




Paper Type: Original Article



Proposing a New Class of Statistical Reliability-Based Methods for Modeling and Decision-Making

Sepideh Etemadi¹, Mehdi Khashei^{1,*} 

¹ Department of Industrial and Systems Engineering, Isfahan University of Technology, Isfahan, Iran; s.etemadi@in.iut.ac.ir; khashei@cc.iut.ac.ir.

Citation:



Etemadi, S., & Khashei, M. (2023). Proposing a new class of statistical reliability-based methods for modeling and decision-making. *Journal of decisions and operations research*, 8(1), 269-282.

Received: 23/10/2021

Reviewed: 21/11/2021

Revised: 12/12/2021

Accepted: 08/01/2022

Abstract

Purpose: The purpose of this paper is to present a new methodology for statistical modeling, which, unlike all commonly developed models and algorithms, maximizes the reliability of the results instead of the resulting accuracy. Accordingly, a new class of statistical modeling approaches has been developed by replacing conventional processes with the proposed process.

Methodology: The multiple linear regression method has been selected to implement the proposed methodology in this paper. To comprehensively evaluate the performance of the proposed regression model, 10 standard datasets from the literature on statistical modeling have been considered.

Findings: Overall, the results show that in 65% of the studied data sets, the proposed model can generalize more than the usual multiple linear regression. The proposed regression model, on average, has been able to improve the accuracy of the modeling by 5.571% and 6.466% in mean absolute error and mean square error, respectively, compared to its classic version. These results clearly show the significant effect of reliability of the results on the degree of generalizability, which is basically not considered in the usual statistical modeling processes.

Originality/Value: Statistical modeling is one of the most important tools for simulating real-world systems and data sets that are often used to make decisions in a wide range of applications. Several different approaches have been developed in the literature with different features to cover real-world issues with the desired accuracy. However, such methods follow a similar concept and idea in the modeling process. The performance basis in all conventional statistical modeling approaches is based on the assumption that maximum accuracy in experimental and inaccessible data will be obtained from models with minimization of error in training data. Although this is a logical and standard procedure in traditional statistical modeling spaces, it is not the unique way to achieve maximum generalizability. In other words, the generalizability of the model simultaneously depends on the model's accuracy and the level of results' reliability. In this paper, a new methodology for statistical modeling is presented, which, unlike all commonly developed models and algorithms, maximizes the reliability of the results instead of the resulting accuracy.

Keywords: Decision-making processes, Statistical modeling, Accuracy and reliability, Generalizability of results, Prediction, Multiple linear regression.

Corresponding Author: khashei@cc.iut.ac.ir

 <http://dorl.net/dor/20.1001.1.25385097.1402.8.5.15.5>


Licensee. **Journal of Decisions and Operations Research**. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0>).



نوع مقاله: پژوهشی



ارایه کلاس جدیدی از روش‌های آماری مبتنی بر قابلیت اعتماد به منظور مدل‌سازی و تصمیم‌گیری

سپیده اعتمادی^۱، مهدی خاشعی^{۱*}

گروه مهندسی صنایع و سیستم‌ها، دانشگاه صنعتی اصفهان، اصفهان، ایران.

چکیده

هدف: در این مقاله، یک متدولوژی جدید برای مدل‌سازی‌های آماری ارایه گردیده است که برخلاف تمامی مدل‌ها و الگوریتم‌های توسعه یافته معمول، قابلیت اعتماد به نتایج را به جای دقت حاصله به حداکثر می‌رساند. بر این اساس در این مقاله، یک دسته جدید از رویکردهای مدل‌سازی آماری با جایگزینی فرآیندهای معمول با فرآیندهای پیشنهادی پیشنهاد گردیده است.

روش‌شناسی پژوهش: در این مقاله، روش رگرسیون خطی چندگانه به منظور پیاده‌سازی متدولوژی پیشنهادی انتخاب شده است. برای ارزیابی جامع عملکرد مدل رگرسیون پیشنهادی، ۱۰ مجموعه داده معیار از ادبیات موضوع مدل‌سازی‌های آماری در نظر گرفته شده است. **یافته‌ها:** به طور کلی، نتایج حاصله نشان می‌دهد که در ۶۵٪ از مجموعه داده‌های بررسی شده، مدل پیشنهادی توانایی تعمیم بیش‌تری نسبت به رگرسیون خطی چندگانه معمول ایجاد نموده است. مدل رگرسیون پیشنهادی، به طور میانگین توانسته است دقت مدل‌سازی‌ها را به ترتیب به میزان ۵/۵۷۱٪ و ۶/۴۶۶٪ در میانگین قدر مطلق خطا و میانگین مربعات خطا نسبت به نسخه کلاسیک خود بهبود بخشد. این نتایج، به وضوح اثر قابل توجه اعتماد به نتایج را بر میزان قابلیت تعمیم نشان می‌دهد که اساساً در فرآیندهای مدل‌سازی آماری معمول لحاظ نمی‌گردد.

اصالت/ارزش افزوده علمی: مدل‌سازی‌های آماری یکی از مهم‌ترین ابزارهای موجود به منظور شبیه‌سازی سیستم‌های تحت مطالعه و مجموعه داده‌های دنیای واقعی می‌باشد که اغلب به منظور تصمیم‌گیری در طیف وسیعی از علوم مورد استفاده قرار می‌گیرد. چندین رویکرد متفاوت در ادبیات موضوع با ویژگی‌های متفاوت برای پوشش مسایل دنیا واقعی با دقت مطلوب توسعه یافته‌اند. با این حال، این‌گونه از روش‌ها از یک مفهوم ایده مشابه در فرآیند مدل‌سازی پیروی می‌کنند. اساس عملکردی در تمامی رویکردهای مدل‌سازی آماری معمول، بر پایه این فرض استوار بوده که حداکثر دقت در داده‌های آزمایش و غیرقابل دسترس از مدل‌هایی با حداقل‌سازی میزان خطا در داده‌های آموزش به دست خواهند آمد. اگرچه، این رویه منطقی و استاندارد در فضاهای مدل‌سازی آماری معمول می‌باشد، اما تنها شیوه منحصر به فرد برای دستیابی به حداکثر قابلیت تعمیم محسوب نمی‌گردد. به عبارت دیگر، قابلیت تعمیم مدل به طور هم‌زمان وابسته به دقت مدل و هم‌چنین سطح قابلیت اعتماد به نتایج حاصله می‌باشد. در این مقاله، یک متدولوژی جدید برای مدل‌سازی‌های آماری ارایه گردیده است که برخلاف تمامی مدل‌ها و الگوریتم‌های توسعه یافته معمول، قابلیت اعتماد به نتایج را به جای دقت حاصله به حداکثر می‌رساند.

کلیدواژه‌ها: فرآیندهای تصمیم‌گیری، مدل‌سازی آماری، دقت و قابلیت اعتماد، قابلیت تعمیم نتایج، پیش‌بینی، رگرسیون خطی چندگانه.

۱- مقدمه

مدل‌سازی یکی از مهم‌ترین ابزارها برای شبیه‌سازی سیستم‌های تحت مطالعه و مجموعه داده‌های موجود در دنیای واقعی است که اغلب در طیف وسیعی از کاربردها با موفقیت استفاده می‌شود. دلیل اصلی محبوبیت مدل‌سازی، درک، تفسیر و تجزیه و تحلیل روابط بین مولفه

* نویسنده مسئول





ها به منظور پیش‌بینی رفتار آینده سیستم مورد مطالعه می‌باشد. به‌ویژه برای مدیران و تصمیم‌گیرندگان بسیار مهم است که تصمیمات مدیریتی، عملیاتی و مالی شایسته‌تری اتخاذ نمایند. در حقیقت، هدف اصلی مدل‌سازی، ساخت متامدلی از سیستم، تحلیل رفتار سیستم، پیش‌بینی آینده و در نهایت تصمیم‌گیری براساس تحلیل‌های انجام شده است [1]. فرآیندهای مدل‌سازی به‌طور کلی برای اهداف مختلف با ناظر و بدون ناظر هم‌چون پیش‌بینی، طبقه‌بندی و خوشه‌بندی توسعه می‌یابند. چندین نوع خطی، غیرخطی، علی و سببی، سری‌های زمانی و مدل‌های آماری و هوشمند در ادبیات موضوع ارایه شده است. این‌گونه از مدل‌ها معمولاً شامل رگرسیون خطی چندگانه^۱، خودرگرسیون میانگین متحرک انباشته^۲، شبکه‌های عصبی مصنوعی^۳، ماشین بردار پشتیبان^۴ در گروه روش‌های با ناظر و k -میانگین^۵ و نقشه‌های خودسازمان‌ده^۶ و تجزیه و تحلیل مولفه‌های اصلی^۷ در گروه رویکردهای بدون ناظر هستند [2].

رگرسیون خطی چندگانه یکی از محبوب‌ترین و پرکاربردترین رویکردهای پیش‌بینی آماری است که کاربردهای بسیاری در حوزه‌های مختلف علمی دارد [3]. رگرسیون خطی چندگانه برای بررسی رابطه میان مجموعه‌ای از متغیرها به‌طور گسترده مورد بهره‌گیری قرار می‌گیرد و تعامل و کنترل بیش‌تر با کاربر بر تجزیه و تحلیل‌های پیش‌بینی کننده را فراهم می‌آورد که از مزایای رقابتی مدل‌های رگرسیون در تقابل با یادگیری ماشین قلمداد می‌شود [4]. این روش بیان می‌کند که چگونه تغییرات در متغیر پاسخ را می‌توان با چندین متغیر توضیح‌دهنده، توصیف نمود. به‌طور کلی، دو دلیل اصلی به‌کارگیری مدل‌های رگرسیون، توضیح و پیش‌بینی است. بهره‌گیری از یک مدل رگرسیون برای اهداف توضیحی بدین‌صورت بیان می‌شود که ضرایب رگرسیون با کنترل سایر عوامل موجود در مدل رگرسیون، برآوردی از تاثیر متغیر توضیحی بر متغیر پاسخ را ارایه می‌دهد؛ بر این اساس رابطه صحیح بین متغیر پاسخ و متغیرهای توضیحی تعیین می‌شود. دومین دلیل اصلی استفاده از مدل رگرسیون به‌منظور پیش‌بینی است. این به معنای پیش‌بینی پاسخ‌های مشاهده نشده یا جدید یا پیش‌بینی مقادیر پاسخ آینده براساس مقادیر فعلی متغیرهای توضیحی است [5]. از رگرسیون خطی چندگانه برای پیش‌بینی در حوزه‌های مختلف هم‌چون پزشکی، مهندسی، انرژی، مالی، مدیریت، محیط‌زیست و... در ادبیات موضوع استفاده می‌شود.

