



## یافتن مشتریان احتمالی خرید بیمه با تکنیک‌های داده‌کاوی

امیر اولیایی، دانشجوی کارشناسی مهندسی کامپیوتر - نرم‌افزار، [stu.amir.oliaei@kashmar.ac.ir](mailto:stu.amir.oliaei@kashmar.ac.ir)

علی حاضری، دانشجوی کارشناسی مهندسی کامپیوتر - نرم‌افزار، [stu.ali.hazeri@kashmar.ac.ir](mailto:stu.ali.hazeri@kashmar.ac.ir)

محمد جمالی، دانشجوی کارشناسی مهندسی کامپیوتر - نرم‌افزار، [stu.mohammad.jamali@kashmar.ac.ir](mailto:stu.mohammad.jamali@kashmar.ac.ir)

عاطفه خزاعی، دکتری مهندسی کامپیوتر - نرم‌افزار، [atefeh.khazaei@kashmar.ac.ir](mailto:atefeh.khazaei@kashmar.ac.ir)

### چکیده

در سال‌های اخیر صنعت بیمه رشدی چشمگیر داشته است و شرکت‌های مختلف با خدمات گوناگون پا به عرصه گذاشته‌اند. بازاریابی موفق یکی از اهداف اصلی شرکت‌های بیمه است؛ پیدا کردن افرادی که احتمال می‌رود از خدمات بیمه استفاده کنند، بسیار مهم است و منجر به مدیریت هرچه بهتر سرمایه و هزینه‌ها می‌شود. در این پژوهش با استفاده از یک مجموعه داده، حاوی اطلاعات مختلف از جمله مشخصات مشتری، اطلاعات فروش، پوشش‌های خاص، موقعیت جغرافیایی و وضعیت اموال، به ارائه مدلی می‌پردازیم که بتواند احتمال اینکه که فرد از خدمات بیمه استفاده می‌کند یا خیر، را پیش‌بینی کند. برای ارائه این مدل تکنیک‌ها و الگوریتم‌های مختلف داده‌کاوی را بکار گرفته‌ایم. پس از پیش‌پردازش‌های لازم بر روی داده‌ها، از الگوریتم‌های جنگل تصادفی، درخت تصمیم و بیزین ساده برای ساخت مدل استفاده کردیم. الگوریتم جنگل تصادفی با بیش از 99 درصد صحت بهترین عملکرد را نشان می‌دهد. نتایج حاصل از این پژوهش نشان می‌دهد که به‌کارگیری تکنیک‌های داده‌کاوی برای پیش‌بینی مشتریان احتمالی می‌تواند به شرکت‌های بیمه کمک کند تا مشتریان احتمالی خود را پیدا و برای فروش خدمات خود بر روی آن افراد سرمایه‌گذاری کنند. این امر باعث کاهش هزینه‌ها از جمله کاهش هزینه تبلیغات می‌شود.

### واژه‌های کلیدی:

مشتری بیمه، داده‌کاوی، جنگل تصادفی، درخت تصمیم، بیز ساده.



### 1- مقدمه

در جهانی که پر است از مخاطرات و حوادث غیرقابل پیش‌بینی، سرمایه‌های مادی افراد دائماً در معرض تهدیدات جدی قرار دارند و این امر ضرورت استفاده از بیمه را دوچندان می‌کند. همین امر منجر به رشد روزافزون شرکت‌های بیمه و ارائه خدمات متنوع شده است. بازاریابی و یافتن مشتری یکی از فعالیت‌های اصلی شرکت‌های بیمه می‌باشد [1].

کسب‌وکار در بستر وب مجموعه داده‌های گول‌پیکری را ایجاد می‌کند؛ از جمله آن‌ها می‌توان به فروش، معاملات، سهام سوابق تجاری، توضیحات محصول، تبلیغات فروش، مشخصات شرکت‌هایی مانند بیمه، عملکرد و بازخورد مشتریان آن‌ها اشاره کرد. مدیران کسب‌وکارهای مختلف می‌دانند که جمع‌آوری اطلاعات و داده‌ها از مشتریان و مخاطبان یک شرکت، یکی از عوامل مورد نیاز برای رشد و توسعه این شرکت‌ها است؛ اما این فقط بخش اول راهی است که باید طی شود. اگر داده‌ها فقط جمع شوند و بلااستفاده بمانند، عملاً هدف اصلی از جمع‌آوری این اطلاعات، برآورده نشده است. گام مهمی که پس از جمع‌آوری داده‌ها می‌بایست برداشته شود، استخراج دانش از این داده‌ها است و این مهم با استفاده از تکنیک‌های داده‌کاوی حاصل خواهد شد [2].

