری بازنشستگی، از کارافتادگی و بازماندگان (۱۸ درصد)		

منبع: سازمان تامین اجتماعی (۱۴۰۰)

هدف این مطالعه بررسی و پیش‌بینی خرید بیمه حرف و مشاغل آزاد سازمان تأمین اجتماعی است. قرار است مدلی ساخته شود که بتواند بر اساس چهار ویژگی سن، جنسیت، سابقه و متوسط درآمد افراد را در دو طبقه متقاضی بیمه‌نامه با مستمری بازنشستگی (حق بیمه ۱۲ درصد) یا متقاضی بیمه‌نامه پوشش دهنده ریسک از کارافتادگی و فوت علاوه بر مستمری بازنشستگی (حق بیمه ۱۴ و ۱۸ درصد) طبقه‌بندی کند. به عبارت دیگر، هدف این است که بر اساس چهار ویژگی در دسترس ثبت شده<sup>۱</sup> (جنسیت، سن، سابقه کار و میانگین درآمد بیمه‌شده) به عنوان متغیرهای ورودی مدل، پیش‌بینی شود که افراد کدام‌یک از بیمه‌نامه‌ها را انتخاب خواهند کرد. ابتدا به بررسی ارتباط متغیر هدف یعنی نوع بیمه‌نامه خریداری شده با متغیرهای توضیحی یا ورودی به تفکیک پرداخته خواهد شد که قطعاً موجب روشن‌تر شدن کار پیش از ورود به یافته‌های پژوهش خواهد شد.

۱. متأسفانه ویژگی‌ها و اطلاعات بیشتری از بیمه‌شدگان حرف و مشاغل آزاد در دسترس محققان قرار نگرفته است.

یکی از متغیرهای ورودی مهم سن بیمه‌شده است. نمودار ۲ توزیع سنی زنان و مردان را به تفکیک نوع بیمه‌نامه نمایش می‌دهد. نکته قابل توجه در این نمودار مربوط به بیمه ۱۲ درصد پوشش دهنده بازنشستگی است که متمایز از توزیع سنی در سایر انواع بیمه‌نامه‌های ۱۴ و ۱۸ درصد است. همانطور که می‌توان دید در میان بیمه‌شدگان مرد، فقط جوانان هستند که این نوع بیمه ۱۲ درصدی را تقاضا می‌کنند در حالی که مردان در سنین بالاتر خریدار بیمه‌نامه‌های ۱۴ و ۱۸ درصدی هستند که علاوه بر بازنشستگی، ریسک فوت و ازکارافتادگی را نیز پوشش می‌دهد. در نمودار می‌توان دید که برای زنان اما این موضوع برقرار نیست و زنان در همه سنین متقاضی بیمه ۱۲ درصدی پوشش دهنده بازنشستگی هستند و توزیع سن زنان در بیمه‌نامه ۱۲ درصدی متقارن است.

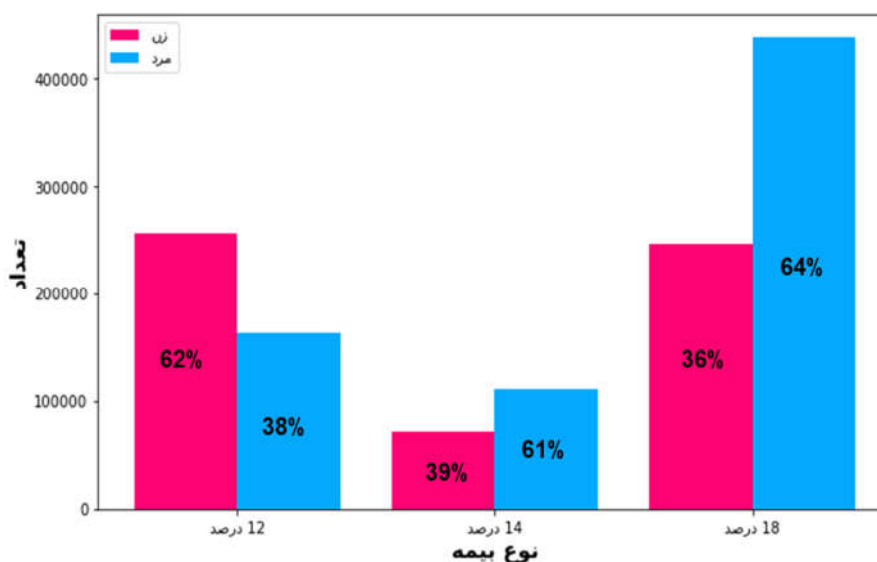


نمودار ۱. وضعیت توزیع سنی جمعیت زنان و مردان در بیمه‌های مختلف حرف و مشاغل آزاد

منبع: یافته‌های پژوهش



همانطور که در بررسی سن نیز دیده شد، متغیر ورودی و توضیحی جنسیت بیمه‌شدگان بسیار اهمیت دارد و از این رو نگاهی به آمار توصیفی خرید بیمه‌نامه‌ها بر اساس جنسیت می‌تواند مفید باشد. برای فهم بهتر می‌توان به سهم مردان و زنان از هر نوع بیمه‌نامه ۱۲ درصد (بازنشستگی)، ۱۴ درصد (بازنشستگی و ازکارافتادگی) و ۱۸ درصد (بازنشستگی، ازکارافتادگی و فوت) که به تفکیک در نمودار ۳ نشان داده شده توجه کرد.



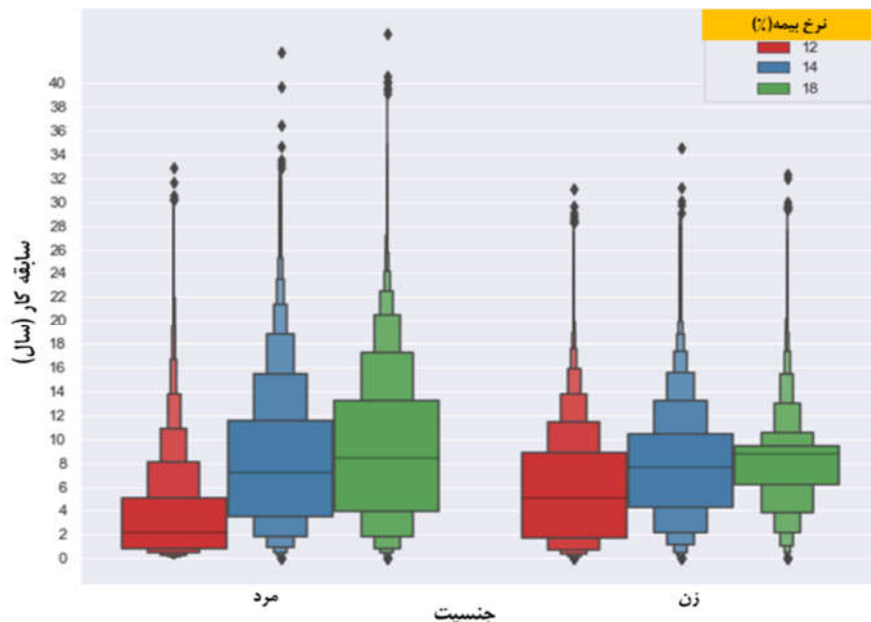
نمودار ۳. سهم مردان و زنان در هر یک از بیمه‌های حرف و مشاغل آزاد

منبع: یافته‌های پژوهش

همانطور که در نمودار ۳ می‌توان دید، مشابه حالت کلی همچنان سهم مردان در بیمه‌نامه‌های ۱۴ و ۱۸ درصدی بیشتر و حدوداً دو برابر زنان است. در حالی که به‌وضوح می‌توان دید این موضوع در بیمه‌نامه‌های با نرخ ۱۲ درصد برعکس است و سهم زنان حدود ۲ برابر مردان است. این نمودار ساده به‌خوبی نشان می‌دهد که زنان از بیمه ۱۲ درصد و مردان از بیمه‌های ۱۴ و ۱۸ درصد استقبال می‌کنند. یک دلیل مهم این امر را باید در قوانین تکفل در ایران

جستجو کرد که تکفل فرزندان و زن را بر عهده شوهر یا پدر قرار داده است. بر این اساس، در صورت فوت بیمه‌شده مرد مستمری بازماندگان به همسر و فرزندانش تعلق می‌گیرد در حالی که در صورت فوت بیمه‌شده زن، در حالت عادی مستمری به همسر و فرزندان تعلق نخواهد گرفت. همین امر باعث شده تا زنان دلیل محکمی برای خرید بیمه‌نامه‌های پوشش دهنده ریسک از کارافتادگی و فوت نداشته باشند و با پرداخت حق بیمه کمتر، بیمه‌نامه ۱۲ درصدی را انتخاب کنند که صرفاً بازنشستگی را پوشش می‌دهد. در مقابل، مردان که عموماً افرادی (شامل همسر و فرزندان) را تحت تکفل خود دارند و غالباً با ریسک بیشتری برای از کارافتادگی مواجه هستند، ترجیح داده‌اند با پرداخت حق بیمه‌های بالاتر بیمه‌نامه‌هایی را انتخاب کنند که علاوه بر بازنشستگی، از کارافتادگی و فوت را نیز ارائه دهد.