رات و همکاران [6] از روش‌های رگرسیون خطی چندگانه برای پیش‌بینی روند روز بعد در موارد فعال بیماری کرونا در آدیشا و هند استفاده کردند. این مدل‌ها در تشخیص کوید ۱۹ دقت چشمگیری به‌دست آوردند. تانگ و همکاران [7] مدل رگرسیون خطی چندگانه را با استفاده از پارامترهای مشخصه موج نبض شریان شعاعی برای ارزیابی پیری عروق ایجاد کردند. هوانگ و همکاران [8] یک مدل رگرسیون خطی چندگانه مبتنی بر k -میانگین برای پیش‌بینی تعداد بستری هفتگی بیماران انسداد مزمن ریوی در اثر آلاینده‌های اصلی هوا ارایه داده‌اند. این مدل پیش‌بینی با برقراری ارتباط میان بیماری مزمن ریوی و آلاینده‌های هوا، به شناسایی زودهنگام، مداخلات فردی برای کند کردن پیشرفت بیماری و کاهش هزینه‌های پزشکی کمک می‌نماید. برای ارزیابی کارایی مدل از میانگین قدر مطلق درصد خطا^۸ استفاده شده است. سیولا و دی آمیکو [9] روش رگرسیون خطی چندگانه را برای تعیین نیاز به انرژی گرمایش یا خنک‌سازی یک ساختمان عمومی در هر شرایط آب‌وهوایی ایجاد کردند. نتایج، استفاده از این مدل را به‌عنوان یک روش جایگزین توجیه می‌کند. در نتیجه موجب تسریع و کمک به برخی از مراحل ارزیابی در برنامه‌ریزی انرژی شده و هم‌چنین معیارهای معتبری که می‌تواند در استانداردها و قوانین مربوط به عملکرد انرژی ساختمان ایجاد شود، ارایه می‌نماید. پارک و همکاران [10] عملکرد گرمایش ساعتی سیستم پمپ گرمایی منبع زمینی در مقیاس بزرگ را با دقت رضایت‌بخش توسط مدل‌های رگرسیون خطی چندگانه و شبکه عصبی مصنوعی پیش‌بینی نمودند. این مدل‌های پیش‌بینی را می‌توان به‌عنوان مبنایی برای اندازه‌گیری و تایید اقدامات احتمالی صرفه‌جویی انرژی در آینده و نظارت بر عملکرد سیستم استفاده کرد. این مطالعه نشان‌دهنده توانایی رگرسیون خطی چندگانه در تحلیل کمی عوامل تاثیرگذار بر عملکرد سیستم پمپ حرارتی منبع زمینی است.

چرچی و هوردوغان [11] مدل‌های رگرسیون خطی چندگانه و شبکه عصبی مصنوعی را برای تخمین دمای لامپ خشک‌کن و مقادیر رطوبت مطلق هوای فرآیند خارج‌شده از خروجی فرآیند چرخ خشک‌کن طراحی کردند که به‌طور گسترده‌ای در تهیه مطبوع، ذخیره‌سازی مواد غذایی و کاربردهای خشک‌کردن به‌ویژه در مناطق گرم و مرطوب استفاده می‌شود. نتایج این مدل‌ها مورد تحلیل و مقایسه قرار گرفت. از معیارهای ضریب تعیین (R^2)، میانگین قدر مطلق خطا^۹ و جذر میانگین مربعات خطا^{۱۰} برای ارزیابی نتایج به‌دست‌آمده از مدل‌های

¹ Multiple Linear Regression (MLR)

² Auto Regressive Integrated Moving Average (ARIMA)

³ Artificial Neural Networks (ANN)

⁴ Support Vector Machine (SVM)

⁵ K-means

⁶ Self-Organizing Maps (SOM)

⁷ Principal Component Analysis (PCA)

⁸ Mean Absolute Percentage Error (MAPE)

⁹ Mean Absolute Error (MAE)

¹⁰ Root Mean Squared Error (RMSE)



مختلف استفاده شده است. خمت و ریچمن [12] با استفاده از مدل رگرسیون خطی چندگانه مقدار نشت هوا در خانه‌ها را براساس متغیرهایی از جمله هندسه ساختمان، مصالح ساختمان، سن ساختمان و آب‌وهوای محلی پیش‌بینی کردند. شاین و همکاران [13] مدل رگرسیون خطی چندگانه را برای پیش‌بینی مصرف ماهیانه برق و آب در مزارع ایرلند براساس تولید شیر، تعداد موجودی، تجهیزات زیرساختی، رویه‌های مدیریتی و شرایط محیطی به‌کار گرفته‌اند. سیاست‌گذاران می‌توانند از مدل‌های توسعه‌یافته رگرسیون خطی چندگانه برای محاسبه تاثیر تولید لبنیات ایرلندی بر منابع طبیعی یا ابزارهای پشتیبانی تصمیم‌گیری برای محاسبه اثرات احتمالی شیوه‌های بالقوه کاهش مصرف در مزرعه استفاده کنند. تریگو-گونزالس و همکاران [14] مدل رگرسیون خطی چندگانه را برای تخمین ساعتی تولید برق فتوولتاییک از نظر بازده انرژی و براساس فناوری‌های مختلف ارائه دادند. مدل پیشنهادی با مدل تخمین فتوولتاییک دیگری که در ادبیات موضوع ارائه شده مورد مقایسه قرار گرفت و مدل پیشنهادی نتایج بهتری را از نظر جذر میانگین مربع خطا نشان داده است. سیاوش و همکاران [15] منحنی توان توربین و سرعت روتور را برای توربین بادی کوچک مجهز به طیف وسیعی از زاویه باز کانال در هر سرعت باد با استفاده از مدل‌های رگرسیون خطی چندگانه و شبکه عصبی مصنوعی پیش‌بینی کردند. چهار مدل رگرسیون خطی چندگانه در اشکال مختلف و یک شبکه عصبی پرسپترون چندلایه برای تخمین قدرت و سرعت زاویه‌ای روتور یک توربین بادی ارائه شده است. دقت مدل‌های پیش‌بینی رگرسیون خطی چندگانه و شبکه عصبی از نظر معیارهای آماری جذر میانگین مربعات خطا و ضریب تعیین ارائه گردیده است.

خو و همکاران [16] مدل رگرسیون خطی چندگانه را برای پیش‌بینی مناطق بالقوه مناسب کشت برای گیاهان خاص پیشنهاد کردند. نتایج حاکی از آن است که این مدل پنج عامل اصلی فعال زیستی در چین را به‌طور موثر پیش‌بینی می‌کند. آبروگی و همکاران [17] رگرسیون خطی چندگانه و شبکه عصبی را برای پیش‌بینی عملکرد محصول سیب‌زمینی آلی با استفاده از سیستم‌های خاکورزی و خواص خاک ارزیابی نمودند. نتایج نشان داد که مدل رگرسیون خطی چندگانه عملکرد محصول را با دقت بیش‌تری نسبت به مدل شبکه عصبی تخمین می‌زند. لی و همکاران [18] از مدل رگرسیون چندگانه برای تخمین توزیع فاصله‌ای رطوبت خاک در کره جنوبی استفاده کردند. ضرایب مدل با توجه به بارش فصلی پنج روز قبل برآورد گردید. زی و همکاران [19] مدل‌های رگرسیون خطی چندگانه و رگرسیون جنگل تصادفی^۱ را برای برآورد فعالیت‌های آمیلاز و اوره آز خاک در یک زمین احیاشده ساحلی اجرا نمودند. پهلوان-راد و همکاران [20] عملکرد مدل‌های رگرسیون خطی چندگانه و جنگل تصادفی را برای پیش‌بینی میزان نفوذ خاک در یک دشت خشک در شرق ایران مقایسه کردند. معیارهای ارزیابی جذر میانگین مربعات خطا و میانگین قدرمطلق خطا بین مدل‌ها مشابه بود. پالمرو و همکاران [21] مدل‌های رگرسیون خطی چندگانه را ایجاد کردند که می‌تواند مقدار کل سمیت و غلظت دیوکسین را در انتشارات جوی به‌درستی تخمین بزند. استویچف و همکاران [22] از مدل رگرسیون خطی چندگانه ابتکاری برای ارزیابی آلودگی فلزات در رسوبات سطح یک تالاب ساحلی استفاده نمودند. یوچی و همکاران [23] از رگرسیون خطی چندگانه و جنگل تصادفی برای مدل‌سازی آلودگی هوای داخلی با ۸۷ متغیر بالقوه پیش‌بینی کننده از داده‌های نظارت بر فضای باز، پرسشنامه‌ها، ارزیابی‌های خانه و مجموعه داده‌های جغرافیایی استفاده کردند. تانگ و همکاران [24] الگوریتم‌های رگرسیون خطی چندگانه و ماشین بردار پشتیبان را برای پیش‌بینی میزان تجزیه بیولوژیکی به‌عنوان یک فرآیند موثر برای حذف مواد شیمیایی آلی از آب، خاک و محیط‌های رسوبی ایجاد نمودند. مدل‌های پیشنهادی می‌توانند به‌عنوان ابزارهای موثر برای پیش‌بینی تجزیه بیولوژیکی مواد شیمیایی آلی و برای ارزیابی ماندگاری زیست‌محیطی مواد شیمیایی آلی استفاده شوند. حسین زاده و همکاران [25] کارایی مدل‌های شبکه عصبی و رگرسیون خطی چندگانه را در پیش‌بینی بازیافت مواد مغذی از مواد زاید جامد، تحت تیمارهای مختلف ورمی کمپوست ارزیابی نمودند. مدل‌های توسعه‌یافته شبکه عصبی و رگرسیون خطی چندگانه از نظر معیارهای آماری از جمله ضریب تعیین، ضریب تعدیل شده، جذر میانگین مربعات خطا و قدرمطلق میانگین انحرافات مقایسه شدند.

مدل‌های رگرسیون خطی چندگانه هم‌چون سایر روش‌های آماری در ادبیات موضوع، تفکر یکسانی در مورد متدولوژی مدل‌سازی دارند. منطقی ایجاد چنین مدل‌هایی به حداکثر رساندن دقت عملکردی داده‌های آموزش برای دستیابی به حداکثر دقت در داده‌های آزمون یا توانایی تعمیم مدل است. بر این اساس، توانایی تعمیم در این نوع مدل‌ها صرفاً مربوط به دقت عملکرد می‌باشد. اگرچه دقت یکی از مهم‌ترین عوامل تاثیرگذار بر توانایی تعمیم مدل است و متدولوژی مورد استفاده برای به‌دست آوردن مناسب‌ترین مدل رگرسیون خطی چندگانه کاملاً منطقی است. با این حال اولاً عامل منحصر به فرد توضیح‌دهنده چگونگی تغییر توانایی تعمیم مدل نیست و ثانیاً تنها راه ممکن به منظور دستیابی به مدلی با حداکثر سطح تعمیم به شمار نمی‌آید. براساس ادبیات موضوع، انتظار می‌رود که یک عامل دیگر

¹ Random Forest (RF)



تاثیرگذار بر توانایی تعمیم مدل، میزان اعتماد به دقت عملکردی یا به عبارت دیگر، تغییر در دقت عملکردی در مقابل شرایط مختلف است که در تفکر متعارف مدل‌سازی رگرسیون خطی چندگانه در نظر گرفته نشده است. از این رو ضعف رویکرد مدل‌سازی معمول نادیده گرفتن سطح اعتماد به دقت به عنوان یکی دیگر از عوامل موثر بر توانایی تعمیم مدل‌ها است.

