



Designing a Hybrid Model for Classification of Imbalanced Data in the Field of Third Party Insurance

Mahnaz Manteqipour^{1*}, Parisa Rahimkhani²

1.PhD., Applied Mathematics (Operation research), Azarbaijan Shahid Madani University, Tabriz, Iran.

(Corresponding Author) mahteqipoor@gmail.com

2. PhD., Numerical Analysis, Alzahra University, Tehran, Iran.

Rahimkhani.parisa@gmail.com

Abstract

Introduction: The major part of Iran's insurance industry portfolio is the field of compulsory civil liability insurance of motor vehicle owners against third parties. Therefore, detecting the behavior of this insurance field will be effective in order to provide better services to the customers of the insurance industry. Predicting the claim rates for insurance policies, based on the features saved for each insurance policy, is one of the problems of the insurance industry that can be solved with the help of data mining techniques. Insurance is designed using the law of large numbers. In simpler words, a sufficient number of insurance policies are issued, and a small part of this number of insurance policies deal with claims. From the sum of the issued insurance premiums, the cost of claims will be compensated. Therefore, the insurance industry is faced with imbalanced data. The imbalances of insurance industry data causes many challenges in data classification. In the field of third-party insurance and in the data set of this research, there are 14 features for every policies and the data imbalance ratio is 1 to 0.0092, which is considered severe imbalanced.

Method: In this research, we deal with the classification of severe imbalanced data in the field of third party insurance. To overcome the problem of imbalanced data, two hybrid models with different architectures based on 5 basic Gaussian Bayes models, support vectors, logistic regression, decision tree and nearest neighbor are designed. First proposed hybrid model is using random sampling from whole dataset and applying a resampling method for classification and second one selects samples from each labels separately and apply a classification model on the whole selected data. The results of these models are compared.

Results:The obtained results show that the proposed hybrid models can predict the occurrence or non-occurrence of traffic accidents better than other data mining algorithms. The popular measures such as precisions and recalls of two proposed hybrid models show that second hybrid model has higher performance. And in ensemble phase, the number of models in simple voting as a hyper parameter can be adjusted based on the company's strategy. Also, the use of decision tree to ensemble basic models to build a combined model provides better results than simple voting of basic models.

Discussion: To do more research on the problem of imbalance data classification more complicated resampling data algorithms could be applied and the results be compared.

Keywords: hybrid model, imbalance data, data mining, third party insurance.

طراحی مدل ترکیبی برای طبقه‌بندی داده‌های نامتوازن در رشته بیمه شخص ثالث

سال سوم، تابستان ۱۴۰۱
شماره دوم، صص: ۹-۱

تاریخ دریافت: ۱۴۰۱/۰۲/۳۱
تاریخ پذیرش: ۱۴۰۱/۰۴/۱۶

مهناز منطقی پور^{۱*}، پریسا رحیم‌خانی^۲

۱. دکترا، ریاضی کاربردی (تحقیق در عملیات)، دانشگاه شهید مدنی آذربایجان، تبریز، ایران. (نویسنده مسئول) mahteqipoor@gmail.com

۲. دکترا، آنالیز عددی، دانشگاه الزهراء، تهران، ایران Rahimkhani.parisa@gmail.com

چکیده

بخش عمده پور تفوی صنعت بیمه کشور ایران را رشته بیمه اجباری مسئولیت مدنی دارندگان وسایل نقلیه موتوری زمینی در مقابل اشخاص ثالث، تشکیل داده است. توانایی پیش‌بینی وقوع و یا عدم وقوع خسارت به‌ویژه خسارت‌های جانی نه تنها برای شرکت‌های بیمه بلکه برای تصمیم‌گیرندگان در حوزه‌های افزایش امنیت جاده‌ها اهمیت بسیاری دارد. به منظور پیش‌بینی برجسب وقوع یا عدم وقوع خسارت از روش‌های طبقه‌بندی استفاده می‌شود که در واقع یک مسئله طبقه‌بندی نامتوازن است. این نامتوازن بودن شدید، ناشی از ماهیت کسب و کار بیمه است. نامتوازن بودن داده‌های صنعت بیمه باعث ایجاد چالش‌های بسیاری در تجزیه و تحلیل داده‌های مربوطه می‌شود. در این پژوهش، ما به طبقه‌بندی داده‌های نامتوازن بیمه شخص ثالث در یک شرکت بیمه معتبر می‌پردازیم. در این راستا دو روش ترکیبی برای رفع مشکل نامتوازن بودن داده‌ها بر اساس ۵ مدل پایه گاوسین بیز، بردارهای پشتیبان، لجستیک رگرسیون، درخت تصمیم، نزدیکترین همسایگی به منظور طبقه‌بندی مؤثرتر داده‌های مربوطه ارائه می‌شود. نتایج به‌دست‌آمده نشان می‌دهد که مدل‌های ترکیبی ارائه‌شده بهتر از سایر الگوریتم‌های داده‌کاوی برای داده‌های مربوطه جواب می‌دهند و استفاده از درخت تصمیم در تجمیع مدل‌های پایه برای ساخت مدل ترکیبی نسبت به رأی‌گیری ساده مدل‌ها نتایج بهتری ارائه می‌کند. همچنین ابر پارامتر تعداد مدل‌های لازم در رأی‌گیری بر اساس استراتژی شرکت قابل تنظیم است. تعداد ویژگی‌های ثبت‌شده از بیمه‌نامه‌ها در شرکت‌های بیمه محدود است با تکمیل این ویژگی‌ها به ویژه اضافه‌شدن سوابق رانندگی و سایر ویژگی‌های فردی می‌توان به مدل بهتری دست یافت.