یکی دیگر از متغیرهای توضیحی و مهم و تأثیرگذار سابقه کار است که در نمودار ۴ به تفکیک جنسیت و برای بیمه‌نامه‌های مختلف نشان داده شده است. می‌توان دید اکثر مردان با سابقه کاری بیشتر، تمایل به استفاده از بیمه‌نامه با نرخ ۱۸ درصد و پس از آن بیمه‌های با نرخ ۱۴ درصد را دارند. اما برای زنان کمی متفاوت است و ارتباط قوی بین سابقه کار و نوع بیمه‌نامه برای زنان وجود ندارد. می‌توان دید اکثر جمعیت زنان با سابقه کاری ۲ تا ۹ سال متقاضی بیمه‌نامه با نرخ ۱۲ درصد هستند و پس از آن، اکثر جمعیت با سابقه کاری ۴ تا ۱۰ سال متقاضی بیمه‌نامه با نرخ ۱۴ درصد هستند.

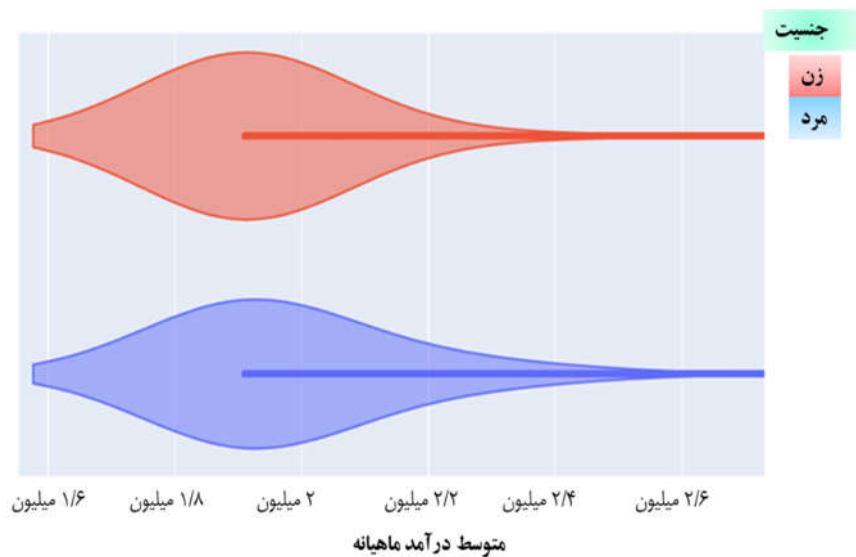


نمودار ۴. وضعیت سابقه کار نسبت به جنسیت در بیمه‌های مختلف حرف و مشاغل آزاد

منبع: یافته‌های پژوهش

یکی دیگر از ویژگی‌های ثابت شده از بیمه‌شدگان متوسط درآمد ماهانه آنهاست که می‌توان از آن برای توضیح و پیش‌بینی خرید بیمه‌نامه‌ها استفاده کرد. نمودار ۵ توزیع درآمد افراد بر اساس جنسیت آن‌ها را نشان می‌دهد. از این نمودار نیز می‌توان دریافت که اکثر افرادی که از بیمه‌های حرف و مشاغل آزاد استفاده می‌کنند در سال ۱۳۹۹ درآمدی بین ۱/۹ میلیون تومان تا ۲/۲ میلیون تومان داشته‌اند که بسیار نزدیک به حداقل دستمزد آن سال معادل ۱۹۱۰۴۲۷ تومان بوده است (خط ممتد در نمودار بازه درآمدی را نشان می‌دهد و در قسمت-هایی که نمودار برآمدگی دارد نشان دهنده وجود افراد بیشتر با آن درآمد است). این موضوع البته به احتمال زیاد به دلیل عدم گزارش صحیح درآمد توسط بیمه‌شده برای پرداخت حق بیمه کمتر است. طبق قوانین بازنشستگی در ایران تنها درآمد دو سال آخر در تعیین

مستمری بازنشستگی تأثیر گذار است و به همین دلیل افراد در طول دوران کاری (به جز دو سال پایانی) انگیزه دارند تا حداقل در آمد را گزارش کنند و تنها سوابق خود را تکمیل کنند.



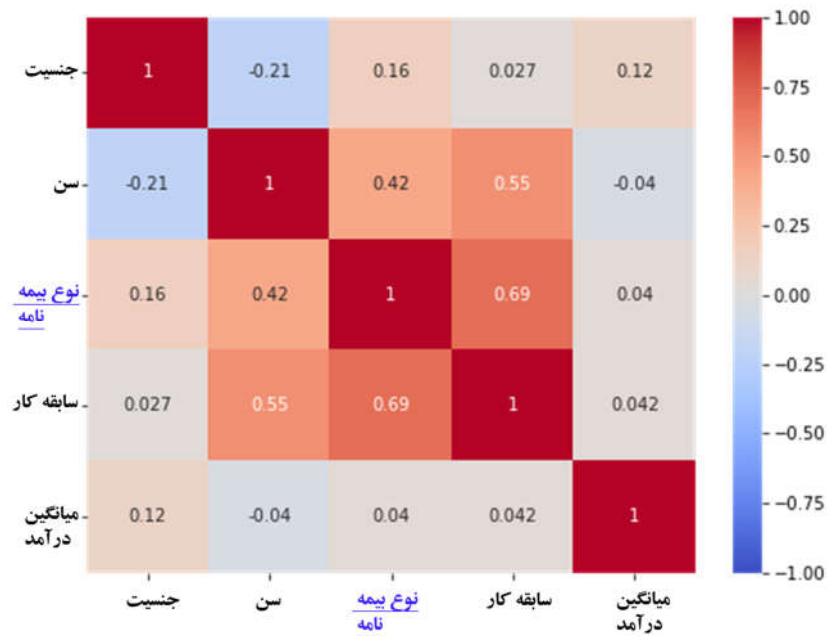
#### نمودار ۵. وضعیت توزیع درآمد مردان و زنان تحت بیمه حرف و مشاغل آزاد

منبع: یافته‌های پژوهش

#### ۲-۳. آماده‌سازی داده‌ها

آماده‌سازی داده‌ها و انتخاب ویژگی‌های مناسب شامل سه مرحله است که در هر قسمت عملیات‌های مختلفی بر روی داده‌ها انجام می‌شود. مرحله اول پیش پردازش داده‌هاست. در این مرحله متغیر جنسیت با جایگذاری مقدار ۰ برای زنان و مقدار ۱ برای مردان به مشخصه با مقادیر عددی تبدیل شد. یکی دیگر از اقدامات مهم در این مرحله شناسایی مقادیر گم شده و مقادیر تکراری است که با رصد تمامی داده‌ها و مشاهدات مشخص شد هیچ مقدار گم شده‌ای برای هیچکدام از مشخصه‌ها وجود ندارد. همچنین در پیش پردازش داده‌ها باید داده‌های پرت نیز مورد بررسی قرار گیرد که البته این مهم بعد از مرحله دوم از آماده‌سازی داده‌ها یعنی انتخاب موثرترین مشخصه‌ها انجام خواهد گرفت. برای اینکه مشخص شود

کدام یک از ویژگی‌های افراد (سن، جنسیت، درآمد و سابقه کار) در تمایل آن‌ها برای انتخاب نوع بیمه‌نامه موثرتر است، از میزان همبستگی بین متغیرهای ورودی با متغیر هدف یا نوع بیمه‌نامه انتخاب‌شده توسط افراد استفاده شده است. معمولاً هدف از انتخاب موثرترین ویژگی‌ها این است که تنها از ویژگی‌های مهم در ساخت مدل و الگوریتم استفاده شود که البته در اینجا این امر موضوعیت ندارد و همانطور که با مرور آمار توصیفی دیده شد هر چهار ویژگی بر انتخاب بیمه‌نامه مؤثرند و مورد استفاده قرار خواهند گرفت. در اینجا هدف از انتخاب موثرترین ویژگی‌ها صرفاً مرتب‌سازی داده‌هاست تا بتوان داده‌های پرت را حذف کرد. با توجه به نمودار حرارتی ۶ می‌توان میزان همبستگی بین مشخصه‌های بیمه‌شدگان و متغیر هدف یا نوع بیمه‌نامه (۱۲ درصدی و بیمه‌نامه‌های (۱۴ و ۱۸ درصدی) با پوشش ریسک بیشتر) را مشاهده کرد. میزان همبستگی بین متغیر «سابقه کار» با متغیر هدف «نوع بیمه‌نامه» تقریباً ۷۰ درصد است که بسیار بیشتر از همبستگی سایر متغیرهای ورودی با متغیر هدف می‌باشد. در نتیجه، از لحاظ آماری می‌توان گفت که سابقه کار افراد بیشترین تاثیر را نسبت به سایر متغیرهای ورودی (سن، جنسیت و میانگین درآمد) در انتخاب بیمه‌نامه دارد. البته تمامی متغیرهای دیگر از جمله جنسیت یا سن نیز تأثیرگذارند.