در این مقاله قصد داریم با به‌کارگیری تکنیک‌های داده‌کاوی بر داده‌های جمع‌آوری‌شده از یک شرکت بیمه بزرگ به دانشی دست‌یابیم که برای شرکت‌های بیمه بسیار حائز اهمیت است. هدف اصلی ما یافتن مشتریان احتمالی که قصد استفاده از خدمات بیمه را دارند، می‌باشد [3]. در ادامه، ابتدا به پیشینه‌ی پژوهش‌های انجام شده در حوزه‌ی بیمه با استفاده از تکنیک‌های داده‌کاوی، پرداخته خواهد شد. پس از آن، داده‌های مورد استفاده در این پژوهش معرفی خواهند شد. در بخش 4 ابعاد مختلف روش مورد استفاده در این پژوهش (پیش‌پردازش‌ها، مدل‌سازی و ارزیابی، و نتایج حاصل‌شده) ارائه می‌شود. بخش پایانی مقاله، بحث و نتیجه‌گیری این پژوهش را در بردارد.

### 2- پیشینه پژوهش

در بازار بیمه نیز همانند بسیاری از حوزه‌های دیگر، داده‌کاوی نفوذ کرده است. پژوهشگران بسیاری با به‌کارگیری تکنیک‌های داده‌کاوی و ارائه‌ی نتایج به مسئولین امر در حوزه‌ی بیمه به رشد روزافزون بیمه‌ها کمک قابل‌توجهی کرده‌اند. در این بخش به‌طور خلاصه به برخی از برجسته‌ترین پژوهش‌های انجام شده این حوزه اشاره می‌کنیم.

کِلداگ مله<sup>2</sup> و کنیتو اسیوتی<sup>3</sup> در پژوهشی که نتایج آن را در سال 2012 منتشر کردند به بررسی داده‌های بیمه درمانی پرداختند. هدف این پژوهش ارائه‌ی روشی برای تشخیص تقلب با استفاده از داده‌کاوی در بیمه‌های درمانی بوده است. برای این هدف روش‌های تشخیص آنومالی و ماشین بردار پشتیبان<sup>4</sup> را بکار گرفتند. تمرکز این پژوهش بر شناسایی افرادی که با فریب عمدی یا بدرفتاری برای به دست آوردن برخی از مزایا، در قالب هزینه‌های بهداشتی تقلب می‌کنند، بوده است. با نتایج حاصل از این مدل، بازرسان بیمه می‌توانند تحقیقات بیشتری را برای افرادی که توسط این مدل‌ها شناسایی شده‌اند، انجام دهند [4].

لینگ<sup>5</sup> و همکارانش در پژوهشی (سال 1998) به بررسی داده‌های بانک، بیمه و شرکت بونس<sup>6</sup> پرداختند. هدف این پژوهش بازاریابی غیرمستقیم و دادن راه‌حل برای حل مشکلات بازاریابی با استفاده از تکنیک‌های داده‌کاوی بوده است. برای رسیدن به این هدف روش‌های بی‌زین ساده<sup>7</sup>، نزدیک‌ترین همسایه<sup>8</sup> و شبکه عصبی<sup>9</sup> را بکار گرفتند. مدل ارائه شده در این پژوهش دارای صحت بیش از 70 درصد می‌باشد [5].

<sup>2</sup> Kirlidog Melih

<sup>3</sup> Cuneyt Asuk

<sup>4</sup> Support Vector Machine(SVM)

<sup>5</sup> Ling

<sup>6</sup> Bonus

<sup>7</sup> Naive Bayes

<sup>8</sup> Nearest neighbour

<sup>9</sup> Neural network



اسمیت کیت<sup>10</sup> و همکارانش در پژوهشی که آن را در سال 2000 منتشر کردند به بررسی داده‌های بیمه پرداختند. هدف این پژوهش تجزیه و تحلیل الگوهای حفظ بیمه با استفاده از تکنیک‌های داده‌کاوی بوده است. در این پژوهش برای نیل به این هدف، پروسه کشف دانش را در یک چارچوب جامع با استفاده از آزمون فرضیه‌ها<sup>11</sup>، خوشه‌بندی، درخت تصمیم<sup>12</sup> و شبکه عصبی انجام دادند. مدل ارائه شده در این پژوهش راهکاری را به شرکت‌های بیمه ارائه می‌کند تا بتوانند مشتریان خود را حفظ کنند [6].