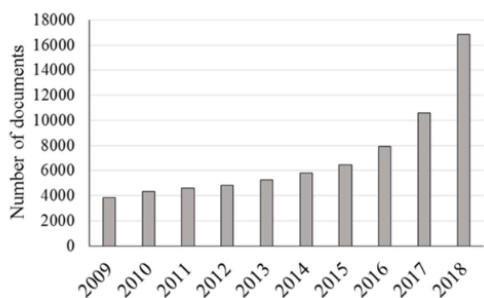
واژه‌های کلیدی: مدل ترکیبی، داده‌های نامتوازن، داده‌کاوی، بیمه شخص ثالث.

۱. راهبر میز داده‌کاوی، گروه بیمه الکترونیک، پژوهشکده بیمه، ایران.

۲. پژوهشگر میز داده‌کاوی، گروه بیمه الکترونیک، پژوهشکده بیمه، ایران.

۱. مقدمه

است. نمودار زیر تعداد متون منتشر شده در مورد این مدل‌ها را نمایش می‌دهد [3].



شکل ۱: افزایش تعداد مقالات مدل‌های ترکیبی در سال‌های اخیر

یک روش ترکیبی در مقاله [4] ارائه شد که در آن عدم توازن داده‌ها با استفاده از ترکیب روش کا-معکوس نزدیکترین همسایه و روش ماشین بردار پشتیبان مدیریت شده است. نویسندگان در مقاله [5] یک روش ترکیبی بر اساس روش SMOTE، اعتبارسنجی متقابل و جنگل تصادفی برای کشف تقلب در زمینه بیمه درمان ارائه دادند. نویسندگان در مقاله [6] روی بهینه‌سازی یادگیرنده برای کشف تقلب در رشته بیمه درمان با استفاده از یک روش نمونه‌گیری مجدد تمرکز کردند.

۲. مسئله پژوهش

هدف این مقاله ساختن مدلی برای پیش بینی خسارت‌های جانی رشته بیمه شخص ثالث است. این نوع از خسارت‌ها فراوانی اندکی دارند و در عین حال به دلیل آنکه منجر به پرداخت دیه می‌گردند هزینه اقتصادی بالایی ایجاد می‌کنند. علاوه بر هزینه‌های مالی، هزینه‌های جانی غیر قابل جبران هستند و تبعات بسیاری برای مردم جامعه ایجاد می‌کنند. از این رو پیش بینی این نوع از خسارت‌ها برای تصمیم‌گیرندگان بسیار مفید خواهد بود.

در این پژوهش، پس از استفاده از روش‌های نمونه‌گیری مجدد به منظور حل چالش عدم توازن داده‌ها، مدل‌های درخت تصمیم، ماشین بردار، رگرسیون لجستیک، نزدیکترین کا-همسایگی، مدل گاوسین بیز و جنگل تصادفی برای پیش‌بینی برچسب وقوع و یا عدم وقوع خسارت مورد استفاده قرار گرفتند و همگی عملکرد ضعیفی داشتند. این مدل‌ها در مرحله تست، عمدتاً نمونه‌ها را متعلق به کلاس اکثریت یعنی کلاس بدون خسارت پیش‌بینی کردند. از این رو به ساختن مدل‌های ترکیبی روی آوردیم. مدل‌های ترکیبی با هدف تقویت نتایج مدل‌های ساده‌تر ایجاد شده‌اند. ایده مدل‌های ترکیبی آن است که از چندین مدل پیش‌بینی کننده به عنوان مدل پایه به منظور دستیابی به مدلی با عملکرد بهتر استفاده می‌شود و مدل نهایی با استفاده از روش‌های متفاوت، جهت جمع‌بندی نتایج مدل‌های پایه، ارائه می‌گردد. در این مقاله با تغییر در ساختار مدل ترکیبی به دنبال مدلی با نتایج قوی‌تر بوده‌ایم.

با پیشرفت‌های تکنولوژی و امکان ذخیره‌سازی حجم عظیم اطلاعات، استفاده از داده‌کاوی به عنوان یک فرآیند کشف الگوهای گوناگون در داده‌ها در حال افزایش است [1] و به یک موضوع مهم در حوزه‌های مختلف از جمله صنعت بیمه تبدیل شده است. رشته بیمه اجباری مسئولیت مدنی دارندگان وسایل نقلیه موتوری زمینی در مقابل اشخاص ثالث که در این مقاله به اختصار رشته بیمه شخص ثالث به آن اطلاق می‌شود، به لحاظ تبعات اجتماعی و اقتصادی نه تنها برای صنعت بیمه که برای آحاد جامعه بسیار با اهمیت است. بر اساس سالنامه آماری ۱۳۹۹ صنعت بیمه، این رشته حدود ۳۳/۶ درصد از کل حق بیمه‌های دریافتی و ۴۰ درصد از کل خسارت‌های پرداختی صنعت بیمه و ضریب خسارت ۱۰۸٫۸ درصد را به خود اختصاص داده است. از این رو نرخ‌گذاری هوشمند و منصفانه در این رشته بیمه‌ای حائز اهمیت است. از آنجا که داده‌کاوی به کاربردها و اهداف موردنظر محقق، مدیر و یا ناظر توجه دارد، تشریح کاربردها و به‌کارگیری آن در تصمیم‌سازی‌های صنعت بیمه مفید و ضروری است.