نمودار ۶. وضعیت همبستگی متغیرهای مورد بررسی با یکدیگر

منبع: یافته های پژوهش

از روش فیلتر که از معروف ترین روش ها در انتخاب ویژگی در مسائل طبقه بندی به شمار می رود نیز برای انتخاب مؤثرترین مشخصه استفاده شد. این روش رابطه ی بین متغیرهای ورودی و متغیر هدف را ارزیابی و نتیجه آن، انتخاب آن دسته از متغیرهای ورودی است که قوی ترین رابطه را با متغیر هدف دارند. در این روش برای انتخاب ویژگی از معیارهای آماری از جمله آزمون مربع کای  $\chi^2 = \frac{(Actual-Expected)^2}{Expected}$  بین متغیرهای ورودی برای امتیاز دادن استفاده می شود. طبق نتایج این روش نیز مشاهده شد که سابقه کار بیشترین امتیاز را دارد و در نتیجه تاثیر آن بر متغیر هدف به مراتب بیشتر از بقیه مشخصه ها می باشد. بر این اساس، مجموعه داده موجود بر اساس سابقه کار تقسیم بندی شد و پس از حذف داده های پرت ۷۲۶۵۸۲ رکورد باقی ماند که در مراحل بعدی و مدل سازی مورد

استفاده قرار می‌گیرد. مرحله سوم از آماده‌سازی داده‌ها، هم‌مقیاس‌سازی ویژگی‌ها است. منظور از هم‌مقیاس‌سازی این است که باید بازه تغییرات در هر ویژگی (متغیر) یکسان شود تا الگوریتم از واحد و مقیاس متغیر متأثر نشود. در اینجا از روش نرمالایز کردن مطابق رابطه 
$$y = \frac{x - \min}{\max - \min}$$
 استفاده می‌شود تا مقادیر همه متغیرها در محدوده جدید ۰ و ۱ قرار گیرند.

### ۳-۳. مدلسازی

برای مدلسازی در این مقاله از روش‌های مبتنی بر یادگیری ماشین استفاده می‌شود که کاربرد زیادی در حوزه‌های بیمه و بانکی پیدا کرده‌اند (متدین و همکاران ۱۴۰۰، متفکر آزاد و همکاران ۱۳۹۰). یادگیری ماشین شاخه‌ای از هوش مصنوعی<sup>۱</sup> و علوم کامپیوتر<sup>۲</sup> است که در آن کامپیوترها از مجموعه‌ای از الگوریتم‌ها استفاده می‌کنند تا الگو و روند داده‌های تاریخی را یاد بگیرند و بتوانند رفتار و مقادیر داده‌های آینده را پیش‌بینی کنند، در واقع رفتاری همانند انسان در یادگیری را دارند که به تدریج آن را بهبود می‌بخشند (هرویتز و کریش<sup>۳</sup>، ۲۰۱۸). یادگیری ماشین در حال حاضر یکی از حوزه‌های رو به رشد در علم داده<sup>۴</sup> است که با استفاده از روش‌های آماری، الگوریتم‌ها را برای طبقه‌بندی یا پیش‌بینی آموزش می‌دهد تا بینش‌هایی کلیدی در پروژه‌های داده‌کاوی را آشکار کنند. به طور کلی می‌توان یادگیری ماشین را به دو دسته یادگیری تحت نظارت<sup>۵</sup> و یادگیری بدون نظارت<sup>۶</sup> تقسیم بندی کرد. یادگیری تحت نظارت معمولاً با مجموعه‌ای از داده‌ها و درک خاصی از نحوه طبقه‌بندی آن داده‌ها آغاز می‌شود، و برای یافتن الگوهایی در داده‌ها در نظر گرفته شده است که می‌تواند در یک فرآیند تحلیلی اعمال شود. یادگیری بدون نظارت زمانی مناسب است که با حجم عظیمی از

---

1. Artificial Intelligence (Ai)  
 2. Computer Science  
 3. Hurwitz & Kirsch  
 4. Data Science  
 5. Supervised Learning  
 6. Unsupervised Learning

داده‌ها سروکار داشته باشیم که برچسب‌گذاری<sup>۱</sup> نشده باشند (ویژگی هدف مشخص نباشد) که در آن یک فرآیند تکراری از تجزیه و تحلیل داده‌ها بدون دخالت انسان انجام می‌شود. در این مطالعه با توجه به آمار قابل اطمینان که از دفتر آمار، اطلاعات و محاسبات سازمان تامین اجتماعی گرفته شده، برای مدل‌سازی از دو الگوریتم درخت تصمیم و جنگل تصادفی استفاده می‌شود که هر دو جزء الگوریتم‌های طبقه‌بندی از یادگیری ماشین می‌باشند. طبقه‌بندی زمانی استفاده می‌شود که داده‌ها بر اساس یک ویژگی هدف به گروه‌های مختلف طبقه‌بندی شوند و یا احتمال نتیجه برچسب هدف را بر اساس سوابق تاریخی بتوان پیش‌بینی کرد. درخت تصمیم در یادگیری ماشین جزء الگوریتم‌های یادگیری تحت نظارت است که برای نمایش طبقه‌بندی نمونه‌ها به کار می‌رود. این الگوریتم با تقسیم داده‌های آموزش به گره‌های<sup>۲</sup> مجزا و در نظر گرفتن ویژگی‌ها به صورت یک به یک ساخته و دارای ساختار درختی سلسله‌مراتبی<sup>۳</sup> است که اجزای اصلی آن شامل: یک گره ریشه، شاخه‌ها، گره‌های داخلی و گره‌های برگ می‌باشد. شروع یک درخت تصمیم با یک گره ریشه (مهم‌ترین ویژگی که در اینجا سابقه کار شناخته شد) است که هیچ شاخه ورودی ندارد. سپس شاخه‌های خروجی از گره ریشه وارد گره‌های داخلی شده که به عنوان گره‌های تصمیم نیز شناخته می‌شوند. گره‌های برگ همه نتایج ممکن را در مجموعه داده نشان می‌دهند. یادگیری درخت تصمیم با انجام یک جستجوی حریصانه برای شناسایی نقاط تقسیم بهینه در یک درخت، از استراتژی تقسیم<sup>۴</sup> و غلبه<sup>۵</sup> استفاده می‌کند، که این فرآیند تقسیم به صورت بازگشتی از بالا به پایین تکرار می‌شود تا زمانی که همه یا اکثر رکوردها تحت برچسب‌های کلاس خاصی طبقه‌بندی شوند. معیارهای متعددی برای تصمیم‌گیری در مورد محل تقسیم و میزان تقسیم وجود دارند که از میان این روش‌ها و معیارها دو روش کسب اطلاعات<sup>۶</sup> و

---

1. Unlabeled  
 2. Node  
 3. Hierarchical  
 4. Divide  
 5. Conquer  
 6. Information Gain



ناخالصی جینی<sup>۱</sup> به عنوان دو معیار محبوب تقسیم شناخته می‌شوند. این روش‌ها به ارزیابی کیفیت شرایط هر آزمون و اینکه چگونه می‌توان نمونه‌ها را در یک کلاس طبقه‌بندی کرد کمک می‌کند.

توضیح روش کسب اطلاعات به مفهوم آنتروپی<sup>۲</sup> مرتبط است که ناخالصی مقادیر نمونه را اندازه می‌گیرد و برای محاسبه آن از رابطه (۱) استفاده می‌شود.

$$(1) \quad Entropy(S) = -\sum_{c \in C} p(c) \log_2(p(c))$$

که عبارت‌های  $S$  مجموعه داده،  $C$  کلاس‌های مجموعه  $S$  و  $P(c)$  نسبت نقاط داده‌ای متعلق به کلاس  $c$  به تعداد کل نقاط داده در مجموعه  $S$  است. مقادیر آنتروپی می‌تواند بین ۰ و ۱ قرار گیرد. اگر تمام نمونه‌های مجموعه داده  $S$  متعلق به یک کلاس باشند، آنتروپی برابر صفر خواهد بود. اگر نیمی از نمونه‌ها به عنوان یک کلاس و نیمی دیگر در کلاس دیگری طبقه‌بندی شوند، آنتروپی به بالاترین حد خود یعنی ۱ خواهد رسید. برای انتخاب بهترین ویژگی جهت تقسیم و یافتن درخت تصمیم بهینه نیز باید از ویژگی (مشخصه یا متغیر) با کمترین مقدار آنتروپی استفاده شود. به تفاوت آنتروپی قبل و بعد از تقسیم در یک ویژگی معین نیز اطلاعات بدست آمده گفته می‌شود. ویژگی با بالاترین بهره اطلاعات بهترین تقسیم را ایجاد می‌کند زیرا بهترین کار را در طبقه‌بندی داده‌های آموزشی بر اساس طبقه‌بندی هدف خود انجام می‌دهد. کسب اطلاعات معمولاً با رابطه (۲) نشان داده می‌شود.