بر این اساس، در این مقاله یک مفهوم مدل‌سازی جدید مربوط به تاثیر قابلیت اعتماد بر قابلیت تعمیم مدل برای توسعه مدل‌های رگرسیون خطی چندگانه مبتنی بر قابلیت اعتماد ارائه شده است. هدف اصلی ایده پیشنهادی ارائه کلاس جدیدی از مدل‌های رگرسیون خطی چندگانه مبتنی بر قابلیت اعتماد است. مفهوم اساسی روش پیشنهادی بر حداکثرسازی قابلیت اعتماد مدل به جای دقت برای افزایش قابلیت تعمیم استوار است. به این ترتیب، مدل‌هایی با تغییر دقت کم‌تر در مواجهه با شرایط مختلف داده‌ای، عملکرد قابل اعتمادتری خواهند داشت و بنابراین، ممکن است دقت عملکرد بالاتری در داده‌های آزمون داشته باشند یا طبق اصطلاحات پیش‌بینی، ممکن است قابلیت تعمیم بالاتری داشته باشند. از سوی دیگر، نقصان اساسی در مدل‌های رگرسیون خطی چندگانه مبتنی بر دقت که به طور معمول توسعه یافته‌اند و شکاف اصلی تحقیقاتی در این زمینه که نوآوری اصلی این مقاله نیز است، در نظر گرفتن قابلیت اعتماد مدل‌های رگرسیون خطی چندگانه به منظور ارائه مدل‌های با قابلیت تعمیم بالاتر است. در این مقاله رگرسیون خطی چندگانه برای اجرای ایده پیشنهادی انتخاب شده است. دلایل اصلی این انتخاب را می‌توان به شرح زیر خلاصه نمود:

۱. الگوریتم یادگیری بهینه مستقیم ریاضی: از آن جاکه یکی از عوامل موثر بر قابلیت تعمیم، الگوریتم یادگیری مورد استفاده در مدل است؛ بنابراین، برای از بین بردن تاثیر الگوریتم یادگیری استفاده شده بر قابلیت تعمیم، در این مقاله، مدل رگرسیون خطی چندگانه که دارای الگوریتم یادگیری بهینه مستقیم ریاضی است، انتخاب شده است. به این ترتیب، تفاوت بین نسخه‌های قابلیت اعتماد محور و دقت محور رگرسیون خطی چندگانه تحت تاثیر الگوریتم یادگیری قرار نمی‌گیرند. از سوی دیگر، به این ترتیب می‌توان اطمینان حاصل کرد که تغییر در قابلیت تعمیم فقط به دلیل تغییر در قابلیت اعتماد/دقت مدل‌ها است.
۲. سهولت به کارگیری و طراحی: از آن جاکه یکی دیگر از عوامل موثر بر قابلیت تعمیم، ساختار طراحی شده مدل مورد استفاده است؛ بنابراین، برای از بین بردن تاثیر روش طراحی بر قابلیت تعمیم، در این مقاله از مدل رگرسیون خطی چندگانه که از ساده‌ترین فرآیند طراحی بهره می‌برد، استفاده شده است. به این ترتیب، تفاوت بین نسخه‌های دقت محور و قابلیت اعتماد محور رگرسیون خطی چندگانه تحت تاثیر ساختار طراحی شده، قرار نمی‌گیرد. به عبارت دیگر، به این ترتیب می‌توان اطمینان حاصل نمود که تغییر در قابلیت تعمیم صرفاً به دلیل تغییر در قابلیت اعتماد/دقت مدل‌ها است.
۳. پیچیدگی: یکی دیگر از عوامل تاثیرگذار بر قابلیت تعمیم، پیچیدگی مدل استفاده شده است؛ بنابراین، به منظور حذف تاثیر پیچیدگی بر تعمیم، در این مقاله مدل رگرسیون خطی چندگانه که دارای کم‌ترین پیچیدگی (خطی و مستقیم) در فرآیند مدل‌سازی است، انتخاب شده است.

اگرچه در این مقاله متدولوژی پیشنهادی تنها بر روی مدل رگرسیون خطی چندگانه پیاده‌سازی شده است، با این وجود می‌توان رویکرد پیشنهادی را با برخی تغییرات در تمام روش‌های مدل‌سازی آماری و حتی غیر آماری با ناظر و بدون ناظر، هم‌چون مدل‌های لاجیت، خودرگرسیون، شبکه‌های عصبی مصنوعی، ماشین بردار پشتیبانی، مدل‌های فازی و ... اجرا نمود. با این حال، باید توجه داشت که در هر دسته فوق‌الذکر، ممکن است تفاوت‌هایی به وجود آید. به عنوان مثال، در مدل‌های غیر خطی، ممکن است یک راه‌حل بهینه ریاضی مستقیم غیر قابل دسترس باشد. به عبارت دیگر، در موقعیت‌های غیر خطی، راه‌حل را صرفاً می‌توان از یک فرآیند غیر مستقیم و تکراری به دست آورد و/یا ممکن است یک شکل راه‌حل ریاضی نداشته باشد و/یا شاید تنها یک راه‌حل غیر بهینه حاصل شود. نهایتاً، به منظور ارزیابی جامع مدل رگرسیون خطی چندگانه مبتنی بر قابلیت اعتماد، ۱۰ مجموعه داده معیار با کاربردهای مختلف در حوزه‌های متفاوت از پایگاه داده UCI انتخاب شده است. این مجموعه داده‌ها با مفاهیم مدل‌سازی جدید مبتنی بر قابلیت اعتماد و مبتنی بر دقت کلاسیک مدل‌سازی شده و نتایج مورد مقایسه و تحلیل قرار می‌گیرند. لازم به ذکر است که در این مقاله، مدل رگرسیون خطی چندگانه مبتنی بر قابلیت اعتماد تنها با مدل رگرسیون خطی چندگانه مبتنی بر دقت کلاسیک مقایسه شده است. دلیل آن این است که رگرسیون خطی چندگانه یک روش مدل‌سازی بهینه در طبقه مدل‌های آماری سببی خطی است. بر این اساس، عملکرد سایر مدل‌های استاندارد در طبقه مدل‌های سببی خطی نمی‌تواند بهتر از رگرسیون خطی چندگانه باشد.

در مدل‌های رگرسیون خطی چندگانه معمول، منطق مدل‌سازی همیشه مبتنی بر به حداقل رساندن خطای مدل‌سازی در داده‌های آموزش بوده است. بر همین پایه، اساس به دست آوردن حداکثر دقت در داده‌های آزمون یا قدرت تعمیم مدل منحصر از طریق به حداقل رساندن دقت در داده‌های آموزش فرض می‌شود. از آن‌جا که کیفیت تصمیمات با توانایی تعمیم مدل در مسایل دنیای واقعی ارتباط نزدیکی دارد، در نظر گرفتن دقت عملکردی به عنوان تنها عامل موثر بر توانایی تعمیم مدل و نادیده گرفتن پتانسیل افزایش تعمیم با بهبود سایر عوامل تاثیرگذار، صحیح به نظر نمی‌رسد. بر این اساس، بهبود توانایی تعمیم مدل‌ها در ادبیات پیش‌بینی و مدل‌سازی به یکی از چالش‌برانگیزترین حوزه‌های تحقیقاتی تبدیل شده است. در این راستا به نظر می‌رسد اعتماد به دقت عملکردی یکی از عوامل تاثیرگذار بر قابلیت تعمیم مدل باشد که در روند مدل‌سازی رگرسیون خطی چندگانه معمول مورد توجه قرار نگرفته است. پایداری در مدل‌سازی یک سیستم با داده‌های مختلف به قابل اعتماد بودن نتایج اشاره دارد. هم‌چنین توانایی تعمیم مدل به معنای دستیابی به پیش‌بینی‌های دقیق و قابل اعتماد است؛ بنابراین، سطح اعتماد دقت و دقت از عوامل موثر بر توانایی تعمیم مدل و مدل‌سازی عدم قطعیت به حساب می‌آیند. در این مقاله، ایده جدید مدل‌سازی رگرسیون خطی چندگانه ارائه شده تا به جای دقت، قابلیت اعتماد مدل را به حداقل برساند. منطق عملکردی پیشنهادی از محاسبه تغییرات در دقت عملکردی مدل‌های اجرا شده بر روی داده‌های اعتبارسنجی الهام می‌گیرد؛ بنابراین، اساس منطق رویکرد پیشنهادی بر این اصل استوار است که انتظار می‌رود مدل‌های رگرسیونی با تغییرات دقت کم‌تر در برابر داده‌های اعتبارسنجی، نتایج قابل اعتمادتری در داده‌های آزمون داشته باشند. به عبارت دیگر، هرچه نوسان عملکردها در داده‌های اعتبارسنجی کم‌تر باشد، ثبات عملکردها در داده‌های غیرقابل دسترس یا آزمون بیش‌تر است. به این ترتیب، منطق ایجاد مدل رگرسیون خطی چندگانه مبتنی بر قابلیت اعتماد توصیف می‌شود.

رگرسیون خطی چندگانه یکی از فرآیندهای متداول با ناظر برای یادگیری روابط خطی بین متغیرهای وابسته (پیش‌بینی شده) و مستقل (پیش‌بینی کننده) بوده به شیوه‌ای که پیش‌بینی متغیر وابسته را بر اساس مقدار جدید متغیرهای مستقل فراهم می‌نماید. از این رو یک مدل رگرسیون خطی k متغیره شامل متغیر وابسته Y و $K-1$ متغیر مستقل یا توضیحی X_1, X_2, \dots, X_k را می‌توان به صورت زیر نوشت:

$$Y_t = \beta_1 + \beta_2 X_{2t} + \beta_3 X_{3t} + \dots + \beta_k X_{kt} + u_t, \quad t = 1, 2, 3, \dots, N. \quad (1)$$

که β_1 عرض از مبدا β_2 تا β_k ضرایب شیب هستند، u_t عبارت خطای تصادفی است و N اندازه نمونه است. مدل حداقل مربعات معمولی یکی از محبوب‌ترین تکنیک‌های مبتنی بر دقت است که اغلب برای تخمین پارامترهای نامعلوم استفاده می‌شود. در تکنیک حداقل مربعات معمولی، پارامترهای نامعلوم به گونه‌ای برآورد می‌شوند تا مجموع مربعات خطاها، یعنی مربع تفاوت بین مقادیر واقعی و برازش شده به حداقل برسد. با این حال، تکنیک حداقل مربعات معمولی هم‌چون سایر برآوردگرهای مبتنی بر دقت بر دقت به عنوان عامل اصلی تاثیرگذار بر توانایی تعمیم تمرکز می‌کند. در حالی که در فرآیند پیشنهادی حداقل مربعات مبتنی بر قابلیت اعتماد، تغییر در مربعات خطاها به حداقل می‌رسد. لازم به ذکر است که متدولوژی پیشنهادی را می‌توان در کلیه الگوریتم‌های یادگیری با ناظر و بدون ناظر اعمال نمود [26].