چو<sup>13</sup> و نای<sup>14</sup> (سال 2003) به بررسی داده‌های بیمه پرداختند. هدف این پژوهش ارائه روشی برای انتخاب عوامل مؤثر در فروش بیمه بوده است. برای رسیدن به این هدف روش‌های تجزیه و تحلیل محرمانه، درخت تصمیم و شبکه عصبی را بکار گرفتند. آن‌ها با پیش‌بینی طول مدت خدمات، حق فروش بیمه و شاخص‌های پایداری عوامل بیمه به مدیران کمک می‌کنند تا عوامل اصلی مؤثر در فروش بیمه را شناسایی کنند [7].

### 3- داده‌ها

داده‌های مورد استفاده در این پژوهش متعلق به شرکت بیمه هوم‌سایت<sup>15</sup> می‌باشد. این شرکت خدمات مختلف بیمه همانند بیمه سرقت، آتش‌سوزی منزل و بیمه درمانی را برای مشتریان خود ارائه می‌کند و برای بازاریابی هرچه بهتر و یافتن مشتریان بیشتر اقدام به جمع‌آوری داده‌های مربوط به افراد مختلف کرده است. این مجموعه داده در سایت کگل<sup>16</sup> موجود است و دارای 130376 نمونه و برای هر نمونه 299 ویژگی می‌باشد. ویژگی‌های این داده‌ها از نوع عددی<sup>17</sup> و اسمی<sup>18</sup> می‌باشند [3].

برچسب مربوط به هر یک از نمونه‌های این مجموعه داده نشان‌دهنده‌ی این است که آیا فرد اقدام به خرید بیمه می‌کند یا خیر. تعداد افرادی که اقدام به خرید بیمه کرده‌اند 24089 نفر و تعداد افرادی اقدام به خرید بیمه نمی‌کنند 106287 نفر می‌باشد. در هنگام مدل‌سازی از اعتبارسنجی متقابل<sup>19</sup> 10 بخشی<sup>20</sup> استفاده می‌کنیم؛ 0.9 از نمونه‌ها برای آموزش<sup>21</sup> و از 0.1 باقی‌مانده برای آزمایش<sup>22</sup> استفاده می‌شوند و این عمل 10 مرتبه اجرا شده و در هر بار اجرا یکی از 10 قسمت برای آزمایش و باقی نمونه‌ها برای آموزش در نظر گرفته می‌شود (جدول 1) [3].

جدول 1- مشخصات داده‌های شرکت بیمه مورد استفاده در این پژوهش

ویژگی‌ها	نمونه	داده‌ها
299	130376	کل داده‌ها
299	117338	آموزش
298	13038	آزمایش

<sup>10</sup> Smith Kate  
<sup>11</sup> Hypothesis testing  
<sup>12</sup> PART  
<sup>13</sup> Cho  
<sup>14</sup> Ngai  
<sup>15</sup> Homesite  
<sup>16</sup> kaggle  
<sup>17</sup> Numeric  
<sup>18</sup> Nominal  
<sup>19</sup> Cross validation  
<sup>20</sup> 10 fold  
<sup>21</sup> Train  
<sup>22</sup> Test



#### 4- روش انجام کار

##### 4-1- پیش پردازش

یکی از مهم ترین مراحل انجام یک پژوهش داده کاوی، پیش پردازش داده ها می باشد. داده های ما داری تعدادی مقادیر از دست رفته<sup>23</sup> هستند. دو ویژگی از مجموع 299 ویژگی این مجموعه داده، دارای مقادیر از دست رفته می باشند. از آنجا که مقادیر بیش از 90 درصد نمونه ها، برای این دو ویژگی یا از دست رفته هستند و یا پرتکرارترین مقدار ثابت را دارند (به عبارت دیگر تنوع مقادیر در این دو ستون پایین است) از این 2 ویژگی در ساخت مدل ها صرف نظر می کنیم. پایین بودن تنوع مقادیر در یک ویژگی باعث می شود که آن ویژگی نتواند در ساخت مدل مؤثر واقع شود [2]. با در نظر داشتن این نکته تعداد 2 ویژگی که دارای مقدار ثابتی هستند نیز از مجموعه ی داده ها حذف شدند. یکی از ویژگی های موجود در این داده ها، از نوع تاریخ می باشد. برای اینکه بتوانیم از این ویژگی در ساخت مدل به شکل مؤثرتری استفاده کنیم به جای ویژگی تاریخ دو ویژگی روز و ماه را از آن استخراج کردیم.