$$(2) \quad Information\ Gain(S, \alpha) = Entropy(S) - \sum_{v \in v\text{values}(\alpha)} \frac{|S_v|}{|S|} Entropy(S_v)$$

1. Gini Impurity  
2. Entropy

که  $\alpha$  یک ویژگی یا برچسب کلاس خاص بوده و  $\frac{|S_{\alpha}|}{|S|}$  نسبت مقادیر در مجموعه داده  $S_V$  به تعداد مقادیر در مجموعه داده  $S$  است. ابتدا آنروپی گره والد محاسبه می‌شود. سپس آنروپی هر گره جداگانه در تقسیم (انشعاب) فرموله و در انتها میانگین وزنی تمام گره‌ها محاسبه می‌شود. در روش دوم یا شاخص جینی، دو آیتم از یک جامعه به طور تصادفی انتخاب و در یک کلاس طبقه‌بندی می‌شوند. این شاخص با متغیر هدف طبقه‌بندی مانند «موفقیت» یا «شکست» کار می‌کند و فقط برای تقسیم (انشعاب) متغیرهای باینری (دودویی) کاربرد دارد. مجموع مجذور احتمال موفقیت و شکست  $(p^2 + q^2)$  برای محاسبه جینی زیرگره‌ها استفاده می‌شود. هرچه مقدار شاخص جینی بیشتر باشد، ارزش همگنی بیشتر است (مائمون و روکاج<sup>۱</sup>، ۲۰۱۴).

الگوریتم دیگری که در این مطالعه استفاده خواهد شد، الگوریتم جنگل تصادفی مجموعه‌ای از  $n$  درخت تصمیم است. هر درخت تصمیم در جنگل بر روی زیرمجموعه‌های مختلف مجموعه داده‌های آموزشی، آموزش داده می‌شود. جنگل تصادفی از انتخاب ویژگی تصادفی در حالی که درخت در حال رشد است استفاده می‌کند. در مورد مجموعه داده‌های چند بعدی این ویژگی بسیار مهم است، به ویژه زمانی که صدها یا هزاران ویژگی وجود دارد، به عنوان مثال در تشخیص پزشکی و اسناد، بسیاری از ویژگی‌های ضعیف مرتبط ممکن است اصلاً در یک درخت تصمیم ظاهر نشوند. این الگوریتم یک راه حل تخصصی برای اکثر مشکلات است و به عنوان یک روش یادگیری گروهی<sup>۲</sup> شناخته می‌شود زیرا گروهی از مدل‌های ضعیف با هم ترکیب می‌شوند تا این مدل قدرتمند را تشکیل دهند (سولیوان<sup>۳</sup>،

1. Maimon & Rokach

۲. یادگیری گروهی (Ensemble method) تشکیل شده از مجموعه‌ای از طبقه‌بندی کننده‌ها - به عنوان مثال درخت‌های تصمیم و پیش‌بینی‌های آنها برای شناسایی محبوب‌ترین نتیجه جمع‌آوری می‌شوند.

3. Sullivan

۲۰۱۷)، فرضیه نهایی در جنگل تصادفی با استفاده از روش رای اکثریت<sup>۱</sup> در بین درختان تولید می‌شود (تنها<sup>۲</sup> و همکاران، ۲۰۱۷).

#### ۴. یافته‌های پژوهش

این مقاله با استفاده از روش کریس و الگوریتم‌های طبقه‌بندی درخت تصمیم و جنگل تصادفی به بررسی رفتار و پیش‌بینی خرید دو نوع بیمه‌نامه حرف مشاغل آزاد پوشش دهنده بازنشستگی (۱۲ درصد) و پوشش دهنده ریسک‌های از کارافتادگی و بازماندگی علاوه بر بازنشستگی (۱۴ و ۱۸ درصد) پرداخت.

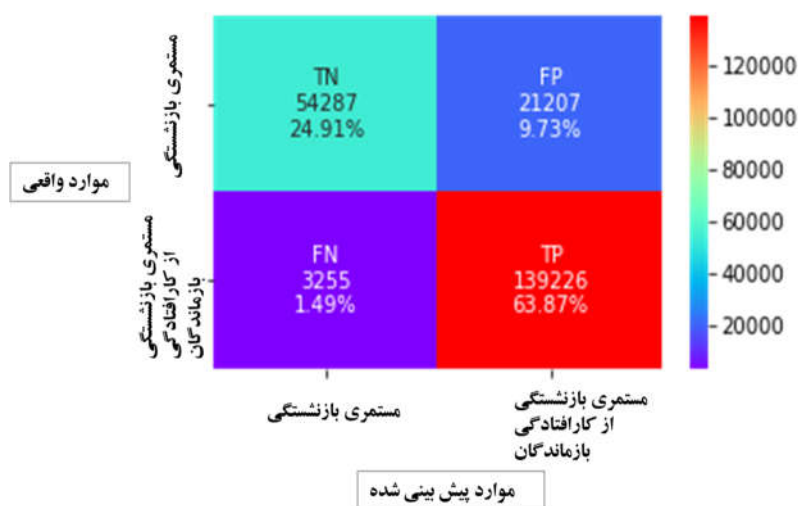
#### ۴-۱. نتایج پیش‌بینی الگوریتم درخت تصمیم و جنگل تصادفی

در این بخش و در ابتدا نتایج حاصل از بکارگیری الگوریتم درخت تصمیم ارائه خواهد شد. در اجرای این الگوریتم ۷۰ درصد از داده‌ها به عنوان داده‌های آموزش<sup>۳</sup> در نظر گرفته شد و از پارامتر اصلی<sup>۴</sup> حداکثر عمق<sup>۵</sup> برابر با ۶ و شاخص انشعاب جینی<sup>۶</sup> استفاده شده است. نتایج حاصل از این الگوریتم در ماتریس اغتشاش نمودار ۷ نمایش داده شده است. ماتریس اغتشاش به عنوان شاخصی از مشخصه‌های یک قانون طبقه‌بندی (تفکیک کننده) استفاده می‌شود. این ماتریس شامل تعداد عناصری است که به درستی یا نادرستی در هر کلاس طبقه‌بندی شده‌اند. تعداد مشاهداتی که به درستی برای هر کلاس طبقه‌بندی شده در قطر اصلی ماتریس و تعداد مشاهداتی که به اشتباه طبقه‌بندی شده در خارج از قطر اصلی قرار می‌گیرند. فرض کنید دو کلاس مثبت و منفی وجود دارد. شیوه کار به این صورت است که برای هر مشاهده در مجموعه تست، کلاس واقعی را با کلاسی که توسط طبقه‌بندی کننده

---

1. Voting  
2. Tanha  
3. Training Set  
4. Hyperparameter  
5. Maximum Depth  
6. Gini Index

آموزش دیده و تخصیص داده شده مقایسه می کنند. برای مشاهداتی که در کلاس مثبت طبقه بندی شده اند، اگر به درستی طبقه بندی شده باشند مثبت واقعی<sup>۱</sup> (TP) و مشاهده منفی که به اشتباه طبقه بندی شده مثبت کاذب (FP) نامیده می شود. به طور مشابه برای مشاهداتی که در کلاس منفی طبقه بندی شده اند، اگر مشاهده مثبتی باشد و به اشتباه طبقه بندی شده باشد منفی کاذب (FN) و بر عکس، یک مشاهده منفی که به درستی طبقه بندی شده باشد منفی واقعی<sup>۲</sup> (TN) نامیده می شود.



نمودار ۷. نتایج ماتریس اغتشاش در الگوریتم درخت تصمیم

منبع: یافته های پژوهش

عملکرد الگوریتم را می توان با معیارهای مختلفی ارزیابی کرد. نخستین معیار، معیار دقت<sup>۳</sup> است که طبق رابطه (۳) بیان می شود.