بر اساس روش پیشنهادی، ابتدا بخشی از داده‌های آموزش به عنوان مجموعه داده‌های اعتبارسنجی در نظر گرفته می‌شود. سپس دقت در این مقاله مجموع مربعات خطا، برای داده‌های آموزش و هم‌چنین داده‌های آموزش به علاوه هر داده از مجموعه اعتبارسنجی به شرح زیر محاسبه می‌گردد:

$$\sum_{t=1}^N \hat{u}_{0t}^2 = \sum_{t=1}^N \left(Y_t - \hat{\beta}_{01} X_{1t} - \hat{\beta}_{02} X_{2t} - \dots - \hat{\beta}_{0k} X_{kt} \right)^2 = \sum_{t=1}^N \left(Y_t - \sum_{j=1}^k \hat{\beta}_{0j} X_{jt} \right)^2. \quad (2)$$





$$\sum_{t=1}^{N+1} \hat{u}_{1t}^2 = \sum_{t=1}^{N+1} \left(Y_t - \hat{\beta}_{11} X_{1t} - \hat{\beta}_{12} X_{2t} - \dots - \hat{\beta}_{1k} X_{kt} \right)^2 = \sum_{t=1}^{N+1} \left(Y_t - \sum_{j=1}^k \hat{\beta}_{1j} X_{jt} \right)^2,$$

$$\sum_{t=1}^{N+2} \hat{u}_{2t}^2 = \sum_{t=1}^{N+2} \left(Y_t - \hat{\beta}_{21} X_{1t} - \hat{\beta}_{22} X_{2t} - \dots - \hat{\beta}_{2k} X_{kt} \right)^2 = \sum_{t=1}^{N+2} \left(Y_t - \sum_{j=1}^k \hat{\beta}_{2j} X_{jt} \right)^2, \quad (۳)$$

$$\dots\dots\dots$$

$$\sum_{t=1}^{N+n} \hat{u}_{nt}^2 = \sum_{t=1}^{N+n} \left(Y_t - \hat{\beta}_{n1} X_{1t} - \hat{\beta}_{n2} X_{2t} - \dots - \hat{\beta}_{nk} X_{kt} \right)^2 = \sum_{t=1}^{N+n} \left(Y_t - \sum_{j=1}^k \hat{\beta}_{nj} X_{jt} \right)^2.$$

که $\sum_{t=1}^{N+i} \hat{u}_{it}^2$, $i = 0, 1, 2, \dots, n, t = 1, 2, \dots, N + i$ مجموع مربعات باقیمانده‌ها و n اندازه مجموعه اعتبارسنجی است. حال، برای یافتن مقدار بهینه پارامترهای نامعلوم در هر وضعیت داده β_{ij} , $i = 0, 1, 2, \dots, n, j = 1, 2, \dots, k$ پارامترها به گونه‌ای تخمین زده می‌شوند که $\sum \hat{u}_{it}^2$ کمینه شود. این کار با مشتق گرفتن هر معادله نسبت به پارامترها در هر موقعیت داده و برابر قرار دادن نتیجه آن با صفر حاصل می‌شود. این روش k معادلات هم‌زمان با k پارامتر نامعلوم، برای هر وضعیت داده، به شرح زیر ارائه می‌دهد. برای داده‌های آموزش:

$$\hat{\beta}_{01} \sum_{t=1}^N X_{1t} + \hat{\beta}_{02} \sum_{t=1}^N X_{2t} + \hat{\beta}_{03} \sum_{t=1}^N X_{3t} + \hat{\beta}_{04} \sum_{t=1}^N X_{4t} + \dots + \hat{\beta}_{0k} \sum_{t=1}^N X_{kt} = \sum_{t=1}^N Y_t,$$

$$\hat{\beta}_{01} \sum_{t=1}^N X_{2t} + \hat{\beta}_{02} \sum_{t=1}^N X_{2t}^2 + \hat{\beta}_{03} \sum_{t=1}^N X_{2t} X_{3t} + \dots + \hat{\beta}_{0k} \sum_{t=1}^N X_{2t} X_{kt} = \sum_{t=1}^N X_{2t} Y_t,$$

$$\hat{\beta}_{01} \sum_{t=1}^N X_{3t} + \hat{\beta}_{02} \sum_{t=1}^N X_{3t} X_{2t} + \hat{\beta}_{03} \sum_{t=1}^N X_{3t}^2 + \dots + \hat{\beta}_{0k} \sum_{t=1}^N X_{3t} X_{kt} = \sum_{t=1}^N X_{3t} Y_t, \quad (۴)$$

$$\dots\dots\dots$$

$$\hat{\beta}_{01} \sum_{t=1}^N X_{kt} + \hat{\beta}_{02} \sum_{t=1}^N X_{kt} X_{2t} + \hat{\beta}_{03} \sum_{t=1}^N X_{kt} X_{3t} + \dots + \hat{\beta}_{0k} \sum_{t=1}^N X_{3t} X_{ki}^2 = \sum_{t=1}^N X_{kt} Y_t.$$

و به طور مشابه برای اولین داده از مجموعه داده اعتبارسنجی:

$$\hat{\beta}_{11} \sum_{t=1}^{N+1} X_{1t} + \hat{\beta}_{12} \sum_{t=1}^{N+1} X_{2t} + \hat{\beta}_{13} \sum_{t=1}^{N+1} X_{3t} + \hat{\beta}_{14} \sum_{t=1}^{N+1} X_{4t} + \dots + \hat{\beta}_{1k} \sum_{t=1}^{N+1} X_{kt} = \sum_{t=1}^{N+1} Y_t,$$

$$\hat{\beta}_{11} \sum_{t=1}^{N+1} X_{2t} + \hat{\beta}_{12} \sum_{t=1}^{N+1} X_{2t}^2 + \hat{\beta}_{13} \sum_{t=1}^{N+1} X_{2t} X_{3t} + \dots + \hat{\beta}_{1k} \sum_{t=1}^{N+1} X_{2t} X_{kt} = \sum_{t=1}^{N+1} X_{2t} Y_t,$$

$$\hat{\beta}_{11} \sum_{t=1}^{N+1} X_{3t} + \hat{\beta}_{12} \sum_{t=1}^{N+1} X_{3t} X_{2t} + \hat{\beta}_{13} \sum_{t=1}^{N+1} X_{3t}^2 + \dots + \hat{\beta}_{1k} \sum_{t=1}^{N+1} X_{3t} X_{kt} = \sum_{t=1}^{N+1} X_{3t} Y_t, \quad (۵)$$

$$\dots\dots\dots$$

$$\hat{\beta}_{11} \sum_{t=1}^{N+1} X_{kt} + \hat{\beta}_{12} \sum_{t=1}^{N+1} X_{kt} X_{2t} + \hat{\beta}_{13} \sum_{t=1}^{N+1} X_{kt} X_{3t} + \dots + \hat{\beta}_{1k} \sum_{t=1}^{N+1} X_{3t} X_{ki}^2 = \sum_{t=1}^{N+1} X_{kt} Y_t.$$

و نهایتاً برای آخرین داده اعتبارسنجی:

$$\hat{\beta}_{n1} \sum_{t=1}^{N+n} X_{1t} + \hat{\beta}_{n2} \sum_{t=1}^{N+n} X_{2t} + \hat{\beta}_{n3} \sum_{t=1}^{N+n} X_{3t} + \hat{\beta}_{n4} \sum_{t=1}^{N+n} X_{4t} + \dots + \hat{\beta}_{nk} \sum_{t=1}^{N+n} X_{kt} = \sum_{t=1}^{N+n} Y_t,$$

$$\hat{\beta}_{n1} \sum_{t=1}^{N+n} X_{2t} + \hat{\beta}_{n2} \sum_{t=1}^{N+n} X_{2t}^2 + \hat{\beta}_{n3} \sum_{t=1}^{N+n} X_{2t} X_{3t} + \dots + \hat{\beta}_{nk} \sum_{t=1}^{N+n} X_{2t} X_{kt} = \sum_{t=1}^{N+n} X_{2t} Y_t,$$

$$\hat{\beta}_{n1} \sum_{t=1}^{N+n} X_{3t} + \hat{\beta}_{n2} \sum_{t=1}^{N+n} X_{3t} X_{2t} + \hat{\beta}_{n3} \sum_{t=1}^{N+n} X_{3t}^2 + \dots + \hat{\beta}_{nk} \sum_{t=1}^{N+n} X_{3t} X_{kt} = \sum_{t=1}^{N+n} X_{3t} Y_t, \quad (۶)$$

$$\dots\dots\dots$$

$$\hat{\beta}_{n1} \sum_{t=1}^{N+n} X_{kt} + \hat{\beta}_{n2} \sum_{t=1}^{N+n} X_{kt} X_{2t} + \hat{\beta}_{n3} \sum_{t=1}^{N+n} X_{kt} X_{3t} + \dots + \hat{\beta}_{nk} \sum_{t=1}^{N+n} X_{3t} X_{ki}^2 = \sum_{t=1}^{N+n} X_{kt} Y_t.$$



با به حداقل رساندن اختلاف بین هر دو مربعات خطاهای ایجاد شده با افزودن داده‌های اعتبارسنجی به داده‌های آموزش، مدل رگرسیون خطی چندگانه مبتنی بر قابلیت اعتماد ایجاد می‌شود. هدف به‌دست آوردن یک تابع رگرسیون خطی چندگانه است که کم‌ترین عدم اعتماد یا عدم قطعیت و یا به‌عبارت‌دیگر حداکثر قابلیت اعتماد را داشته باشد. بر این اساس برای رسیدن به یک خط رگرسیون قابلیت اعتماد محور با حداقل تغییر مربعات خطاها در تمام نقاط داده اعتبارسنجی، پارامترهای نامعلوم تمام خطوط رگرسیون مبتنی بر دقت باید با یکدیگر برابر باشند؛ بنابراین داریم:

$$\hat{\beta}_{ij} = \hat{\beta}_{i'j} \quad \text{for all } i, i' = 0, 1, 2, \dots, n; \quad j = 1, 2, \dots, k, \tag{۷}$$

$$\hat{\beta}e_j = \hat{\beta}_{ij} \quad \text{for all } i = 0, 1, 2, \dots, n.$$