مهمترین چالش مواجهه با مسئله مدل‌سازی داده‌های صنعت بیمه، نامتوازن بودن داده‌ها است. داده‌هایی که تعداد اعضای کلاس‌های آن متوازن نیست را مسئله طبقه‌بندی نامتوازن گویند [2]. مسائل طبقه‌بندی بر اساس برچسب ادعای خسارت در صنعت بیمه با توجه به ماهیت کسب و کار بیمه بازرگانی منجر به مسائل طبقه‌بندی نامتوازن می‌گردند. اغلب الگوریتم‌های داده‌کاوی در مواجهه با داده‌های نامتوازن عملکرد خود را از دست می‌دهند و تمام داده‌ها را متعلق به کلاس با برچسب اکثریت پیش‌بینی می‌کنند. از این رو چندین راهکار برای حل این مسائل در مقالات متعدد ارائه شده است.

به منظور طبقه‌بندی داده‌های نامتوازن، دو رویکرد اصلی وجود دارد: ۱- روش‌های نمونه‌گیری مجدد ۲- روش‌های حساسیت هزینه. روش‌های نمونه‌گیری مجدد به دنبال دستیابی به مجموعه داده‌های متوازن هستند و روش‌های حساسیت هزینه، حساسیت تابع زیان مسئله طبقه‌بندی را برای خطای کلاس اقلیت بالایی برند تا مینیمم‌سازی تابع زیان منجر به پیش‌بینی همه نمونه‌ها به کلاس حداکثر نشود.

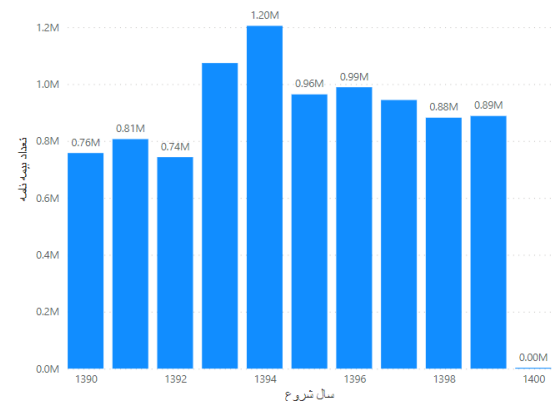
مدل‌های مرسوم یادگیری ماشین با گذر زمان توسعه یافته‌اند و مدل‌های ترکیبی و گروهی خلق شده‌اند. جنگل تصادفی به‌واقع مدلی ترکیبی است که چندین درخت را ترکیب می‌نماید. این روش‌ها به منظور افزایش قدرت مدل‌ها از جنبه‌های محاسباتی، عملکرد، استواری و قدرت مدل‌ها خلق شده‌اند. در مدل‌های ترکیبی، خروجی یک مدل جهت بهبود به مدل بعدی ارسال می‌شود و مدل قوی‌تر به این شکل ساخته می‌شود و در مدل‌های گروهی، مدل‌های مختلف به صورت مجزا ساخته می‌شوند و با استفاده از روش‌های متفاوت خروجی مدل گروهی، تعیین می‌شود. در دهه گذشته استفاده از مدل‌های گروهی و ترکیبی به سرعت افزایش یافته

۳. داده‌های تحقیق

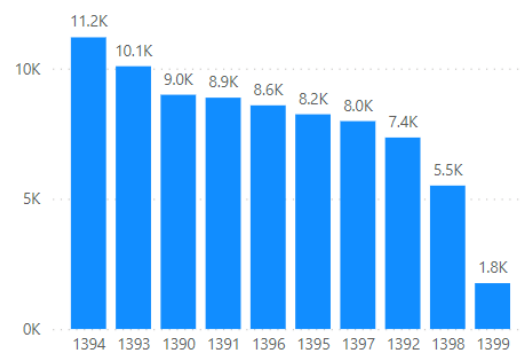
داده‌های این پژوهش مربوط به سال‌های ۱۳۹۰ الی ۱۴۰۰ در رشته بیمه شخص ثالث می‌باشند که متعلق به یک شرکت بیمه معتبر در ایران است.

داده‌ها در قالب یک فایل پشتیبان پایگاه داده SQL server 2019 دریافت شده‌اند و پس از توصیف و بررسی، بصری‌سازی داده‌ها انجام شده است. داده‌های پرت حذف شده‌اند، پاکسازی داده‌ها انجام شده است و داده‌ها جهت به کارگیری الگوریتم‌های یادگیری ماشین آماده شده‌اند. داده‌ها شامل ۹۲۸۳۳۱۶ بیمه‌نامه شخص ثالث و خسارت‌های مربوطه می‌باشند. تعداد بیمه‌نامه‌ها و تعداد خسارت‌ها به تفکیک سال شروع بیمه‌نامه مطابق شکل ۲ و شکل ۳ می‌باشد. این داده‌ها ابتدای سال ۱۴۰۰ از پایگاه داده‌های شرکت استخراج شده است و بنابراین تعداد کمی از بیمه‌نامه‌ها متعلق به سال ۱۴۰۰ می‌باشند.

نرخ نامتوازن بودن داده‌های این پژوهش ۱ به ۰,۰۰۹۲ است. به عبارتی به ازای هر ۱۰۰۰ بیمه‌نامه ۹/۲ بیمه‌نامه با خسارت جانی مواجه شده‌اند.



شکل ۲: تعداد بیمه‌نامه‌ها به تفکیک سال شروع بیمه‌نامه



شکل ۳: تعداد خسارت‌های جانی به تفکیک سال شروع

به دلیل تفاوت رفتاری مشاهده شده در گروه‌های خودرو محاسبات این مقاله بر داده‌های خودروهای سواری و منحصرأ بر خسارت‌های جانی انجام شده است.

۴. مدل‌های ترکیبی

پس از پاکسازی و آماده‌سازی داده‌ها و شناخت داده‌ها، تلاش نمودیم با استفاده از الگوریتم‌های یادگیری ماشین و داده‌کاوی مدلی برای پیش

بینی وقوع خسارت‌ها بسازیم. از آنجا که ظرفیت ویژگی‌ها محدود است، به جای پیش‌بینی نرخ خسارت‌ها ترجیح دادیم وقوع یا عدم وقوع خسارت‌ها را بررسی کنیم.