1. True Positive
2. True Negative
3. Accuracy Score

$$(۳) \quad Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

این معیار درصد طبقه‌بندی صحیح در هر دو بیمه‌نامه را اندازه‌گیری می‌کند. مطابق نمودار ۷ می‌توان دید که ۱۳۹۲۲۶ نفر به درستی متقاضی بیمه‌نامه ۱۴ و ۱۸ درصدی و ۵۴۲۸۷ نفر نیز به درستی متقاضی بیمه‌نامه ۱۲ درصدی طبقه‌بندی شده‌اند. بنابراین، از مجموع ۲۱۷۹۷۵ رکورد داده‌های تست، ۸۸ درصد آن‌ها به درستی طبقه‌بندی شده‌اند. نرخ دقت ۸۸ درصد نشان می‌دهد که مدل می‌تواند خروجی یا خریداری بیمه‌نامه ۱۲ درصدی در مقابل بیمه‌نامه ۱۴ و ۱۸ درصدی را با دقت ۸۸ درصد پیش‌بینی کند. معیار دوم، معیار پوشش<sup>۱</sup> است که طبق رابطه (۴) به ما می‌گوید از میان افرادی که واقعاً متقاضی بیمه ۱۴ و ۱۸ درصدی بوده‌اند (سطر پایین ماتریس اغتشاش)، چند درصد درست تشخیص داده شده‌اند. در واقع تمرکز اصلی معیار پوشش بر روی داده‌هایی است که واقعاً متقاضی بیمه ۱۴ و ۱۸ درصدی (با خدمات مستمری بازنشستگی، از کارافتادگی و بازنماندگان) بوده‌اند و هدف آن پوشش حداکثری این گروه است تا هیچ یک از آن‌ها به اشتباه در دسته دیگر قرار نگیرند. از نتایج ماتریس اغتشاش می‌توان دید که از مجموع ۱۴۲۴۸۱ نفری که متقاضی بیمه ۱۴ و ۱۸ درصدی بوده‌اند، تعداد ۱۳۹۲۲۶ نفر یا ۹۷/۷ درصد به درستی تشخیص داده شده‌اند.

$$(۴) \quad Recall = \frac{TP}{TP + FN}$$

معیار سوم، معیار صحت<sup>۲</sup> است که طبق رابطه (۵) تمرکز اصلی آن بر روی مواردی است که توسط الگوریتم در دسته ۱۴ و ۱۸ درصدی طبقه‌بندی و متقاضی این بیمه‌نامه پیش‌بینی شده‌اند (ستون دوم ماتریس اغتشاش) و به ما می‌گوید صحت این پیش‌بینی چقدر بوده است. از نتایج ماتریس اغتشاش می‌توان دید که از مجموع ۱۶۰۴۳۳ نفری که پیش‌بینی شده متقاضی بیمه ۱۴ و ۱۸ درصدی باشند، تعداد ۱۳۹۲۲۶ نفر یا ۸۶/۷ درصد واقعاً متقاضی بوده و این پیش‌بینی در مورد آن‌ها صحیح است.

---

1. Recall  
2. Precision

$$(۵) \quad Precision = \frac{TP}{TP + FP}$$

بین معیارهای صحت و پوشش یک بدهستان<sup>۱</sup> وجود دارد و تلاش برای بهبود یک معیار اغلب منجر به بدتر شدن معیار دوم می گردد. ساده ترین راه حل این مشکل این است که در یک تصمیم گیری چند معیاره (MCDM) هر دو را ترکیب کرده و جمع وزنی این معیارها را مدنظر قرار داد. این معیار ترکیبی امتیاز-F است که طبق رابطه (۶) مقدار آن برابر ۹۱/۹ درصد محاسبه می شود. این معیار ترکیب متعادلی از دو معیار صحت و پوشش است. هر چه معیار F بیشتر باشد نشان دهنده تایید دقت مدل است به این معنی که حد آستانه در بهترین حالت خود قرار گرفته و پیش بینی با کمترین خطا همراه است.

$$(۶) \quad F = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

در محاسبه معیارهای یادشده تمرکز بر بیمه نامه ۱۴ و ۱۸ درصدی بود و دقت، صحت و پوشش الگوریتم در طبقه بندی این گروه ارزیابی شد. به عبارت دیگر، عملکرد مدل و الگوریتم در پیش بینی افرادی که متقاضی خرید بیمه نامه پوشش ریسک بالاتر (بیمه ۱۴ و ۱۸ درصدی) هستند مدنظر بوده است. اما از آنجائی که هر دو مقدار متغیر هدف یا هر دو نوع بیمه نامه برای ما به یک اندازه اهمیت دارند، می توان توجه را به سمت گروه دیگر یعنی بیمه ۱۲ درصدی گرداند و با استفاده از اطلاعات ماتریس اغتشاش عملکرد الگوریتم در پوشش این گروه را بررسی کرد. برای این منظور، این بار، بر افرادی که واقعاً متقاضی بیمه ۱۲ درصدی بوده اند (سطر بالای ماتریس اغتشاش) تمرکز می شود تا درستی تشخیص الگوریتم برای این گروه تعیین شود. این معیار در اصطلاح «تشخیص پذیری»<sup>۲</sup> خوانده می شود و همانطور که از ماتریس اغتشاش می توان دید، از مجموع ۷۵۴۹۴ نفر که واقعاً متقاضی بیمه نامه ۱۲ درصدی (با خدمات مستمری بازنشستگی) بوده اند، تعداد ۵۴۲۸۷ یا ۷۱/۹ درصد به درستی طبقه بندی و پیش بینی

1. Trade-off  
2. Specificity

شده‌اند. یا توجه را به مواردی که متقاضی بیمه‌نامه ۱۲ درصدی پیش‌بینی شده‌اند (ستون اول ماتریس اغتشاش) متمرکز کرد و صحت این پیش‌بینی را سنجید. این شاخص «ارزش پیش‌بینی شده منفی»<sup>۱</sup> نام دارد و می‌توان دید که طبق نتایج ماتریس اغتشاش، از ۵۷۵۴۲ نفری که پیش‌بینی شده متقاضی بیمه‌نامه ۱۲ درصدی باشند تعداد ۵۴۲۸۷ یا ۹۴/۳ درصد واقعاً متقاضی خرید این بیمه‌نامه بوده و این پیش‌بینی در مورد آن‌ها صحیح است. در جدول ۲ به طور خلاصه نتایج حاصل از معیارهای ارزیابی الگوریتم درخت تصمیم ارائه شده است.

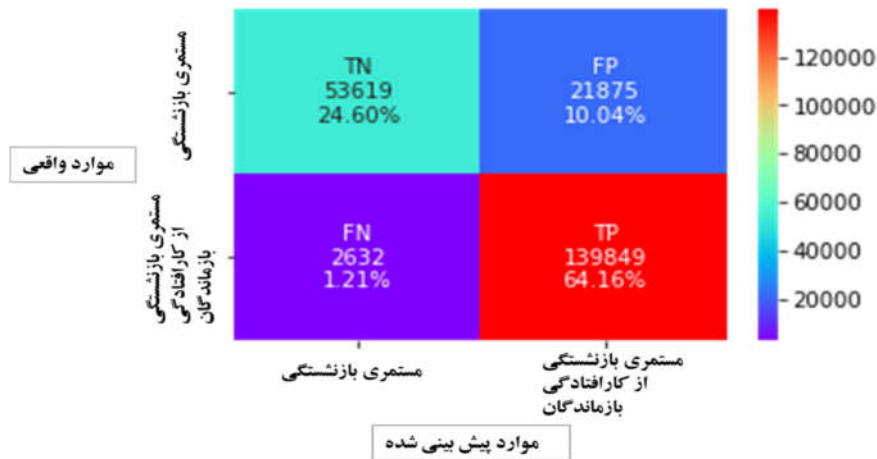
جدول ۲. نتایج حاصل از الگوریتم درخت تصمیم

نام معیار	معیار (درصد)
دقت	۸۸/۷
پوشش	۹۷/۷
صحت	۸۶/۷
امتیاز-F	۹۱/۹
تشخیص پذیری	۷۱/۹
ارزش پیش‌بینی شده منفی	۹۴/۳

منبع: یافته‌های پژوهش

علاوه بر درخت تصمیم، در این مقاله از الگوریتم جنگل تصادفی نیز برای پیش‌بینی خرید بیمه‌نامه‌ها استفاده شده است. نتایج بدست آمده از بکارگیری الگوریتم جنگل تصادفی با پارامترهای اصلی، حداکثر عمق ۸ و شاخص جینی و تعداد ۱۰۰ درخت مورد نیاز برای تخمین<sup>۲</sup>، در نمودار ۸ نشان داده شده است.