که $\hat{\beta}e_k$ پارامتر خط رگرسیونی پیشنهادی است. به این ترتیب، معادله (۴) تا معادله (۶) را می‌توان به شرح زیر نشان داد:

$$\begin{aligned} \hat{\beta}e_1 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{1t} + \hat{\beta}e_2 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t} + \hat{\beta}e_3 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t} + \dots + \hat{\beta}e_k \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt} &= \sum_{i=0}^n \sum_{t=1}^{N+i} X_{1t} Y_t, \\ \hat{\beta}e_1 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t} + \hat{\beta}e_2 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t}^2 + \hat{\beta}e_3 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t} X_{3t} + \dots + \hat{\beta}e_k \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t} X_{kt} &= \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t} Y_t, \\ \hat{\beta}e_1 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t} + \hat{\beta}e_2 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t} X_{2t} + \hat{\beta}e_3 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t}^2 + \dots + \hat{\beta}e_k \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t} X_{kt} &= \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t} Y_t, \end{aligned} \tag{۸}$$

$$\hat{\beta}e_1 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt} + \hat{\beta}e_2 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt} X_{2t} + \hat{\beta}e_3 \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt} X_{3t} + \dots + \hat{\beta}e_k \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt}^2 = \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt} Y_t.$$

در فرم ماتریس، می‌توان آن را به‌صورت زیر توصیف نمود:

$$\begin{bmatrix} \sum_{i=0}^n \sum_{t=1}^{N+i} X_{1t} & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t} & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t} & \dots & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt} \\ \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t} & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t}^2 & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t} X_{3t} & \dots & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t} X_{kt} \\ \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t} & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t} X_{2t} & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t}^2 & \dots & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t} X_{kt} \\ \dots & \dots & \dots & \dots & \dots \\ \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt} & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt} X_{2t} & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt} X_{3t} & \dots & \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt}^2 \end{bmatrix} \begin{bmatrix} \hat{\beta}e_1 \\ \hat{\beta}e_2 \\ \hat{\beta}e_3 \\ \dots \\ \hat{\beta}e_k \end{bmatrix} = \begin{bmatrix} \sum_{i=0}^n \sum_{t=1}^{N+i} X_{1t} Y_t \\ \sum_{i=0}^n \sum_{t=1}^{N+i} X_{2t} Y_t \\ \sum_{i=0}^n \sum_{t=1}^{N+i} X_{3t} Y_t \\ \dots \\ \sum_{i=0}^n \sum_{t=1}^{N+i} X_{kt} Y_t \end{bmatrix}. \tag{۹}$$

نهایتاً، پارامترهای نامعلوم خط رگرسیونی پیشنهادی را می‌توان با حل معادله (۹) به‌دست آورد. به‌عنوان مثال، برای یک مدل رگرسیون خطی ۳ متغیره، پارامترهای مدل پیشنهادی را می‌توان به‌صورت زیر نشان داد:

$$\begin{aligned} \hat{\beta}e_1 &= \frac{(A_{22}A_{33} - A_{23}^2)B_1 - (A_2A_{33} - A_3A_{23})B_2 + (A_2A_{23} - A_3A_{22})B_3}{A_1A_{22}A_{33} - A_1A_{23}^2 - A_2^2A_{33} + 2A_2A_3A_{23} - A_{22}A_3^2}, \\ \hat{\beta}e_2 &= \frac{(A_3A_{23} - A_2A_{33})B_1 + (A_1A_{33} - A_3^2)B_2 + (A_2A_3 - A_1A_{23})B_3}{A_1A_{22}A_{33} - A_1A_{23}^2 - A_2^2A_{33} + 2A_2A_3A_{23} - A_{22}A_3^2}, \\ \hat{\beta}e_3 &= \frac{(A_2A_{23} - A_3A_{22})B_1 + (A_2A_3 - A_1A_{23})B_2 + (A_1A_{22} - A_2^2)B_3}{A_1A_{22}A_{33} - A_1A_{23}^2 - A_2^2A_{33} + 2A_2A_3A_{23} - A_{22}A_3^2}, \end{aligned} \tag{۱۰}$$

$$\sum_{i=0}^n \sum_{t=1}^{N+i} X_{jt} X_{j't} = A_{j,j'}, \quad j, j' = 1, 2, \dots, k, \quad \sum_{i=0}^n \sum_{t=1}^{N+i} X_{jt} Y_t = B_j, \quad j = 1, 2, \dots, k.$$

به‌طورکلی مزیت اصلی مدل رگرسیون خطی چندگانه قابلیت اعتماد محور نسبت به سایر مدل‌های خطی در نظر گرفتن تاثیر قابلیت اعتماد بر قابلیت تعمیم مدل است. هم‌چنین با به حداقل رساندن قابلیت اعتماد مدل پیشنهادی، عدم قطعیت مدل پیشنهادی در مقایسه با سایر مدل‌های خطی به حداقل می‌رسد. مدل پیشنهادی برای انواع مختلف مسایل تصمیم‌گیری در مقایسه با سایر مدل‌های خطی قابل اجرا است و در مسایل با عدم قطعیت بیش‌تر نتایج کارآمدتر و دقیق‌تری را در داده‌های دیده نشده به‌دست می‌آورد. بر این اساس، کیفیت تصمیمات به‌دست‌آمده از مدل رگرسیون خطی چندگانه اعتماد محور بیش‌تر از سایر مدل‌های خطی است. هم‌چنین، ضرایب



مدل پیشنهادی که براساس فرآیند برآورد پارامتر مبتنی بر قابلیت اعتماد به‌دست آمده‌اند، بهینه هستند. از سوی دیگر، اگرچه پیچیدگی رویکرد پیشنهادی مبتنی بر قابلیت اعتماد نسبت به رویکرد مبتنی بر دقت بیشتر است، با این حال معادله مستقیم و خطی جهت تخمین پارامترهای مدل مبتنی بر قابلیت اعتماد مطابق رابطه (۹) به‌دست آمده است.

۳- مجموعه داده‌ها

در این مقاله، ۱۰ مجموعه داده معیار که از بخش رگرسیون پایگاه داده UCI انتخاب شده‌اند، به‌منظور ارزیابی جامع عملکرد رگرسیون خطی چندگانه پیشنهادی در مقایسه با نسخه رگرسیون معمول در نظر گرفته شده‌اند. این مجموعه داده‌ها شامل مثال‌های واقعی و یا شبیه‌سازی‌شده در زمینه‌های مختلف مانند محیط‌زیست، زیست‌شناسی، انرژی، حمل‌ونقل، کامپیوتر و آموزش می‌باشند که طی سال‌های ۱۹۸۷ تا ۲۰۱۹ میلادی جمع‌آوری شده‌اند. اندازه نمونه این مجموعه داده‌ها از ۲۰۹ تا ۴۵۷۳۰ داده متفاوت است. علاوه‌براین، تعداد متغیرها در مدل‌ها از ۶ تا ۳۱ متغیر توضیحی تغییر می‌نماید. مشخصات ۱۰ مجموعه معیار موردنظر در جدول ۱ خلاصه گردیده‌اند. این مجموعه داده‌ها دارای متغیرهای مختلفی هستند که شامل دو دسته اصلی: ۱- منفرد و ۲- مختلط (حقیقی، صحیح) و دو زیرمجموعه از ویژگی‌های منفرد شامل حقیقی و صحیح می‌باشند.

جدول ۱- اطلاعات عمومی داده‌های معیار سببی در پایگاه داده UCI.

Table 1- General information of causal criterion data in UCI database.

عنوان	سال انتشار	تعداد نمونه / متغیر	نوع متغیر	حوزه کاربرد
1 Student performance	2014	31.395	صحیح	آموزش (ارزیابی)
2 UJIIndoorLoc-Mag	2015	6.540	حقیقی/صحیح	کامپیوتر (جهت‌یابی)
3 Electrical grid stability simulated data	2018	11.10000	حقیقی	انرژی (پیش‌بینی پایداری شبکه هوشمند)
4 Physicochemical properties of protein tertiary structure	2013	9.45730	حقیقی	علوم زیست‌شناسی (بررسی خصوصیات فیزیکوشیمیایی)
5 Beijing multi-site air-quality	2019	15.32907	حقیقی/صحیح	محیط‌زیست (ارزیابی کیفیت هوا)
6 Bike sharing	2013	12.731	حقیقی/صحیح	حمل‌ونقل (پیش‌بینی تعداد دوچرخه‌های اجاره‌ای)
7 Appliances energy prediction	2017	27.19735	حقیقی	انرژی (پیش‌بینی مصرف برق لوازم‌خانگی)
8 EEG-steady-state visual evoked potential signals	2018	13.2944	صحیح	کامپیوتر (سنجش سیگنال‌های مرتبط با خواب)
9 SML2010	2014	20.4137	حقیقی	محیط‌زیست (پیش‌بینی دمای هوا)
10 Computer hardware	1987	7.209	صحیح	کامپیوتر (پیش‌بینی عملکرد پردازنده)

۴- تحلیل نتایج

در این بخش در ابتدا دو مجموعه داده به‌عنوان نمونه از کل داده‌ها انتخاب شده و فرآیند روش پیشنهادی بر روی آن‌ها اجرا گردیده است. سپس نتایج حاصل از اجرای مدل رگرسیون پیشنهادی بر روی ۱۰ مجموعه داده معیار مورد تحلیل و بررسی قرار گرفته است.

در این مقاله، مجموعه داده “Bike Sharing” به‌عنوان نمونه اول انتخاب شده و فرآیند روش پیشنهادی و هم‌چنین تحلیل‌های مربوطه در مورد آن بیان می‌شود. این مجموعه داده مربوط به تعداد دوچرخه‌های کرایه‌ای ساعتی بین سال‌های ۲۰۱۱ تا ۲۰۱۲ میلادی با بهره‌گیری از اطلاعات هواشناسی و فصلی است. متغیرهای ورودی شامل روز، فصل، سال، ماه، تعطیلات، روز هفته، روز کار، وضع هوا، دما، دمای حسی، رطوبت، سرعت باد به ترتیب X_1 تا X_{12} می‌باشند که برای پیش‌بینی تعداد کل دوچرخه‌های کرایه‌ای اعم از ثبت‌نام‌شده و تصادفی است [27]. نمودار تعداد کل دوچرخه‌های کرایه‌ای (Y) در شکل ۱ و نمودار مقادیر متغیرهای توضیحی در مقابل متغیر وابسته در شکل ۲ ترسیم گردیده‌اند. هم‌چنین از آن جایی که متغیرهای ۱ تا ۸ شامل روز، فصل، سال، ماه، تعطیلات، روز هفته، روز کار و وضع هوا هستند؛

بر همین اساس گزارش خصوصیات آماری این دسته از متغیرها بی معنا است؛ لذا صرفاً خصوصیات آماری متغیرهای توضیحی نهم تا دوازدهم و متغیر وابسته در جدول ۲ گزارش شده است.