ویژگی‌های به کاررفته در مدل در جدول زیر آورده شده است:

جدول ۱: ویژگی‌های بیمه‌نامه‌ها

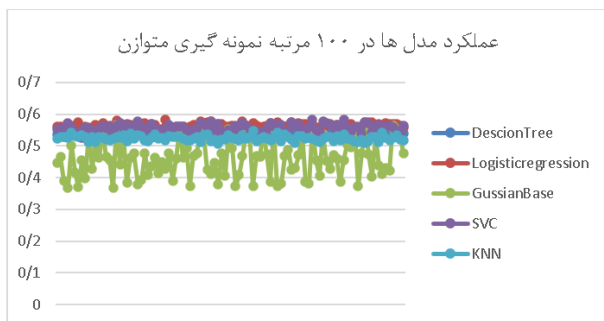
شمارش	نام متغیر
۱	یت
۲	ذآر
۳	ذآر
۴	
۵	
۶	
۷	
۸	
۹	
۱۰	رد
۱۱	ت
۱۲	ت
۱۳	ت
۱۴	در

۱۴ ویژگی ذکر شده در جدول ۱: ویژگی‌های بیمه‌نامه‌ها را به منظور پیش‌بینی وقوع یا عدم وقوع خسارت‌های جانی به کار گرفته ایم. بنابراین قصد داریم برچسبی که مقادیر یک و صفر اختیار می‌کند، پیش‌بینی کنیم.

پس از مرحله پاکسازی داده‌ها [7]، مهمترین چالش این مسئله مواجهه با نامتوازن بودن داده‌ها است.

برای طبقه‌بندی داده‌های مذکور ابتدا به پیاده‌سازی الگوریتم درخت تصمیم، جنگل تصادفی، رگرسیون لجستیک، نایو بیس نمودیم [8]. نامتوازن بودن داده‌ها منجر به پیش‌بینی برچسب ۰ به معنی عدم وقوع خسارت برای تمام بیمه‌نامه‌ها شد که این نتیجه ناشی از ناتوانی مدل‌ها در رویارویی با داده‌های نامتوازن است. زیرا مدل‌ها تلاش بر بیشینه نمودن دقت مدل دارند که با پیش‌بینی برچسب کلاس اکثریت برای تمام نمونه‌ها محقق می‌شود.

سپس از تکنیک‌های مختلف روش‌های نمونه‌گیری مجدد و روش حساسیت هزینه استفاده نمودیم. روش حساسیت هزینه عملکرد مطلوبی نداشت. در روش نمونه‌گیری مجدد ابتدا این روش‌ها را بر داده‌ها به کار بردیم و سپس داده‌ها را به دو بخش آموزش و تست تقسیم نمودیم. مدل‌های حاصل بر داده‌های تست دارای عملکرد ۹۹ درصد به بالا داشتند. اما از آنجا که روش‌های نمونه‌گیری مجدد داده‌ها را تغییر می‌دهند بخشی از داده‌های آموزش در داده‌های تست موجودند و بنابراین تصمیم‌گرفتم داده‌های آموزش و تست را قبل از به کارگیری الگوریتم‌های نمونه‌گیری مجدد انتخاب‌نمائیم و روش‌های نمونه‌گیری مجدد را صرفأ بر داده‌های آموزش اعمال کنیم. با این اقدام نتایج مدل بر داده‌های تست به شدت افت کرد. از این رو بر آن شدیم تا به منظور بهبود



شکل ۵: مقایسه عملکرد مدل ها بر اساس معیار F1

همان طور که مشاهده می شود ضعیفترین پیش بینی ها متعلق به روش گاو سین بیز است. روش نزدیکترین همسایگی، پیش بینی هایی بهتر از پیش بینی های تصادفی دارد و به ترتیب سه مدل رگرسیون لجستیک، بردار پشتیبان و درخت تصمیم بهترین عملکرد را داشته اند. بنابراین مدل ترکیبی را با استفاده از این سه مدل می سازیم. از این رو مجدد مدل را بررسی می کنیم و این بار مدل های پایه را شامل این سه مدل در نظر می گیریم.

عملکرد مدل ترکیبی ۱ با استفاده از سه مدل پایه، در صد مرتبه نمونه گیری منجر به ۳۰۰ مدل شده است که ترکیب آنها با درصد رای گیری ۵۷ درصد، عملکرد زیر را داشته است:

	precision	recall	f1-score	support
0	0.990682	0.988622	0.989651	857761
1	0.011045	0.013482	0.012142	8085
accuracy			0.979516	865846
macro avg	0.500863	0.501052	0.500896	865846
weighted avg	0.981534	0.979516	0.980523	865846

مطابق ماتریس طبقه بندی و معیارهای ارزیابی مدل نمایش داده شده فوق، می توان گفت اندکی نتایج بهتر است اما همچنان ضعیف است. زیرا تنها ۱/۳٪ از کل خسارت ها شناسایی شده اند و از کل خسارت های پیش بینی شده توسط مدل ۱/۱٪ درست بوده اند. عملکرد الگوریتم به ازای انتخاب مدل SVC به عنوان مدل پایه با درصد رای گیری ۵۲٪ به صورت زیر بوده است:

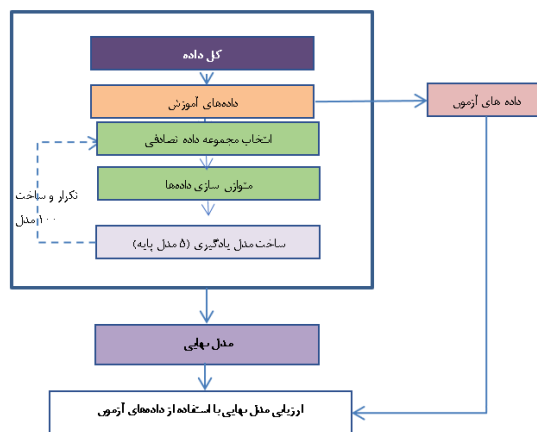
	precision	recall	f1-score	support
0	0.990979	0.677319	0.804663	857761
1	0.010001	0.345826	0.019439	8085
accuracy			0.674224	865846
macro avg	0.500490	0.511572	0.412051	865846
weighted avg	0.981818	0.674224	0.797331	865846

مشاهده می شود که حدود ۳۴،۵٪ از یک های شناسایی شده درست بوده اند. در عین حال ۵۲۸۹ بیمه نامه بدون خسارت به عنوان بیمه نامه های دارای خسارت جانی شناسایی شده اند. بنابراین به کارگیری این مدل با کاهش ۳۲،۳ درصدی پرتفو منجر به کاهش ۳۴،۵٪ از

عملکرد مدل ها، این مدل ها را به صورت ترکیبی محاسبه و دقت را اندازه گیری کنیم.

جهت ساخت مدل ترکیبی از ۵ مدل پایه گاو سین بیز، بردارهای پشتیبان، لجستیک رگرسیون، درخت تصمیم، نزدیکترین همسایگی ۱۰۰ مدل ساخته ایم و دقت مدل ها را محاسبه کرده ایم.

۱.۴. ساخت مدل ترکیبی (۱) انتخاب مجموعه داده تصادفی



شکل ۴: معماری مدل ترکیبی ۱

با انتخاب روش Smote به عنوان روش نمونه گیری مجدد [2]، در مدل ترکیبی ۱، ماتریس طبقه بندی و معیارهای ارزیابی مدل حاصل با درصد رای گیری ۵۰٪ به صورت زیر است:

	precision	recall	f1-score	support
0	0.990818	1.000000	0.995388	857896
1	0.000000	0.000000	0.000000	7950
accuracy			0.990818	865846
macro avg	0.495409	0.500000	0.497694	865846
weighted avg	0.981721	0.990818	0.986249	865846

بر اساس نتایج حاصله بهبودی در عملکرد مدل نسبت به مدل جنگل تصادفی و درخت تصمیم و الگوریتم هایی که ذکر شد، حاصل نشده است. حتی می توان ادعا کرد عملکرد مدل ضعیف تر نیز شده است. با بررسی عملکرد الگوریتم ها متوجه شدیم عملکرد سه الگوریتم از دیگر الگوریتم ها بهتر است در شکل زیر عملکرد پنج الگوریتم قابل مشاهده است.

پس از ساخت ۵۰۰ مدل، عملکرد مدل ها با استفاده از معیار F1-score به صورت زیر بوده است:

خسارت‌های جانی می‌شود که می‌تواند برای شرکت سودآور باشد. اما کاهش ۳۲٫۲ درصدی پرتفو برای کمتر شرکتی قابل قبول است. سایر ترکیب‌های مدل‌های یادگیری از ۵ مدل به عنوان مدل‌های پایه عملکرد ضعیفتری داشته‌اند. از ذکر نتایج آن صرف نظر می‌کنیم. در مدل ترکیبی دوم نحوه انتخاب نمونه داده‌ها را در هر تکرار تغییر داده‌ایم که در بخش بعد توضیح داده شده‌است.

۲.۴. مدل ترکیبی ۲ (با نمونه‌گیری مجزا از هر کلاس)

مراحل انجام کار به صورت زیر طراحی شده‌است:

تبدیل متغیرهای اسمی به متغیرهای عددی
جداسازی داده‌ها به داده‌های آزمون و تست

هم مقیاس نمودن داده‌ها بدون استفاده از داده‌های آزمون و صرفاً تبدیل داده‌های آزمون با استفاده از الگوی هم مقیاسی داده‌های آموزش
جداسازی داده‌های آزمون بر اساس برچسب‌های دارای خسارت جانی،
۱ و بدون خسارت جانی

مراحل ۵ و ۷، ۱۰۰ مرتبه تکرار شده‌اند:

انتخاب ۱۰۰۰۰ نمونه تصادفی با برچسب ۰ و انتخاب ۱۰۰۰۰ نمونه تصادفی با برچسب ۱

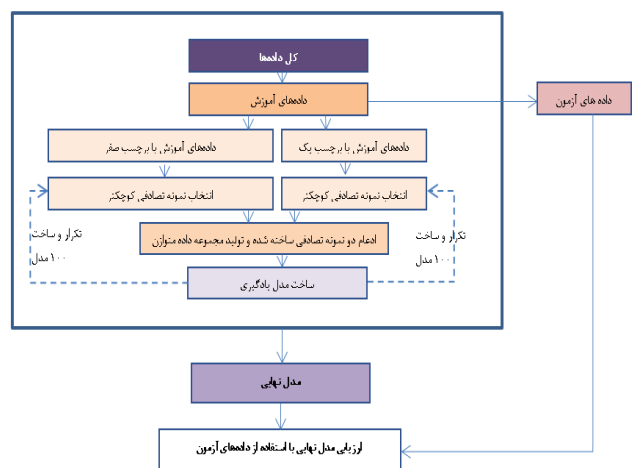
ادغام نمونه‌های برچسب‌ها و ایجاد مجموعه داده نمونه

محاسبه مدل مربوطه بر اساس مجموعه داده نمونه ۲۰۰۰۰ تایی

خروجی مدل نهایی که بر اساس رأی‌گیری تمام مدل‌های ساخته شده در مرحله ۷ حاصل می‌شود.