به کمک محاسبه ی مقادیر همبستگی<sup>24</sup> می توان از چگونگی ارتباط میان ویژگی ها مطلع شد. هر چه این مقادیر بیشتر باشند همبستگی میان ویژگی ها بیشتر است و به این معنا است که این دو ویژگی حاوی اطلاعات تقریباً یکسانی هستند و رفتاری مشابه دارند [2]. از بین ویژگی هایی که دارای همبستگی بالای 0/8 با یکدیگر هستند، آن ویژگی که همبستگی کمتری با ستون هدف دارد را حذف می کنیم. علاوه بر این، ویژگی هایی که مقدار همبستگی آن ها با ستون هدف کمتر از 0/01 بودند، نیز حذف شدند. با انجام این پردازش، 30 ویژگی از مدل سازی ما حذف شدند. لازم به ذکر است که این مقادیر آستانه به صورت تجربی برای این پژوهش حاصل شده اند.

نویزها<sup>25</sup> به دلایل مختلفی، چون بروز اشتباه در فرآیند جمع آوری و یا وارد کردن اطلاعات به سیستم بوجود می آید؛ شناسایی آن ها به ما کمک می کند تا با آگاهی و دقت بیشتری به طراحی مدل بپردازیم [2]. پس از یافتن نمونه های نویزدار، آن ها را حذف کردیم؛ این کار باعث شد که تعداد 20 ویژگی دارای مقادیری تقریباً ثابت شوند، همان طور که پیش از این اشاره شده ویژگی های دارای مقادیر ثابت در مدل سازی مؤثر نیستند در نتیجه از این ویژگی ها صرف نظر شد.

لازم به ذکر است که در پایان مرحله ی پیش پردازش، در مجموع 245 ویژگی برای ساخت مدل باقی مانده است.

##### 4-2- مدل سازی و ارزیابی

در این پژوهش برای مدل سازی از سه الگوریتم معروف داده کاوی استفاده کرده ایم. این الگوریتم ها عبارتند از:

- درخت تصمیم، یکی از الگوریتم های پر کاربرد داده کاوی است. در این الگوریتم که از نوع تکرار شونده و حریصانه است، یک درخت تصمیم با یافتن بهترین ویژگی جداکننده در هر گام و سپس جداسازی متوالی داده ها به گروه های مجزا، ساخته می شوند [2].
- بیزین ساده، بر پایه ی قضیه بیز برای مدل سازی پیش گوینه ارائه شده است. قضیه بیز از روشی برای دسته بندی پدیده ها بر پایه احتمال وقوع یا عدم وقوع یک پدیده استفاده می کند و احتمال رخ دادن یک پدیده محاسبه و دسته بندی می شود [2].
- جنگل تصادفی<sup>26</sup>، یک الگوریتم یادگیری ماشین است که اغلب اوقات نتایج بسیار خوبی را حتی بدون تنظیم فرآیندهای آن، فراهم می کند. این الگوریتم درخت های تصمیم زیادی را تولید می کند. هر درخت یک طبقه بند مجزا است و نتیجه ی نهایی در جنگل تصادفی با رأی گیری از این درخت ها حاصل می شود. رأی گیری از بین چندین طبقه بند باعث می شود که مدل نهایی در برابر نویز مقاوم باشد [2].

<sup>23</sup> Null

<sup>24</sup> Correlation

<sup>25</sup> Noise

<sup>26</sup> Random forest



برای ارزیابی و مقایسه‌ی مدل‌ها از معیارهای صحت<sup>27</sup> - بازخوانی<sup>28</sup> - دقت<sup>29</sup> - معیار اف<sup>30</sup> استفاده کرده‌ایم:

- صحت: اولین معیار یا سنج‌های که به ذهن می‌رسد، معیار صحت یا میزان تشخیص درست مدل است؛ یعنی نسبت تشخیص‌های درست (مثبت واقعی<sup>31</sup> + منفی واقعی<sup>32</sup>) به کل داده‌ها. این معیار برای ارزیابی مدل‌ها در زمان استفاده از داده‌های نامتوازن (یعنی تفاوت زیادی در تعداد نمونه‌های دسته‌ها وجود دارد) کافی نیست؛ زیرا این عدم توازن باعث می‌شود مدل‌های متمایل به دسته پرتعداد، شناسایی نشوند و به اشتباه یک مدل بد، مدل خوبی معرفی شود.
- بازخوانی: ارزیابی اینکه کل نمونه‌های واقعاً مثبت شامل نمونه‌هایی است که درست، مثبت شناسایی شده‌اند (مثبت واقعی) و نمونه‌هایی که مثبت بوده‌اند اما نادرست، منفی شناسایی شده‌اند (منفی اشتباه<sup>33</sup>). در این معیار بر تعداد نمونه‌های مثبت شناسایی شده به کل نمونه‌های مثبت تمرکز می‌شود.
- دقت: در کنار معیار بازخوانی معیار دیگری را به نام دقت، برابر تعداد نمونه‌های تشخیصی مثبت واقعی به کل نمونه‌های مثبت اعلام شده تعریف می‌شود تا میزان مثبت‌های اشتباه<sup>34</sup> هم در نظر گرفته شود.
- معیار اف: میانگین هارمونیک بازخوانی و دقت می‌باشد [2].

### 3-4- نتایج

در بخش‌های پیشین به مشخصات داده‌های این پژوهش، پیش‌پردازش‌های انجام شده بر روی داده‌ها، الگوریتم‌های مورد استفاده برای ساخت مدل و معیارهای ارزیابی در نظر گرفته شده اشاره شد. پس از پیاده‌سازی و انجام آزمایش‌ها نتایج قابل‌توجهی حاصل شد. خلاصه‌ای از نتایج این آزمایش‌ها در جدول 2 ارائه شده است. همان‌طور که در این جدول مشخص است، الگوریتم جنگل تصادفی در همه‌ی معیارهای ارزیابی توانسته است نتایج بهتری را ارائه کند. جهت ارائه‌ی جزئیات بیشتر، ماتریس سردرگمی<sup>35</sup>، الگوریتم جنگل تصادفی در جدول 3 ارائه شده است.

جدول 2- مقایسه الگوریتم‌های مختلف معیارهای ارزیابی صحت، بازخوانی، دقت و معیار اف

معیار اف	دقت	بازخوانی	صحت	
0,997648	0,999263	0,996039	0,996172	جنگل تصادفی
0,916715	0,955845	0,880663	0,869968	بیزین ساده
0,971946	0,958397	0,985885	0,953703	درخت تصمیم

جدول 3- نتیجه مدل‌سازی با الگوریتم جنگل تصادفی

	0	1	جمع
0	105866	421	106287
1	78	24011	24089
جمع	105944	24432	130376

<sup>27</sup> Accuracy

<sup>28</sup> Recall

<sup>29</sup> Precision

<sup>30</sup> F-measure

<sup>31</sup> True Positive (TP)

<sup>32</sup> True Negative (TN)

<sup>33</sup> False Negative (FN)

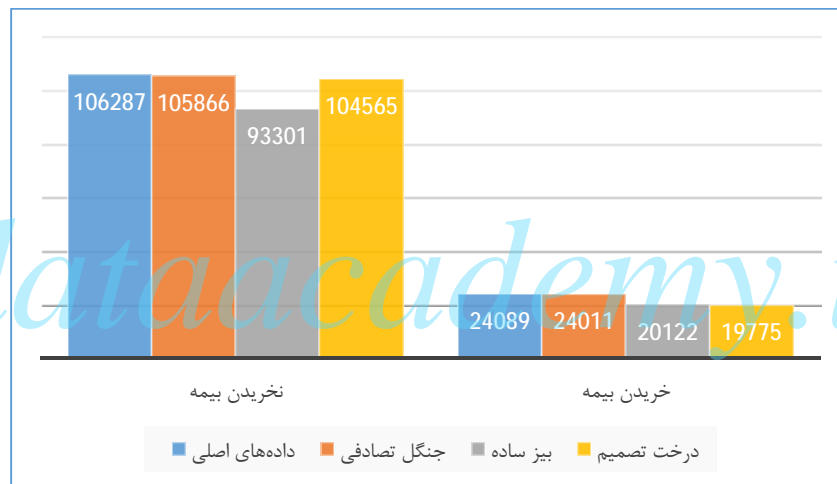
<sup>34</sup> False Positive (FP)