این مجموعه داده ابتدا به دو زیرمجموعه داده‌های آموزش و آزمون تفکیک می‌شود. سپس مطابق با مدل پیشنهادی، بخشی از مجموعه داده‌های آموزش به‌عنوان داده‌های اعتبارسنجی لحاظ می‌گردند. در این مقاله، ۸۵٪ از داده‌های خام به‌طور تصادفی به‌عنوان مجموعه آموزش و ۱۵٪ باقیمانده به‌عنوان مجموعه آزمون انتخاب شده‌اند. به‌علاوه، ۱۰٪ از داده‌های آموزش به‌طور تصادفی به‌عنوان مجموعه داده‌های اعتبارسنجی انتخاب گردیده‌اند. از آنجایی که نسبت شکستن مجموعه داده‌های آموزش/اعتبارسنجی و آزمون به‌صورت ۷۵٪، ۱۰٪ و ۱۵٪ در ادبیات موضوع رویکردهای مدل‌سازی جز نسبت‌های معمول و پرکاربرد است، لذا در این مقاله نیز این نسبت در نظر گرفته شده است. بر این اساس، تابع رگرسیونی برآوردشده از روش رگرسیون خطی چندگانه کلاسیک که توسط الگوریتم بهینه‌سازی حداقل مربعات معمولی محاسبه گردیده است، مطابق با معادله (۱۱) حاصل می‌شود.

$$\hat{Y} = 139531.251 - 3.403 X_1 + 295.262 X_2 + 3307.985 X_3 + 137.165 X_4 - 424.103 X_5 + 47.201 X_6 + 97.512 X_7 - 542.982 X_8 + 3106.924 X_9 + 2035.499 X_{10} - 943.181 X_{11} - 2370.935 X_{12} \quad (11)$$

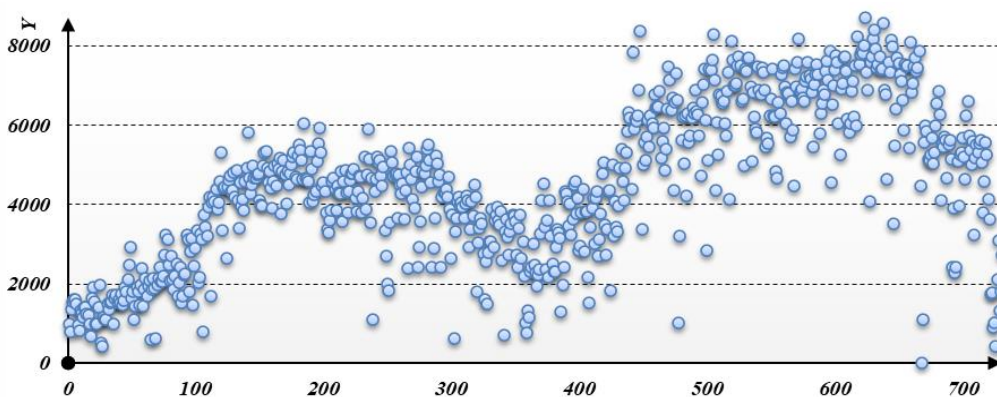
به‌طور مشابه، تابع رگرسیونی برآوردشده از روش رگرسیون خطی چندگانه پیشنهادی که توسط حداکثرسازی قابلیت اعتماد محاسبه گردیده است، مطابق با معادله (۱۲) به‌دست آمده است.

$$\hat{Y} = 176121.747 - 4.310 X_1 + 309.982 X_2 + 3631.063 X_3 + 161.411 X_4 - 399.604 X_5 + 49.398 X_6 + 69.772 X_7 - 524.995 X_8 + 815.052 X_9 + 4710.939 X_{10} - 956.308 X_{11} - 2243.079 X_{12} \quad (12)$$

جدول ۲- خصوصیات آماری مجموعه داده Bike Sharing

Table 2- Statistical characteristics of Bike Sharing dataset.

Y	X ₁₂	X ₁₁	X ₁₀	X ₉	نمونه	حوزه کاربرد
431	0.022	0	0.079	0.059	داده‌های	مینیمم
8362	0.507	0.973	0.841	0.862	داده‌های	ماکسیمم
1162	0.111	0.613	0.655	0.266	آموزش	مد
4338.500	0.184	0.623	0.512	0.528	آزمون	میانه
4312.984	0.192	0.625	0.485	0.509		میانگین
1861.314	0.076	0.146	0.168	0.189		انحراف معیار
22	0.047	0.333	0.220	0.216	داده‌های	مینیمم
8714	0.407	0.925	0.610	0.658	داده‌های	ماکسیمم
-	0.083	0.570	0.608	0.599	آموزش	مد
5611	0.166	0.637	0.391	0.389	آزمون	میانه
5573.234	0.182	0.642	0.412	0.421		میانگین
2012.443	0.086	0.123	0.113	0.124		انحراف معیار



شکل ۱- نمودار تعداد کل دوچرخه‌های کرایه‌ای (Y).

Figure 1- Chart of the total number of rental bicycles (Y).





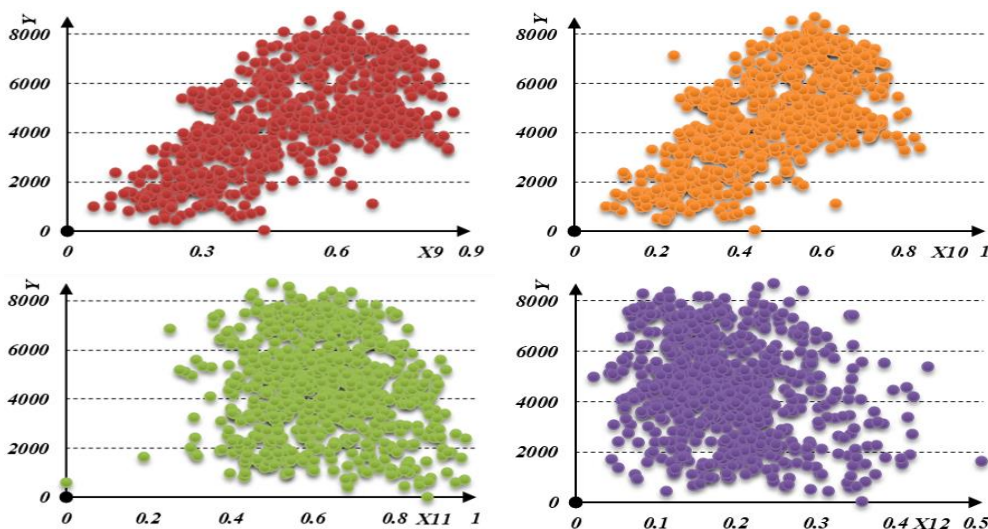
سپس، به منظور حذف اثر فرآیندهای احتمالی بر عملکرد مدل‌ها، هر مدل ۱۰۰ مرتبه تکرار می‌شود. شاخص‌های عملکردی، شامل میانگین قدرمطلق خطا و میانگین مربعات خطای مدل‌ها و میزان بهبود مدل رگرسیون خطی چندگانه پیشنهادی قابلیت اعتماد محور در مقایسه با نسخه کلاسیک دقت محور خود به ترتیب $953/507\%$ ، 156518% ، $967/994\%$ ، 1609922% ، $1/497\%$ و $2/779\%$ به دست می‌آیند. مجموعه داده با عنوان “*UJIIndoorLoc-Mag*” به عنوان نمونه دوم انتخاب شده و تحلیل‌ها در مورد آن بیان می‌شود. این مجموعه داده مربوط به آزمایش سیستم موقعیت‌یابی داخلی بوده که وابسته به تغییر در میدان مغناطیسی زمین می‌باشد. متغیرهای توضیحی شامل مقادیر مغناطیسی در سه محور X ، Y و Z و همچنین مقادیر شتاب‌سنج در این سه محور مختصاتی به ترتیب X_1 تا X_6 می‌باشند که برای پیش‌بینی مقدار سنسور جهت‌یابی (Y) به کار گرفته می‌شوند. نمودار مقادیر سنسور موقعیت‌یابی (Y) در شکل ۳ و نمودار مقادیر متغیرهای توضیحی در مقابل متغیر وابسته در شکل ۴ ترسیم گردیده‌اند. هم‌چنین خصوصیات آماری متغیرهای توضیحی و متغیر وابسته در جدول ۳ گزارش شده است. تابع رگرسیونی برآوردشده از روش رگرسیون خطی چندگانه کلاسیک که توسط الگوریتم بهینه‌سازی حداقل مربعات معمولی محاسبه گردیده است، در ادامه مطابق با معادله (۱۳) حاصل می‌شود.

$$\hat{Y} = 1.547 - 0.035 X_1 - 0.028 X_2 + 0.036 X_3 + 0.283 X_4 + 0.181 X_5 - 0.025 X_6. \quad (13)$$

به‌طور مشابه، تابع رگرسیونی برآوردشده از روش رگرسیون خطی چندگانه پیشنهادی که توسط حداکثرسازی قابلیت اعتماد محاسبه گردیده است، مطابق با معادله (۱۴) به دست آمده است.

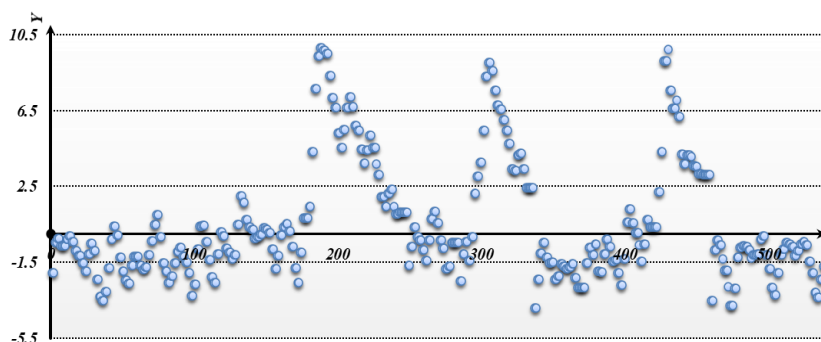
$$\hat{Y} = 1.096 - 0.014 X_1 - 0.068 X_2 + 0.029 X_3 + 0.118 X_4 + 0.204 X_5 - 0.062 X_6. \quad (14)$$

شاخص‌های عملکردی، شامل میانگین قدرمطلق خطا و میانگین مربعات خطای مدل‌ها و میزان بهبود مدل رگرسیون خطی چندگانه پیشنهادی قابلیت اعتماد محور در مقایسه با نسخه کلاسیک دقت محور خود به ترتیب $1/023\%$ ، $2/162\%$ ، $1/965\%$ ، $5/940\%$ و $52/692\%$ و $66/913\%$ به دست می‌آیند.