به عبارتی با فرآیند فوق اگر ۵ مدل یادگیری به عنوان مدل‌های پایه انتخاب شوند، در پایان مرحله ۷، الگوریتم طراحی شده ۵۰۰ مدل محاسبه شده‌است. خروجی مدل نهایی از رأی‌گیری نتایج تمام مدل‌ها حاصل خواهد شد.

معماری این مدل ترکیبی در شکل ۶ خلاصه شده‌است:



شکل ۶: معماری مدل ترکیبی ۲

به منظور انتخاب مجموعه داده از هر کلاس صفر و یک، برای ساخت هر مدل ۱۰۰۰۰ نمونه انتخاب شده‌است و بنابراین مجموعه داده برای

ساخت هر مدل دارای ۲۰۰۰۰ نمونه می‌باشد. برای هر مجموعه داده پنج مدل یادگیری اجرامی‌شود. مدل‌های یادگیری انتخاب شده شامل درخت تصمیم، رگرسیون لجستیک، نایو بیز گوسین، بردارهای پشتیبان، ک- نزدیکترین همسایگی می‌باشند. بنابراین پس از پایان برنامه ۵۰۰ مدل خواهیم داشت. در ادامه نتایج اجرای این برنامه را شرح خواهیم داد.

جدول ۲: تعداد نمونه‌ها به تفکیک برچسب‌های صفر و یک

مجموعه داده	تعداد نمونه‌ها	تعداد نمونه‌های هر کلاس
کل داده‌ها	۴۳۲۹۳۰	۰: ۴۲۸۹۵۹۲ ۱: ۳۹۶۳۸
داده‌های آموزش	۳۴۶۳۳۸۴	۰: ۳۴۳۱۸۳۱ ۱: ۳۱۵۵۳
داده‌های آزمون	۸۶۵۸۴۶	۰: ۸۵۷۷۶۱ ۱: ۸۰۸۵

ماتریس طبقه‌بندی مدل حاصل از رأی‌گیری این پنج مدل با نسبت رأی گیری ۵۶٪ به صورت زیر است. (به این معنی که اگر بیش از ۵۶٪ از مدل‌ها برچسب یک را پیش‌بینی کرده‌اند، پیش‌بینی مدل نهایی برچسب ۱ است):

	precision	recall	f1-score	support
0	0.990954	0.747986	0.852496	857761
1	0.010202	0.275572	0.019675	8085
accuracy			0.743574	865846
macro avg	0.500578	0.511779	0.436085	865846
weighted avg	0.981796	0.743574	0.844719	865846

که نتایج فوق به این معنی است که ۲۷٪ از کل خسارت‌های جانی به درستی دارای خسارت جانی پیش‌بینی شده‌اند و ۷۴٪ از بدون خسارت‌ها نیز به درستی بدون خسارت پیش‌بینی شده‌اند. با پیاده سازی این روش برای صدور بیمه‌نامه‌ها با کاهش ۲۵/۲ درصدی پرتفو، خسارت‌ها ۲۷/۵٪ کاهش خواهند یافت. این نتایج در مقایسه با مدل‌های ساده غیر ترکیبی که نمونه‌گیری مجدد Smote انجام شده بود عملکرد بهتری است. حتی از ۱۰۰ مدل ترکیبی SVC نیز عملکرد بهتر است.

با کاهش درصد رأی‌گیری به ۵۵٪، نتایج به صورت زیر است. اختلاف رأی‌گیری ۵۵٪ و ۵۶٪، به اندازه اثرگذاری ۵ مدل از ۵۰۰ مدل است. در این مدل، ۷۱ درصد از خسارت‌ها به درستی شناسایی شده‌اند و پرتفو ۶۹٫۴ کوچک شده‌است. در شرایطی که استراتژی شرکت، کوچک شدن پرتفوی بیمه‌نامه‌های رشته بیمه شخص ثالث و کاهش ضریب خسارت باشد به کارگیری چنین مدلی مناسب است.

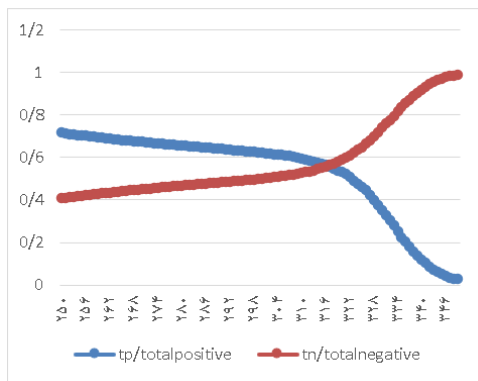
	precision	recall	f1-score	support
0	0.991120	0.916813	0.952519	857872
1	0.012823	0.116253	0.023098	7974
accuracy			0.909440	865846
macro avg	0.501971	0.516533	0.487809	865846
weighted avg	0.982110	0.909440	0.943960	865846

	precision	recall	f1-score	support
0	0.303663	0.991191	0.464898	262785
1	0.713667	0.009568	0.018883	603061
accuracy			0.307491	865846
macro avg	0.508665	0.500379	0.241890	865846
weighted avg	0.589231	0.307491	0.154249	865846

جدول ۳: تعداد مدل در رای گیری و نسبت پیشبینی ها

ت	تپاه	ری	ت
۱۱/۶	۸/۳	۳۴۰	۱
۲۲/۳	۱۶/۷	۳۳۵	۲
۲۵/۱	۲۶/۲	۳۳۰	۳
۷۱	۵۹	۲۵۰	۴

نکته قابل توجه آن است که با افزایش تعداد نمونه‌های هر مدل از ۲۰۰۰۰ نمونه به ۵۵۰۰۰ نمونه عملکرد مدل نهایی بهتر شده است. به منظور شناخت مدل به نمودار زیر توجه کنید. با افزایش تعداد مدل‌ها در رای گیری برای تشخیص خسارت‌ها، تعداد صفرهایی که به درستی شناخته می‌شوند، افزایش و تعداد خسارت‌های تشخیص داده شده کاهش می‌یابد.