1. Negative Predictive Value  
2. N-Estimators



#### نمودار ۸. نتایج ماتریس اغتشاش در الگوریتم جنگل تصادفی

موارد واقعی و پیش‌بینی شده انتخاب دو نوع بیمه‌نامه را می‌توان در ماتریس اغتشاش دید. مانند قبل عملکرد الگوریتم را می‌توان با معیارهای مختلفی ارزیابی کرد. معیار دقت درصد طبقه‌بندی صحیح در هر دو بیمه‌نامه را اندازه‌گیری می‌کند. مطابق نمودار ۸ می‌توان دید که ۱۳۹۸۴۹ نفر به درستی متقاضی بیمه‌نامه ۱۴ و ۱۸ درصدی و ۵۳۶۱۹ نفر نیز به درستی متقاضی بیمه‌نامه ۱۲ درصدی طبقه‌بندی شده‌اند. بنابراین، از مجموع ۲۱۷۹۷۵ رکورد داده‌های تست، ۸۸ درصد آن‌ها به درستی طبقه‌بندی شده‌اند. معیار پوشش نشان می‌دهد از میان افرادی که واقعاً متقاضی بیمه ۱۴ و ۱۸ درصدی بوده‌اند (سطر پایین ماتریس اغتشاش)، چند درصد درست تشخیص داده شده‌اند. در واقع تمرکز اصلی معیار پوشش بر روی داده‌هایی است که واقعاً متقاضی بیمه ۱۴ و ۱۸ درصدی (با خدمات مستمری بازنشستگی، از کارافتادگی و بازماندگان) بوده‌اند و هدف آن پوشش حداکثری این گروه است تا هیچ یک از آن‌ها به اشتباه در دسته دیگر قرار نگیرند. از نتایج ماتریس اغتشاش می‌توان دید که از مجموع ۱۴۲۴۸۱ نفری که متقاضی بیمه ۱۴ و ۱۸ درصدی بوده‌اند، تعداد ۱۳۹۸۴۹ نفر یا ۹۸/۱ درصد به درستی تشخیص داده شده‌اند. معیار سوم، معیار صحت است که بر روی مواردی که توسط الگوریتم در دسته ۱۴ و ۱۸ درصدی طبقه‌بندی و متقاضی این بیمه‌نامه پیش‌بینی شده‌اند (ستون دوم ماتریس اغتشاش) تمرکز دارد و نشان می‌دهد



صحت این پیش‌بینی چقدر بوده است. از نتایج ماتریس اغتشاش می‌توان دید که از مجموع ۱۶۱۷۲۴ نفری که پیش‌بینی شده متقاضی بیمه ۱۴ و ۱۸ درصدی باشند، تعداد ۱۳۹۸۴۹ نفر یا ۸۶/۴ درصد واقعاً متقاضی بوده و این پیش‌بینی در مورد آن‌ها صحیح است. معیار امتیاز-F که ترکیبی است از معیارهای صحت و پوشش نیز برابر ۹۱/۹ درصد محاسبه می‌شود.

اگر این بار توجه را معطوف به گروه دیگر یعنی بیمه ۱۲ درصدی کرده و با استفاده از اطلاعات ماتریس اغتشاش عملکرد الگوریتم در پوشش این گروه را بررسی کنیم می‌توان معیار «تشخیص پذیری» را محاسبه کرد که نشان می‌دهد چند درصد این افراد توسط الگوریتم درست تشخیص داده شده‌اند. همانطور که از ماتریس اغتشاش می‌توان دید، از مجموع ۷۵۴۹۴ نفر که واقعاً متقاضی بیمه‌نامه ۱۲ درصدی (با خدمات مستمری بازنشستگی) بوده‌اند، تعداد ۵۳۶۱۹ یا ۷۱ درصد به درستی طبقه‌بندی و پیش‌بینی شده‌اند. یا می‌توان به ستون اول ماتریس اغتشاش توجه کرد و صحت این پیش‌بینی یعنی شاخص «ارزش پیش بینی شده منفی» را سنجید. طبق نتایج ماتریس اغتشاش، از ۵۶۲۵۱ نفری که پیش‌بینی شده متقاضی خرید بیمه‌نامه ۱۲ درصدی باشند تعداد ۵۳۶۱۹ یا ۹۵/۳ درصد واقعاً متقاضی این بیمه‌نامه بوده و این پیش‌بینی در مورد آن‌ها صحیح است. در جدول ۳ به طور خلاصه نتایج حاصل از معیارهای ارزیابی جنگل تصادفی ارائه شده است.

جدول ۳. نتایج حاصل از الگوریتم جنگل تصادفی

نام معیار	دقت معیار(درصد)
دقت	۸۸/۷
پوشش	۹۸/۱
صحت	۸۶/۴
امتیاز-F	۹۱/۹
تشخیص پذیری	۷۱
ارزش پیش‌بینی شده منفی	۹۵/۳

منبع: یافته‌های پژوهش

۴-۲. تحلیل و تفسیر ویژگی‌های فردی مؤثر بر انتخاب بیمه‌نامه‌های حرف و مشاغل آزاد با اتکا به نتایج به دست آمده می‌توان گفت که سابقه کار افراد از موثرترین عوامل در انتخاب بیمه‌نامه افراد است. الگوریتم درخت تصمیم برای همه افراد از هر جنسیتی در صورتی که سابقه کار کمتر از ۴/۵ سال داشته باشند پیش‌بینی کرده که متقاضی بیمه ۱۲ درصدی خواهند بود و همچنین این الگوریتم نشان می‌دهد افراد تنها با افزایش سابقه کار تمایل به بیمه‌های ۱۴ و ۱۸ درصدی با پوشش بیشتر از جمله از کارافتادگی و ریسک فوت پیدا می‌کنند. البته در اینجا جنسیت تاثیرگذار است و طبق نتایج بدست آمده مردان بسیار زودتر به خریداری بیمه‌های با پوشش بیشتر تمایل پیدا می‌کنند. جداول ۴ و ۵ به تفکیک میزان تمایل زنان و مردان در سوابق کاری، سن و سطوح درآمدی مختلف برای انواع بیمه‌نامه‌ها را نمایش می‌دهند.

همانطور که از جدول ۴ پیداست، زنان با سابقه کمتر از ۵/۵ سال پیش‌بینی شده که متقاضی بیمه‌نامه ۱۲ درصدی با پوشش بازنشستگی و زنان با سابقه کار بالاتر از ۹/۵ سال پیش‌بینی شده که خریدار بیمه‌نامه‌های ۱۴ و ۱۸ درصدی باشند. اما پیش‌بینی تقاضا برای زنانی که سابقه کاری بین ۵/۵ تا ۹/۵ سال دارند کمی سخت‌تر بوده و به دیگر متغیرها بستگی دارد. با توجه به نتایج به دست آمده می‌توان گفت زنانی که ۵/۵ تا ۷/۵ سال سابقه کاری، سن کمتر از ۳۷/۵ سال و درآمد کمتر از ۲ میلیون و ۲۰۰ هزار تومان دارند، با احتمال بیشتر متقاضی بیمه‌نامه ۱۲ درصدی خواهند بود اما برای زنان با سن و درآمد بیشتر احتمال خریداری بیمه ۱۴ و ۱۸ درصدی بیشتر خواهد بود. نتیجه مهمی که می‌توان گرفت این است که سابقه کار، سن و درآمد تاثیر مثبتی بر انتخاب بیمه‌نامه ۱۴ و ۱۸ درصدی با پوشش از کار افتادگی و فوت علاوه بر بازنشستگی توسط زنان دارند. در واقع، در سنین و سوابق بالا تقاضا برای بیمه‌نامه با خدمات بیشتر و در سنین و سوابق پایین تقاضا برای بیمه‌نامه با خدمات کمتر اما ارزان‌تر وجود دارد.

جدول ۴. خلاصه نتایج برآورد درخت تصمیم برای زنان بیمه‌شده حرف و مشاغل آزاد

سابقه (سال)	سن	متوسط درآمد (میلیون ریال)	پیش‌بینی (درصد احتمال)
-------------	----	---------------------------	------------------------

بیمه‌نامه ۱۲ درصد (احتمال ۱۰۰)	در هر سطحی از درآمد	در هر سن	۵ تا ۵/۵
بیمه‌نامه ۱۲ درصد (احتمال ۵۵)	کمتر از ۲۲	کوچکتر از ۳۷/۵	۷/۵ تا ۵/۵
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۶۳)	بیشتر از ۲۲		
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۵۷)	کمتر از ۲۱/۴۷	بزرگتر از ۳۷/۵	۷/۵ تا ۵/۵
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۷۹)	بیشتر از ۲۱/۴۷		
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۶۱)	کمتر از ۲۱/۹۷	کوچکتر از ۳۹/۵	۸/۵ تا ۷/۵
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۷۹)	بیشتر از ۲۱/۹۷		
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۶۵)	در هر سطحی از درآمد	بین ۳۹/۵ تا ۵۹/۵	۸/۵ تا ۷/۵
بیمه‌نامه ۱۲ درصد (احتمال ۵۹)	در هر سطحی از درآمد	بزرگتر از ۵۹/۵	
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۷۳)	در هر سطحی از درآمد	کوچکتر از ۳۹/۵	۹/۵ تا ۸/۵
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۸۷)	در هر سطحی از درآمد	بزرگتر از ۳۹/۵	
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۱۰۰)	در هر سطحی از درآمد	در هر سن	بیشتر از ۹/۵