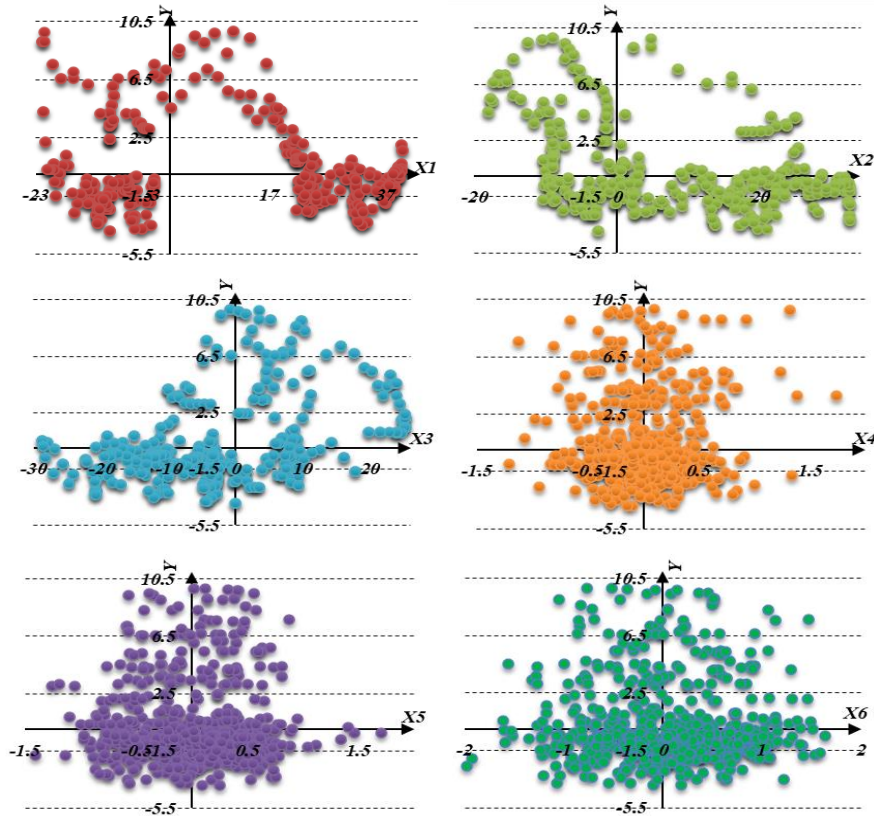
شکل ۲- نمودار متغیر تعداد دوچرخه‌های کرایه‌ای (محور عمودی) در برابر متغیرهای توضیحی X_9 تا X_{12} (محور افقی).

Figure 2- Variable diagram of the number of rental bicycles (vertical axis) against explanatory variables X_9 to X_{12} (horizontal axis).



شکل ۳- نمودار متغیر سنسور جهت‌یابی (Y).

Figure 3- Variable diagram of orientation sensor (Y).



شکل ۴- نمودار متغیر سنسور جهت‌یابی (محور عمودی) در برابر متغیرهای توضیحی X_1 تا X_6 (محور افقی).
Figure 4- Variable diagram of the orientation sensor (vertical axis) against explanatory variables X_1 to X_6 (horizontal axis).

نتایج به‌دست آمده از این مطالعه موردی بیان‌گر این است که مدل رگرسیون خطی چندگانه مبتنی بر قابلیت اعتماد پیشنهادی می‌تواند به‌طور قابل‌توجهی عملکرد مدل رگرسیون مبتنی بر دقت کلاسیک را بهبود بخشد. با این حال، به‌طور کلی در ادبیات موضوع مدل‌سازی نشان‌داده شده که مشخصات داده‌ها، اساساً عملکرد مدل‌ها را تحت تاثیر قرار می‌دهد و نوع و میزان این تاثیر در مدل‌های مختلف متفاوت است؛ بنابراین، به‌منظور حذف تاثیر مشخصات داده‌ها بر عملکرد مدل‌ها، فرآیند مدل‌سازی پیشنهادی و کلاسیک که پیش‌تر بیان گردید برای ۱۰ مجموعه داده معیار دیگر نیز تکرار می‌شود. نتایج حاصله به تفکیک هر مجموعه داده در جدول ۴ خلاصه گردیده‌اند.

در ادامه نتایج مدل‌های رگرسیون پیشنهادی و کلاسیک از نظر کیفیت، یعنی تعداد و میزان برتری و هم‌چنین دقت، به معنی درصد بهبود توانایی تعمیم نتایج، بررسی و باهم مقایسه می‌شوند. نتایج حاصل از پیاده‌سازی مدل‌های پیشنهادی و کلاسیک در ۱۰ مجموعه داده معیار نشان می‌دهد که در ۶ مورد (۶۰٪ از موارد)، از نظر معیارهای ارزیابی میانگین قدرمطلق خطا و میانگین مربعات خطا و در ۷ مورد (۷۰٪ از موارد) از نظر حداقل یکی از معیارهای مذکور، مدل پیشنهادی برتر از مدل رگرسیون کلاسیک بوده است. به‌علاوه، در ۷ مورد (۷۰٪ از موارد) از نظر میانگین قدرمطلق خطا و در ۶ مورد (۶۰٪ از موارد) در معیار میانگین مربعات خطا، مدل پیشنهادی نسبت به مدل رگرسیون کلاسیک برتر می‌باشد. نهایتاً، با نرخ ۶۵٪، مدل رگرسیون پیشنهادی به‌طور متوسط از مدل کلاسیک بهتر عمل نموده است.

از نظر کمی، این نتایج به‌وضوح تاثیر قابلیت اعتماد را به‌عنوان یک عامل تاثیرگذار، علاوه بر دقت، بر توانایی تعمیم مدل مبتنی بر رگرسیون نشان می‌دهد. هم‌چنین، نتایج نشان می‌دهد که وقتی قابلیت تعمیم مدنظر باشد، اهمیت نسبی قابلیت اعتماد به‌طور متوسط بیش‌تر از دقت است. به‌طور مشابه از نظر بهبود دقت نتایج، مدل پیشنهادی، مدل رگرسیون کلاسیک را در معیار میانگین قدرمطلق خطا به میزان ۵/۵۷۱٪، در میانگین مربعات خطا به میزان ۶/۴۶۶٪ و به‌طور متوسط از نظر هر دو معیار مذکور به میزان ۶/۰۱۸٪ ارتقا داده



است. این نتایج حاکی از برتری مدل پیشنهادی نسبت به مدل رگرسیون کلاسیک است. براساس این پیامدها، می‌توان استنباط نمود که قابلیت اعتماد نسبت به دقت بر توانایی تعمیم عامل موثرتری می‌باشد. از این رو، ممکن است معقول‌تر باشد که در شرایط کاملاً ناشناخته یا در انتخاب نا آگاهانه یک مدل مبتنی بر رگرسیون برای اهداف مدل‌سازی، مدل رگرسیون پیشنهادی بر مدل رگرسیون کلاسیک ترجیح داده شود.

جدول ۳- خصوصیات آماری مجموعه داده UJIIndoorLoc-Mag

Table 3- Statistical characteristics of the UJIIndoorLoc-Mag data set.

Y	X ₆	X ₅	X ₄	X ₃	X ₂	X ₁	نمونه	حوزه کاربرد
-3.91	-1.98	-1.41	-1.21	-28.80	-17.94	-21.78	داده‌های آموزشی	مینیمم
9.82	1.63	5.31	1.73	85.80	28.68	95.40	داده‌های آموزشی	ماکسیمم
1.12	0.77	-0.21	-0.14	-2.82	-1.50	-10.32	داده‌های آموزشی	مد
-0.17	0.0004	0.01	-0.01	-0.06	0.48	21.84	داده‌های آموزشی	میانه
0.92	0.003	0.02	0.01	-0.17	3.64	13.56	داده‌های آموزشی	میانگین
3.18	0.68	0.58	0.39	12.79	12.65	20.68	داده‌های آموزشی	انحراف معیار
-3.78	-1.23	-0.87	-0.79	-21.30	15.42	-12.42	داده‌های آزمایش	مینیمم
-0.10	1.42	0.81	0.87	-3.30	32.52	-2.70	داده‌های آزمایش	ماکسیمم
-1.94	-	-	-	-10.08	16.98	-7.26	داده‌های آزمایش	مد
-1.11	-0.05	0.02	-0.002	-10.08	26.04	-7.08	داده‌های آزمایش	میانه
-1.48	0.01	-0.02	-0.01	-10.72	25.20	-7.15	داده‌های آزمایش	میانگین
1.04	0.68	0.49	0.31	6.00	5.96	2.59	داده‌های آزمایش	انحراف معیار

جدول ۴- معیارهای عملکردی در مدل رگرسیون خطی چندگانه پیشنهادی و کلاسیک.

Table 4- Performance criteria in the proposed and classical multiple linear regression model.

میزان بهبود (%)		مدل دقت محور		مدل اعتماد محور		عنوان	
MSE	MAE	MSE	MAE	MSE	MAE		
0.453	4.488	6.068	1.502	6.040	1.435	Student performance	1
66.913	52.692	5.940	2.162	1.965	1.023	UJIIndoorLoc-Mag	2
-0.040	0.078	0.14515	0.3211	0.14521	0.3209	Electrical grid stability simulated data	3
0.033	0.005	5.583	1.61603	5.581	1.61596	Physicochemical properties of protein tertiary structure	4
1.272	0.551	584.838	17.519	577.400	17.422	Beijing multi-site air-quality	5
2.779	1.497	1609922	967.994	1565180	953.507	Bike sharing	6
-0.535	-0.001	7422	50.5270	7462	50.5273	Appliances energy prediction	7
-8.031	-3.243	391.351	16.316	422.782	16.845	EEG-steady-state visual evoked potential signals	8
-7.874	-4.479	2.335	1.197	2.519	1.250	SML2010	9
9.687	-4.118	2929	28.912	2645	27.721	Computer hardware	10
6.466	5.571					میانگین	

۵- نتیجه‌گیری

توانایی تعمیم مدل که ارتباط مستقیمی با کیفیت تصمیمات دنیای واقعی دارد، یکی از چالش‌برانگیزترین موضوعات در حیطه مدل‌سازی است. بر این اساس، کشف عوامل موثر و ارزیابی تاثیر هر عامل بر میزان بهبود تعمیم، برای پژوهشگران بسیار حایز اهمیت است. در این راستا، در این مقاله، تاثیر قابلیت اعتماد بر توانایی تعمیم مدل رگرسیون خطی چندگانه در طبقه مدل‌های آماری کلاسیک خطی سببی بررسی شده است. در این مقاله، برخلاف روش‌های معمول مدل‌سازی آماری، هدف به حداکثر رساندن قابلیت اعتماد و ارایه یک مدل رگرسیون خطی چندگانه مبتنی بر قابلیت اعتماد است. درحالی‌که در ادبیات موضوع، دقت مهم‌ترین و منحصر به فردترین عامل تاثیرگذار بر توانایی تعمیم مدل در نظر گرفته شده است، نتایج تجربی به دست آمده از این مقاله نشان‌دهنده برتری توانایی تعمیم مدل مبتنی بر قابلیت اعتماد نسبت به مدل مبتنی بر دقت است.