شکل ۷: نسبت صفرهای درست و یک‌های درست به تفکیک تعداد مدل‌های رای گیری

شکل فوق بررسی عملکرد مدل نهایی با پارامتر تعداد رای گیری، نمودار آبی نسبت تعداد پیشبینی‌های یک درست به تعداد اعضای کلاس یک‌ها و نمودار قرمز نسبت تعداد پیشبینی‌های صفر درست به تعداد اعضای کلاس صفر را نشان می‌دهد.

در شکل، مشاهده می‌شود که به منظور انتخاب تعداد مدل‌های مورد نیاز رای گیری، لازم است استراتژی شرکت لحاظ شود. اگر استراتژی شرکت کاهش خسارت پرداختی است می‌تواند از تعداد مدل‌های کمتری برای رای گیری استفاده کند و در شرایطی که بزرگ بودن پورتنو برای شرکت در اولویت است از تعداد مدل‌های کمتری برای رای گیری استفاده می‌شود. به عنوان مثال؛ انتخاب تعداد ۳۳۵ مدل برای رای گیری، ۲۲ درصد از بیمه‌نامه‌های منجر به خسارت را به درستی پیش‌بینی می‌کند و

به عنوان مثال شرکتی که در رشته بیمه شخص ثالث عملکرد خوبی ندارد، نهاد ناظر می‌تواند آن شرکت را صرفاً برای صدور بیمه‌نامه‌های بدون خسارت پیش‌بینی شده توسط مدل فوق مجاز بداند. بنابراین به نظر می‌رسد متوازن‌سازی داده‌ها با روش تصادفی ساده به شکل الگوریتم فوق کارایی بهتری دارد.

۱،۲،۴. تغییر تعداد نمونه‌های هر برچسب در مدل ترکیبی

به منظور اصلاح عملکرد مدل فوق در هر بار نمونه‌گیری از کلاس یک، ۲۰۰۰۰ نمونه و از کلاس صفر، ۳۵۰۰۰ نمونه انتخاب می‌کنیم. اجرای این برنامه ۳۵۰۱۳۲،۱۰۹۳۷۵ ثانیه معادل ۹۷،۲۵ ساعت برابر با حدود ۴ شبانه روز به طول انجامیده است. نتایج مدل بر داده‌های آزمون، به صورت زیر است: به ازای رای گیری ۵۰ درصدی (رای گیری با سهم ۲۵۰ مدل):

	precision	recall	f1-score	support
0	0.990863	0.990835	0.990849	857872
1	0.017004	0.017055	0.017030	7974
accuracy			0.981867	865846
macro avg	0.503934	0.503945	0.503940	865846
weighted avg	0.981894	0.981867	0.981881	865846

رای گیری ۳۳۰ مدل:

	precision	recall	f1-score	support
0	0.991888	0.737605	0.846053	857872
1	0.012282	0.351016	0.023733	7974
accuracy			0.734045	865846
macro avg	0.502085	0.544311	0.434893	865846
weighted avg	0.982866	0.734045	0.838480	865846

رای گیری ۳۳۵ مدل:

	precision	recall	f1-score	support
0	0.991403	0.832916	0.905275	857872
1	0.012252	0.222975	0.023228	7974
accuracy			0.827298	865846
macro avg	0.501828	0.527945	0.464252	865846
weighted avg	0.982386	0.827298	0.897152	865846

رای گیری ۳۴۰ مدل:

این موضوع با هزینه شناسایی نادرست ۱۶٪ از بیمه‌نامه‌های بدون خسارت به عنوان بیمه‌نامه‌های دارای خسارت همراه است. واضح است که هزینه ۲۲٪ از خسارت‌ها بیشتر از مجموع حق بیمه ۱۶٪ از بیمه‌نامه‌ها است.

بررسی نتایج مدل ترکیبی با کاهش مدل‌ها به کمتر از ۵ مدل پایه اولیه، نتایج بهتری ارائه کرد و عملکرد مدل نهایی کاهش یافت.

۲.۲.۴. جایگزینی روش رأی‌گیری مدل ترکیبی ۲ با استفاده از مدل‌های یادگیری ماشین

با به کارگیری مدل درخت تصمیم بر خروجی ۵۰۰ مدل، عملکرد مدل حاصل به صورت زیر است:

	precision	recall	f1-score	support
0	0.991246	0.999995	0.995602	428949
1	0.989362	0.046804	0.089380	3974
accuracy			0.991246	432923
macro avg	0.990304	0.523400	0.542491	432923
weighted avg	0.991229	0.991246	0.987283	432923

این مدل بر داده‌های تست ساخته شده است و به نظر می‌رسد استفاده از درخت تصمیم بر خروجی نهایی مدل‌ها می‌تواند عملکرد مدل را بهبود بخشد. ۴/۶٪ از کل خسارت‌ها و ۹۹٪ از صفرها به درستی، شناسایی شده‌اند. ۹۸٪ از یک‌های شناسایی شده صحیح‌اند. استفاده از این مدل در انتخاب ریسک، با کاهش کمتر از ۱٪ از پرتفو، ۴/۶٪ از خسارت‌ها شناسایی می‌شود که منطقی است و برای شرکت‌ها قابل پیاده‌سازی می‌باشد.