منبع: یافته‌های پژوهش

برای مردان گرایش به انتخاب بیمه‌نامه‌های ۱۴ و ۱۸ درصدی با پوشش کامل بسیار سریع‌تر اتفاق می‌افتد. آنچنان که در جدول ۵ نشان داده شده خریداری بیمه‌نامه ۱۲ درصدی تنها به مردان با سابقه کمتر از ۴/۵ محدود بوده و در سطوح سابقه بالاتر از ۶/۵ مردان بطور کامل پیش‌بینی می‌شود که متقاضی بیمه‌نامه ۱۴ و ۱۸ درصدی باشند. برای مردان با سابقه کاری بین ۴/۵ تا ۶/۵ نیز طبق نتایج به دست آمده می‌توان دید که خریداری بیمه‌نامه ۱۲ درصدی به افراد جوان با سن کمتر از ۲۴/۵ محدود خواهد بود و با افزایش سن تقاضا برای بیمه‌نامه ۱۴ و ۱۸ درصدی بالاتر رفته و الگوریتم احتمال بیشتری برای انتخاب این بیمه‌نامه پیش‌بینی می‌کند. طبق نتایج بدست آمده، پیش‌بینی شده که مردان سنین بین ۲۴/۵ سال تا ۲۷/۵ به احتمال ۶۱ درصد متقاضی بیمه ۱۴ و ۱۸ درصد باشند و این احتمال برای سنین ۲۷/۵ سال تا ۳۰/۵ به ۷۰ درصد، برای سنین ۳۰/۵ سال تا ۳۳/۵ به ۷۴ درصد، برای سنین ۳۳/۵ سال تا ۳۶/۵ به ۸۰ درصد، برای سنین بین ۳۶/۵ سال تا ۳۹/۵

به ۸۴ درصد، برای سنین ۳۹/۵ سال تا ۴۹/۵ به ۸۹ درصد و برای سنین بیشتر از ۴۹/۵ سال نیز به بیش از ۹۳ درصد می‌رسد.

در تمامی این نتایج و با توجه به ساختار الگوریتم درخت تصمیم، می‌توان به کم اهمیت بودن درآمد افراد در تعیین نوع بیمه‌نامه آن‌ها اشاره کرد. اما به خوبی می‌توان دید که سه عامل دیگر یعنی سن، جنسیت و سابقه از توضیح دهندگی و قابلیت پیش‌بینی بالایی برخوردارند. به طور کلی، یافته‌ها نشان می‌دهند که برای افراد با سن و سابقه کم بیمه‌نامه ۱۲ درصد و برای افراد با سن و سابقه بالا بیمه‌نامه ۱۴ و ۱۸ درصد جذابیت بیشتری خواهد داشت.

جدول ۵. خلاصه نتایج برآورد درخت تصمیم برای مردان بیمه‌شده حرف و مشاغل آزاد

پیش‌بینی (درصد احتمال)	متوسط درآمد (میلیون ریال)	سن	سابقه (سال)
بیمه‌نامه ۱۲ درصد (احتمال ۱۰۰)	در هر سطحی از درآمد	در هر سن	۴/۵ تا ۱۰
بیمه‌نامه ۱۲ درصد (احتمال ۵۵)	در هر سطحی از درآمد	کوچکتر از ۲۴/۵	۴/۵ تا ۶/۵
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۶۱)	در هر سطحی از درآمد	۲۴/۵ تا ۲۷/۵	
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۷۰)	در هر سطحی از درآمد	۲۷/۵ تا ۳۰/۵	
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۷۴)	در هر سطحی از درآمد	۳۰/۵ تا ۳۳/۵	
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۸۰)	در هر سطحی از درآمد	۳۳/۵ تا ۳۶/۵	
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۸۴)	در هر سطحی از درآمد	۳۶/۵ تا ۳۹/۵	
بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۸۹)	در هر سطحی از درآمد	۳۹/۵ تا ۴۹/۵	
بیمه‌نامه ۱۲ درصد (احتمال ۹۳)	در هر سطحی از درآمد	بزرگتر از ۴۹/۵	

بیشتر از ۶/۵	در هر سن	در هر سطحی از درآمد	بیمه‌نامه ۱۴ و ۱۸ درصد (احتمال ۱۰۰)
--------------	----------	---------------------	-------------------------------------

منبع: یافته‌های پژوهش

## ۵. جمع‌بندی و نتیجه‌گیری

در این مطالعه سعی شد با استفاده از تجزیه و تحلیل داده‌های مربوط به بیمه‌شدگان حرف و مشاغل آزاد سازمان تامین اجتماعی در سال ۱۳۹۹، درکی از مشخصه‌های افراد این گروه حاصل شود. این درک و فهم از داده‌های مربوط به افراد بیمه‌شده می‌تواند در جهت ایجاد برنامه‌هایی به منظور افزایش نرخ نفوذ این نوع بیمه‌ها مفید باشد زیرا با توجه به نتایج حاصل شده می‌توان برای هر وضعیت یک استراتژی مخصوص را در پیش گرفت و با هدف قرار دادن گروه مد نظر به افزایش پوشش این نوع بیمه‌ها کمک کرد. در بین چهار مشخصه اصلی از افراد، نتایج بدست آمده نشان دادند که سه ویژگی سابقه کار، سن افراد و جنسیت به ترتیب بیشترین تاثیر را بر انتخاب بیمه‌نامه با نرخ‌های متفاوت دارند و پس از آن درآمد هم می‌تواند تا حدودی تاثیرگذار باشد.

مهمترین نتایج بدست آمده عبارتند از:

- جنسیت: از بررسی رابطه بین تعداد و جنسیت بیمه‌شدگان مشخص شد که مردان با ۵۵/۴ درصد بیشترین متقاضیان انواع بیمه‌نامه‌های حرف و مشاغل آزاد هستند که البته بیش از همه بیمه‌نامه با نرخ ۱۸ درصد پوشش دهنده ریسک از کارافتادگی و فوت (مستمری بازماندگان) علاوه بر بازنشستگی برای آن‌ها جذابیت دارد. در مقابل اما زنان به طور عمده متقاضی بیمه‌های ۱۲ درصد هستند. برای توضیح و درک این نتیجه باید به این واقعیت توجه کرد که در ایران بار تکفل همسر و فرزندان بر عهده مردان است و به همین دلیل است که مردان بیشتر متقاضی بیمه‌های با پوشش بالاتر ۱۴ و ۱۸ درصدی هستند و در مقابل این بیمه‌نامه‌ها با خدمات کامل‌تر از جمله مستمری

بازماندگان جذابیت زیادی برای زنان ندارد چون بازماندگان زنان بیمه‌شده عموماً واجد شرایط دریافت مستمری بازماندگی نیستند. این نتیجه با نتایج مطالعه عبدی (۱۳۸۵) تطابق دارد که مردان را عمده مشتریان بیمه‌نامه‌های حرف و مشاغل آزاد و اختیاری معرفی می‌کند که انگیزه و اهداف اصلی آنان از بیمه به ترتیب بازنشستگی، درمان، از کارافتادگی و فوت بوده است.

● سن: جوانان به طور عمده متقاضی بیمه‌های ۱۲ درصد هستند که این یعنی آن‌ها نگرانی از بابت ریسک فوت و مستمری بازماندگان ندارند و در نتیجه متقاضی آن نیستند. این امر تا حدودی به خصلت نزدیک‌بین بودن<sup>۱</sup> افراد نیز برمی‌گردد. با این وجود مشاهده شد که افزایش سن خرید بیمه‌نامه‌های ۱۴ و ۱۸ درصدی را به‌طور ویژه برای مردان جذابتر می‌کند. مشاهده شد که مردان با سابقه کاری بالا و افزایش سن تمایل به استفاده از خدمات کامل بیمه‌ای یعنی نرخ ۱۴ و ۱۸ درصد دارند اما در بین زنان متغیر سن به طور عمده تاثیر چندانی بر تمایل آن‌ها ندارد و زنان با افزایش سن و در هر گروه سنی همچنان متقاضی بیمه ۱۲ درصد هستند. این یعنی چون در ایران طبق قوانین بازنشستگی همسر و فرزندان تحت تکفل مردان هستند، در نتیجه بیمه‌نامه ۱۴ و ۱۸ درصدی با پوشش مستمری بازماندگان و ریسک فوت با افزایش سن برای مردان جذابیت پیدا می‌کند اما برای زنان در هر گروه سنی جذابیت ندارد چراکه همانطور که گفته شد بازماندگان زنان بیمه‌شده عموماً واجد شرایط دریافت مستمری بازماندگی نخواهند بود. با توجه به این نتایج می‌توان گفت که در ارتباط با مردان، اولویت سازمان در بیمه حرف و مشاغل آزاد باید برجسته‌سازی و تبلیغ خدمات بیمه‌نامه با نرخ ۱۸ درصد و در ارتباط با زنان باید بر تبلیغ خدمات بیمه‌نامه با نرخ ۱۲ درصد تمرکز کند. علاوه بر این موارد، نتایجی که از برآورد الگوریتم درخت تصمیم حاصل شد نشان می‌دهند که:

• سابقه کاری: افراد با سابقه کاری کمتر از ۴/۵ سال با هر سطحی از درآمد، سن یا جنسیت متقاضی بیمه‌های ۱۲ درصد با خدمات مستمری بازنشستگی هستند، اما مردان زمانی که سابقه کاری آن‌ها از ۶/۵ سال بیشتر می‌شود و زنان وقتی سابقه کاری آن‌ها بالاتر از ۹/۵ سال است به سمت بیمه‌های ۱۴ و ۱۸ درصدی (خدمات کامل‌تر) روی می‌آورند.