بر این اساس، انتظار می‌رود که با تمرکز بر رویکرد مبتنی بر قابلیت اعتماد، کیفیت تصمیمات اتخاذ شده افزایش یابد. هم‌چنین رویکرد مدل‌سازی مبتنی بر قابلیت اعتماد با توانایی مدل‌سازی عدم قطعیت، روش مناسب‌تری برای حل مسایل تصمیم‌گیری در دنیای واقعی



محسوب می‌شود. در این مقاله، نتایج پیاده‌سازی مدل رگرسیون پیشنهادی بر روی ۱۰ داده معیار مورد تجزیه و تحلیل قرار گرفته است. از دیدگاه کلی، در ۶۰٪ از موارد مدل رگرسیون پیشنهادی از نظر معیارهای میانگین قدر مطلق خطا و میانگین مربعات خطا توانایی تعمیم بهتری دارد. در ۷۰٪ از موارد مدل رگرسیون پیشنهادی از نظر حداقل یکی از دو معیار ارزیابی مذکور نسبت به مدل‌های مبتنی بر دقت توانایی تعمیم بهتری کسب نموده است. نتایج تجربی حاصله حاکی از این است که قابلیت اعتماد و دقت از عوامل موثر در توانایی تعمیم مدل‌های رگرسیون خطی چندگانه هستند. با این حال، بر اساس نتایج به دست آمده، قابلیت اعتماد نسبت به دقت عامل موثرتری در بهبود توانایی تعمیم مدل به حساب می‌آید. بر همین اساس، رویکرد مدل‌سازی مبتنی بر قابلیت اعتماد با تمرکز بر به حداکثر رساندن اعتماد نتایج به منظور افزایش قدرت تعمیم مدل، دارای کارایی بالایی می‌باشد. به این ترتیب، می‌توان با جایگزینی فرآیندهای متداول با فرآیند پیشنهادی در انواع دسته‌بندی‌های مدل‌سازی و مدل‌ها، یک کلاس جدید از رویکردهای مدل‌سازی ایجاد نمود. به عنوان پیشنهادها پژوهشی آتی می‌توان به اجرای رگرسیون خطی چندگانه پیشنهادی بر مجموعه داده‌ها در حوزه‌های تصمیم‌گیری مختلف اشاره نمود. همچنین می‌توان فرآیند پیشنهادی مبتنی بر قابلیت اعتماد را بر مدل‌های دیگری در دسته مدل‌های طبقه‌بندی خطی، مدل‌های سری زمانی خطی، مدل‌های هوشمند غیر خطی و سایر مدل‌ها پیاده‌سازی نمود. نهایتاً می‌توان به شناسایی عوامل مختلف داده‌ای و مدل بر برتری هر دو رویکرد مدل‌سازی مبتنی بر قابلیت اعتماد و دقت در مسایل تصمیم‌گیری متفاوت اقدام نمود.

منابع

- [1] Khashei, M., & Bijari, M. (2010). An artificial neural network (p, d, q) model for timeseries forecasting. *Expert systems with applications*, 37(1), 479–489.
- [2] Khashei, M., & Bijari, M. (2011). A novel hybridization of artificial neural networks and ARIMA models for time series forecasting. *Applied soft computing*, 11(2), 2664–2675.
- [3] Catalina, T., Iordache, V., & Caracaleanu, B. (2013). Multiple regression model for fast prediction of the heating energy demand. *Energy and buildings*, 57, 302–312.
- [4] Yildiz, B., Bilbao, J. I., & Sproul, A. B. (2017). A review and analysis of regression and machine learning models on commercial building electricity load forecasting. *Renewable and sustainable energy reviews*, 73, 1104–1122.
- [5] Fitzmaurice, G. M. (2016). Regression. *Diagnostic histopathology*, 22(7), 271–278.
- [6] Rath, S., Tripathy, A., & Tripathy, A. R. (2020). Prediction of new active cases of coronavirus disease (COVID-19) pandemic using multiple linear regression model. *Diabetes & metabolic syndrome: clinical research & reviews*, 14(5), 1467–1474.
- [7] Tang, Q., Huang, L., & Pan, Z. (2019). Multiple linear regression model for vascular aging assessment based on radial artery pulse wave. *European journal of integrative medicine*, 28, 92–97. <https://doi.org/10.1016/j.eujim.2019.05.006>
- [8] Huang, Z., Lin, S., Long, L., Cao, J., Luo, F., Qin, W., ... Gregersen, H. (2020). Predicting the morbidity of chronic obstructive pulmonary disease based on multiple locally weighted linear regression model with K-means clustering. *International journal of medical informatics*, 139, 104141. <https://doi.org/10.1016/j.ijmedinf.2020.104141>
- [9] Ciulla, G., & D'Amico, A. (2019). Building energy performance forecasting: a multiple linear regression approach. *Applied energy*, 253, 113500. <https://doi.org/10.1016/j.apenergy.2019.113500>
- [10] Park, S. K., Moon, H. J., Min, K. C., Hwang, C., & Kim, S. (2018). Application of a multiple linear regression and an artificial neural network model for the heating performance analysis and hourly prediction of a large-scale ground source heat pump system. *Energy and buildings*, 165, 206–215. <https://doi.org/10.1016/j.enbuild.2018.01.029>
- [11] Çerçi, K. N., & Hürdoğan, E. (2020). Comparative study of multiple linear regression (MLR) and artificial neural network (ANN) techniques to model a solid desiccant wheel. *International communications in heat and mass transfer*, 116, 104713. <https://doi.org/10.1016/j.icheatmasstransfer.2020.104713>
- [12] Khemet, B., & Richman, R. (2018). A univariate and multiple linear regression analysis on a national fan (de) Pressurization testing database to predict airtightness in houses. *Building and environment*, 146, 88–97. <https://doi.org/10.1016/j.buildenv.2018.09.030>
- [13] Shine, P., Scully, T., Upton, J., & Murphy, M. D. (2018). Multiple linear regression modelling of on-farm direct water and electricity consumption on pasture based dairy farms. *Computers and electronics in agriculture*, 148, 337–346. <https://doi.org/10.1016/j.compag.2018.02.020>
- [14] Trigo-González, M., Batlles, F. J., Alonso-Montesinos, J., Ferrada, P., Del Sagrado, J., Martínez-Durbán, M., Cortés, M., Partillo, C., & Marzo, A. (2019). Hourly PV production estimation by means of an exportable multiple linear regression model. *Renewable energy*, 135, 303–312. <https://doi.org/10.1016/j.bse.2020.104052>
- [15] Siavash, N. K., Ghobadian, B., Najafi, G., Rohani, A., Tavakoli, T., Mahmoodi, E., & Mamat, R. (2021). Prediction of power generation and rotor angular speed of a small wind turbine equipped to a controllable duct using artificial neural network and multiple linear regression. *Environmental research*, 196, 110434. <https://doi.org/10.1016/j.envres.2020.110434>



- [16] Xu, N., Meng, F., Zhou, G., Li, Y., Wang, B., & Lu, H. (2020). Assessing the suitable cultivation areas for *Scutellaria baicalensis* in China using the Maxent model and multiple linear regression. *Biochemical systematics and ecology*, 90, 104052. <https://doi.org/10.1016/j.bse.2020.104052>
- [17] Abrougui, K., Gabsi, K., Mercatoris, B., Khemis, C., Amami, R., & Chehaibi, S. (2019). Prediction of organic potato yield using tillage systems and soil properties by artificial neural network (ANN) and multiple linear regressions (MLR). *Soil and tillage research*, 190, 202–208. <https://doi.org/10.1016/j.still.2019.01.011>
- [18] Lee, Y., Jung, C., & Kim, S. (2019). Spatial distribution of soil moisture estimates using a multiple linear regression model and Korean geostationary satellite (COMS) data. *Agricultural water management*, 213, 580–593. <https://doi.org/10.1016/j.agwat.2018.09.004>
- [19] Xie, X., Wu, T., Zhu, M., Jiang, G., Xu, Y., Wang, X., & Pu, L. (2021). Comparison of random forest and multiple linear regression models for estimation of soil extracellular enzyme activities in agricultural reclaimed coastal saline land. *Ecological indicators*, 120, 106925. <https://doi.org/10.1016/j.ecolind.2020.106925>
- [20] Pahlavan-Rad, M. R., Dahmardeh, K., Hadizadeh, M., Keykha, G., Mohammadnia, N., Gangali, M., Keikha, M., Davatgar, N., & Brungard, C. (2020). Prediction of soil water infiltration using multiple linear regression and random forest in a dry flood plain, eastern Iran. *Catena*, 194, 104715. <https://doi.org/10.1016/j.catena.2020.104715>
- [21] Palmer, D., Pou, J. O., Gonzalez-Sabaté, L., & Diaz-Ferrero, J. (2018). Multiple linear regression based congener profile correlation to estimate the toxicity (TEQ) and dioxin concentration in atmospheric emissions. *Science of the total environment*, 622, 510–516. <https://doi.org/10.1016/j.scitotenv.2017.11.344>
- [22] Stoichev, T., Coelho, J. P., De Diego, A., Valenzuela, M. G. L., Pereira, M. E., de Chanvalon, A. T., & Amouroux, D. (2020). Multiple regression analysis to assess the contamination with metals and metalloids in surface sediments (Aveiro Lagoon, Portugal). *Marine pollution bulletin*, 159, 111470. <https://doi.org/10.1016/j.marpolbul.2020.111470>
- [23] Yuchi, W., Gombojav, E., Boldbaatar, B., Galsuren, J., Enkhmaa, S., Beejin, B., ... & Allen, R. W. (2019). Evaluation of random forest regression and multiple linear regression for predicting indoor fine particulate matter concentrations in a highly polluted city. *Environmental pollution*, 245, 746-753. <https://doi.org/10.1016/j.envpol.2018.11.034>
- [24] Tang, W., Li, Y., Yu, Y., Wang, Z., Xu, T., Chen, J., ... Li, X. (2020). Development of models predicting biodegradation rate rating with multiple linear regression and support vector machine algorithms. *Chemosphere*, 253, 126666. <https://doi.org/10.1016/j.chemosphere.2020.126666>
- [25] Hosseinzadeh, A., Baziar, M., Alidadi, H., Zhou, J. L., Altaee, A., Najafpoor, A. A., & Jafarpour, S. (2020). Application of artificial neural network and multiple linear regression in modeling nutrient recovery in vermicompost under different conditions. *Bioresource technology*, 303, 122926. <https://doi.org/10.1016/j.biortech.2020.122926>
- [26] Etemadi, S., & Khashei, M. (2021). Etemadi multiple linear regression. *Measurement*, 186, 1–19. <https://doi.org/10.1016/j.measurement.2021.110080>
- [27] Fanaee-T, H., & Gama, J. (2014). Event labeling combining ensemble detectors and background knowledge. *Progress in artificial intelligence*, 2, 113–127. <https://doi.org/10.1007/s13748-013-0040-3>