روش‌های SVC، جنگل تصادفی و کا- نزدیکترین همسایگی در این رویکرد عملکرد ضعیفتری داشتند.

۵. نتیجه‌گیری

مدل‌های ترکیبی عملکرد بهتری نسبت به الگوریتم‌های جنگل تصادفی، درخت تصمیم، ماشین بردار پشتیبان و رگرسیون لجستیک به همراه روش‌های مختلف متوازن‌سازی در این پژوهش داشته‌اند.

در این مقاله، دو مدل ترکیبی برای تعیین وقوع یا عدم وقوع خسارت در رشته بیمه شخص ثالث ارائه شده است. در هر دو مدل ترکیبی از پنج روش طبقه‌بندی شامل نزدیکترین همسایگی، نایو بیس، درخت تصمیم، لجستیک رگرسیون و نزدیکترین همسایگی به عنوان مدل پایه استفاده شده است. در هر دو مدل ترکیبی پیشنهادی، ۱۰۰ نمونه از داده‌های انتخاب شده است و مدل‌های پایه با استفاده از این صد نمونه در مرحله آموزش، آموزش داده شده‌اند. بنابراین هر مدل ترکیبی دارای ۵۰۰ مدل طبقه‌بندی پایه است. در مدل ترکیبی اول انتخاب داده‌های نمونه برای ساخت هر مدل پایه به صورت تصادفی انتخاب شده‌اند و نمونه‌های انتخاب شده با استفاده از روش smote متوازن شده‌اند. در مدل ترکیبی دوم نمونه داده‌ها، صد مرتبه به صورت متوازن از هر برچسب به صورت جداگانه انتخاب ترکیب شده‌اند و مدل‌های پایه بر هر یک از این صد نمونه

داده آموزش دیده‌اند. مقایسه نتایج دو مدل ترکیبی مذکور نشان می‌دهد مدل دوم عملکرد بهتری دارد. به عبارتی انتخاب داده‌ها به صورت متوازن از هر برچسب مدل‌های بهتری برای پیش‌بینی می‌سازد.

همچنین روش‌های متعددی برای تجمیع مدل‌های پایه وجود دارد. روش رأی‌گیری ساده در این مقاله مورد استفاده قرار گرفته است و نتایج نشان می‌دهد تعداد مدل‌های پیش‌بینی‌کننده برچسب یک که بر چسب اقلیت است می‌تواند به عنوان ابر پارامتر مدل ترکیبی لحاظ شود. این ابر پارامتر بر اساس استراتژی شرکت در حفظ حجم پرتفوی شرکت و یا کاهش ضریب خسارت قابل تنظیم است.

لازم به ذکر است ویژگی‌های ثبت شده در پایگاه داده‌های یک شرکت بیمه می‌تواند افزایش یابد از جمله سوابق تخلفات رانندگی و ویژگی‌هایی مانند شغل راننده می‌تواند در بهبود نتایج مدل‌ها تاثیر چشم‌گیری داشته باشد.

۶. مطالعات آتی

به منظور دستیابی به مدلی قویتر برای فائق آمدن بر مسئله طبقه‌بندی داده‌های نامتوازن لازم است در ساخت مدل ترکیبی سایر روش‌های نمونه‌گیری مجدد را نیز به کار گرفت و یا از ترکیب آن‌ها استفاده کرد. همچنین نحوه تجمیع مدل‌های پایه نیز موضوعی است که می‌تواند به تفصیل مورد بحث و بررسی قرار گیرد و از ابتکارات و تکنیک‌های بیشتری برای تجمیع مدل‌ها استفاده شود. از جمله استفاده از الگوریتم مورچگان در تجمیع درختان می‌تواند مورد استفاده قرار گیرد [9].

سپاسگزاری

در پایان از مدیران واحد طرح و توسعه شرکت بیمه البرز که ما را در انجام مطالعه حاضر حمایت کردند، به ویژه جناب آقای عابدینی، سرکار خانم کمالخانی و سرکار خانم ابوطالبی صمیمانه سپاسگزاریم.

Reference

- [1] K. P. Murphy., Probabilistic Machine Learning: An Introduction, MIT Press, 2022.
- [2] A. Fernández, . S. García and M. Galar, R, Learning from Imbalanced Data Sets, Springer, 2018.
- [3] S. Ardabili, A. Mosavi and . A. R. Varkonyi-Koczy, "Advances in Machine Learning Modeling Reviewing Hybrid and Ensemble Methods," *Preprints*, 2019.
- [4] G. . G. Sundarkumar and V. Ravi, "A novel hybrid under sampling method for mining unbalanced datasets in banking and insurance," *Engineering Applications of Artificial Intelligence*, vol. 37, p. 368–377, 2015.
- [5] S. I. V. Shamitha, S. K. Shamitha and V. Ilango, "A hybrid technique for health insurance fraud detection on highly imbalanced dataset," *International Journal of Innovative technology and exploring engineering (IJITEE)*, vol. 8, no. 11, pp. 2278–3075, 2019.
- [6] S. Kotekani and I. Velchamy, "An Effective Data Sampling Procedure for Imbalanced Data learning on

health insurance fraud detection, CIT," *Journal of Computing and Information Technology*, vol. 28, no. 4, p. 269–285, (2020)..

- [7] J. Brownlee, *Data Preparation for machine learning*, Jason Brownle, 2020.
- [8] A. Géron, *Hands-on Machine learning with scikit-learn, keras, tensorflow*, Beijing, Boston, Farnham, Sebastopol, Tokyo: O'Reilly Media, Inc, 2019.
- [9] J. Kozak, *Decision Tree and ensemble learning based on ant colony algorithm*, Katowice, Poland: Springer, 2019.