• درآمد: طبق نتایج به دست آمده اگرچه متغیر درآمد به دلیل عدم اظهار مقادیر واقعی توسط بیمه‌شدگان چندان اثر گذار نبود، اما باز هم مشاهده شد که افراد با درآمد بالا متقاضی بیمه‌نامه ۱۴ و ۱۸ درصدی هستند و در درآمدهای پایین متقاضی بیمه با خدمات ارزان‌تر (۱۲ درصد) هستند. البته اثر درآمد چندان تعیین کننده نیست و در بازه‌های بالا یا پایین سوابق یا سن، دیگر قدرت پیش‌بینی کنندگی ندارد و در هر سطحی از درآمد، در سنین و سوابق بالا تقاضا برای بیمه‌نامه با خدمات بیشتر و در سنین و سوابق پایین تقاضا برای بیمه‌نامه با خدمات کمتر اما ارزان‌تر مشاهده می‌شود.

دقت حاصل از الگوریتم درخت تصمیم و جنگل تصادفی به ۸۸ درصد رسید، یعنی به احتمال ۸۸ درصد به درستی می‌توان پیش‌بینی کرد که افراد با توجه به مشخصه‌هایی که دارند مانند سن، جنسیت، سابقه کار و درآمد، متقاضی بیمه حرف و مشاغل آزاد با خدمات مستمری بازنشستگی هستند یا متقاضی بیمه با خدمات مستمری بازنشستگی، از کار افتادگی و بازماندگان خواهند بود. در کل می‌توان گفت که دو الگوریتم به خوبی عمل کرده‌اند و دقت مشابهی در پیش‌بینی تقاضای بیمه‌نامه دارند.

نتایج این پژوهش، از طریق ارائه تصویری از مشخصه‌های کلیدی بیمه‌شدگان حرف و مشاغل آزاد و انگیزه آن‌ها برای انتخاب هر یک از انواع این بیمه‌ها، می‌تواند به سازمان تأمین اجتماعی در گسترش پوشش بیمه و بهبود مدیریت ارتباط با این دسته از بیمه‌شدگان کمک کند. برای مثال، با توجه به انگیزه پایین زنان و جوانان برای انتخاب بیمه‌های با خدمات

گسترده، می‌توان از طریق ارائه مشوق‌ها یا خدمات کوتاه‌مدت، جذابیت این نوع بیمه‌ها را برای این گروه از افراد افزایش داد.

#### ۶. تقدیر و تشکر

نویسندگان وظیفه خود می‌دانند از کارشناسان سازمان تأمین اجتماعی که در اختیار قرار دادن داده‌ها همکاری داشته‌اند و همچنین داوران محترم به سبب ارائه نقطه نظرات ارزشمند که به غنای مقاله کمک شایانی داشته است تشکر و قدردانی نمایند.



## منابع

- Abdi, A. (2006). Examining the issues and problems of insured persons and the free and optional jobs in interaction with the social security organization. *Social Security Quarterly*, 8(1). 255-282. (In Persian)
- Abdi, F., Khalili-Damghani, K., & Abolmakarem, S. (2017). Solving customer insurance coverage sales plan problem using a multi-stage data mining approach. *Kybernetes*.1, 2-19.
- Abdul-Rahman, S., Arifin, N. F. K., Hanafiah, M., & Mutalib, S. (2021). Customer Segmentation and Profiling for Life Insurance using K-Modes Clustering and Decision Tree Classifier. *International Journal of Advanced Computer Science and Applications (IJACSA)*, 12(9)
- Azzone, M.; Barucci, E.; Mancayo, G.G.; Marazzina, D. (2022). A Machine Learning Model for Lapse Prediction in Life Insurance Contracts. *Expert Systems with Applications*, 191.
- Chapman, P., Clinton, J., Kerber, R., Khabaza, T., Reinartz, T., Shearer, C., & Wirth, R. (2000). CRISP-DM 1.0: Step-by-step data mining guide. *SPSS inc*, 9,13.
- Chen, I. J., & Popovich, K. (2003). Understanding customer relationship management (CRM): People, process and technology. *Business process management journal*. 672-688.
- Chen, Y., & Hu, L. (2005). Study on data mining application in CRM system based on insurance trade. *In Proceedings of the 4th international conference on Electronic commerce*. 839-841.
- Hurwitz, J., & Kirsch, D. (2018). Machine learning for dummies. *IBM Limited Edition*, 75.
- Hosseini, S. M. S., Maleki, A., & Gholamian, M. R. (2010). Cluster analysis using data mining approach to develop CRM methodology to assess the customer loyalty. *Expert Systems with Applications*, 37(7), 5259-5264.
- Kracklauer, A. H., Mills, D. Q., & Seifert, D. (2004). Customer management as the origin of collaborative customer relationship management. *In Collaborative Customer Relationship Management (pp.3-6)*. Springer, Berlin, Heidelberg
- Khalili-Damghani, K., Abdi, F., & Abolmakarem, S. (2019). Solving customer insurance coverage recommendation problem using a two-stage

- clustering-classification model. *International Journal of Management Science and Engineering Management*, 14(1), 9-19.
- Kong, H.; Yun, W.; Joo, W.; Kim, J.H.; Kim, K.K.; Moon, I.C.; & Kim, W.C. (2022). Constructing a personalized recommender system for life insurance products with machine-learning techniques. *Intelligent Systems in Accounting, Finance and Management*, 29 (4). 242-253.
  - Maimon, O. Z., & Rokach, L. (2014). Data mining with decision trees: *theory and applications (Vol. 81)*. World scientific
  - Mau, S., Pletikosa, I., & Wagner, J. (2018). Forecasting the next likely purchase events of insurance customers: A case study on the value of data-rich multichannel environments. *International Journal of Bank Marketing*, 1123-1144.
  - Mollamohammadi, R., & Mostofi, M.R. (2014). The Factors Affecting the Success of the Social Security Organization in Paying Retirement Pension to Those Insured by Qom'First Branch of Social Security Organization. *Organizational Culture Management*, 12(2), 299-323. (In Persian)
  - Motdin, N.; Nazarian, R.; Daman-Kshideh, M., & Seifipour, R. (2021). Designing a Comparative Model of Bank Credit Risk Using Neural Network Models, Survival Probability Function and Support Vector Machine. *Journal of Economic Modeling Research*, 11 (45), 199-230. (In Persian)
  - Motafakkerzad, M.A.; & Ghafarnejad Mehraban, A. (2011). Intelligent Modeling of Asymmetric Effects of Monetary Shocks on Output in Iran (Neural Network Application). *Journal of Economic Modeling Research*, 2 (4), 83-102. (In Persian)
  - Najafi, A. (2019). Predictability of loyalty and separation of self-insurance Persons of Social Security Organization based on data mining method. *Social Security Quarterly*, 15(1). 88-109. (In Persian)
  - Ngai, E. W., Xiu, L., & Chau, D. C. (2009). Application of data mining techniques in customer relationship management: A literature review and classification. *Expert systems with applications*, 36(2), 2592-2602.
  - Parmah, S.; Mardomdar, S., & Heidari, A. (2020). Macroeconomic Variables and Demand for Self-employment Insurance in the Social Security Organization. *Social Security Quarterly*, 16(1). 41-59. (In Persian)
  - Rahman, S., Arefin, K. Z., Masud, S., Sultana, S., & Rahman, R. M. (2017, April). Analyzing Life Insurance Data with Different Classification

Techniques for Customers' Behavior Analysis. *In Asian Conference on Intelligent Information and Database Systems*. 15-25.

- Rygielski, C., Wang, J. C., & Yen, D. C. (2002). Data mining techniques for customer relationship management. *Technology in society*, 24(4), 483-502.
- Severino, Matheus Kempa, and Yaohao Peng. "Machine learning algorithms for fraud prediction in property insurance: Empirical evidence using real-world microdata." *Machine Learning with Applications* 5 (2021): 100074.
- Shokohyar, S.; Rezaeian, A., & Boroufar, A. (2017). Identifying the customer behavior model in life insurance Sector using data mining. *Management Research in Iran*, 20(4). 65-94. (In Persian)
- Sullivan, W. (2017). Machine Learning For Beginners Guide Algorithms: Supervised & Unsupervised Learning. Decision Tree & Random Forest Introduction. *Healthy Pragmatic Solutions Inc*
- Tanha, J., Van Someren, M., & Afsarmanesh, H. (2017). Semi-supervised self-training for decision tree classifiers. *International Journal of Machine Learning and Cybernetics*, 8(1), 355-370.
- Wirth, R., & Hipp, J. (2000, April). CRISP-DM: Towards a standard process model for data mining. *In Proceedings of the 4th international conference on the practical applications of knowledge discovery and data mining*. 29-40.