<sup>35</sup> Confusion matrix



## 5- بحث و نتیجه گیری

هدف پژوهش ما، ارائه‌ی یک مدل پیشگو برای مشتریان احتمالی استفاده از خدمات بیمه می‌باشد. در سال‌های اخیر صنعت بیمه رشدی چشمگیر داشته است و بازاریابی موفق یکی از اهداف اصلی شرکت‌های بیمه است؛ پیدا کردن افرادی که احتمال می‌رود از خدمات بیمه استفاده کنند، بسیار مهم است و منجر به مدیریت هرچه بهتر سرمایه و هزینه‌ها می‌شود. نتایج حاصل از آزمایش‌های ما نشان داد که تکنیک‌های داده‌کاوی می‌توانند رویکردی مناسب برای دست‌یابی به چنین دانش سودمندی باشند. داده‌های این پژوهش بسیار نامتوازن هستند اما ما با پیش‌پردازش‌هایی مناسب و استفاده از الگوریتم‌هایی کارا توانستیم به مدل پیش‌گویی با صحت بیش از 99 درصد دست‌یابیم. نمودار ارائه شده در شکل 1 فروانی مثبت صحیح (افراد خریدار بیمه) و منفی صحیح (افراد که بیمه نمی‌خرند) را برای الگوریتم‌های مختلف و داده‌های اصلی با هم مقایسه می‌کند. همان‌طور که در بخش پیشین اشاره شد، الگوریتم جنگل تصادفی توانسته است بهترین نتایج را ارائه دهد. با استفاده از نتایج این پژوهش مسئولین بیمه می‌توانند، مشتریان احتمالی خود را پیدا و برای فروش خدمات خود بر روی آن افراد سرمایه‌گذاری کنند.



شکل 1: مقایسه‌ی فروانی دو کلاس خریدن/نخریدن بیمه برای الگوریتم‌های مختلف و داده‌های اصلی

## مراجع

- [1] J. Hsia, E. Kemper, C. Kiefe, J. Zapka, S. Sofaer, M. Pettinger, D. Bowen, M. Limacher, L. Lillington, E. Mason, and Women's Health Initiative Investigators, "The importance of health insurance as a determinant of cancer screening: evidence from the Women's Health Initiative," *Preventive medicine*, vol. 3, no. 31, pp. 261-270, 2000.
- [2] J. Han, J. J. Pei, and M. Kamber, . Data mining: concepts and techniques, 3rd ed., Morgan Kaufmann, 2011.
- [3] Kaggle Inc, "kaggle," 1 April 2001. [Online]. Available: <https://www.kaggle.com>. [Accessed 24 January 2019].
- [4] M. Kirlidog, and C. Asuk, , "A fraud detection approach with data mining in health insurance," *Procedia-Social and Behavioral Sciences*, vol. 62, pp. 989-994, 2012.
- [5] C. X. Ling, and C. Li., "Data mining for direct marketing: Problems and solutions," in *Kdd*, Newyork, 1998.
- [6] K.A. Smith, R.J. Willis, and M. Brooks, , "An analysis of customer retention and insurance claim patterns using data mining: A case study," *Journal of the operational research society*, vol. 5, no. 51, pp. 532-541, 2000.
- [7] V. Cho, and E.W. Ngai., "Data mining for selection of insurance sales agents," *Expert systems*, vol. 3, no. 20, pp.



دومین کنفرانس ملی کامپیوتر ، فناوری اطلاعات و

کاربردهای هوش مصنوعی

1 اسفند ماه 1397 - دانشگاه شهید چمران اهواز



123-132, 2003.

## Finding Potential Insurance Customers Using Data-mining Techniques

**Amir Oliaei**, Undergraduate student, Kashmar Higher Education Institute

**Ali Hazeri**, Undergraduate student, Kashmar Higher Education Institute

**Mohammad Jamali**, Undergraduate student, Kashmar Higher Education Institute

**Atefeh Khazaei**, PhD in Software Engineering, Kashmar Higher Education Institute

### Abstract

In recent years, the insurance industry has grown dramatically, with various companies and services. The successful marketing is one of the main goals of insurance companies; finding people who would like to use insurance services is very important and leads to better management of the fund and expense. In this research, using a dataset containing various information such as customer characteristics, sales information, specific coverage, geographic location and status of properties, we present a model that predicts the probability of buying insurance by persons. To provide this model, we have applied various data-mining techniques and algorithms. After data preprocessing, we used Random Forest, Decision Trees, and Naïve Bayes algorithms to construct the model. Random Forest algorithm shows the best performance with more than 99% accuracy. The results of this research show that the data-mining techniques can be useful to predict potential customers and help insurers to find their customers and invest on them to sell their services. This model reduces the costs, such as the advertising cost.

**Keywords:** Insurance customer, Data-mining, Random Forest, Decision Tree, Naïve Bayes.