# Prediction of building carbon emission based on grey Markov theory

**Chenhao Zhao**

Tianjin Chengjian University, Tianjin, 300192, China

1852400343@qq.com

**Abstract.** The increase in carbon emissions in the construction industry is one of the factors that lead to environmental problems and hinder the sustainable development of society and the construction industry. Therefore, there is an urgent need to rationally control carbon emissions by formulating scientific and effective energy-saving and emission reduction policies so as to safeguard the development of the industry. However, because of the uncertainty and complexity of $CO_2$ emissions, more reliable prediction and assessment tools are needed to comprehensively analyze and predict $CO_2$ emissions. This study combines the gray system theory and Markov principle to give full play to the advantages of the two methods, using Markov theory to determine the state transition probability, find out the characteristics of carbon emissions in the construction industry, and form a suitable gray Markov prediction model. The feasibility of the prediction model is demonstrated by calculating and examining the prediction model using data related to carbon emissions from the U.S. buildings from 2002 to 2022, which greatly improves the accuracy of $CO_2$ emissions prediction.

**Keywords:** carbon emission, building, construction industry, Grey Markov prediction model.

## 1. Introduction

Many countries are facing challenges with water and energy supply, as well as air pollution, due to factors such as population growth, climate change, droughts, increased water and energy usage, depletion of resources, and inconsistent use of fossil fuels. In the coming years, many countries will face greater challenges in supplying water and energy, while the increasing use of fossil fuels will exacerbate air pollution and climate change.

Climate change currently poses a threat to the entire world, causing financial challenges for countries while seriously affecting the normal lives of people, and human beings, societies and countries will also pay a heavy price for their own deleterious impacts on the environment. Shifting climate patterns have led to an increasing frequency of extreme weather events, while at the same time greenhouse gas emissions have reached their highest levels to date [3]. The year 2020 is an important year for the countries to increase their Nationally Determined Contributions (NDCs), which will essentially consider further measures to diminish energy usage and emissions comprising embodied emissions in the construction sector and buildings [4].

In developed countries, buildings are responsible for 50% of energy consumption and 30% of greenhouse gas emissions [1]. The construction industry, while beneficial for employment and economic growth, imposes considerable environmental costs. It is among the sectors globally

contributing to $CO_2$ emissions [2]. Recent experiments and studies have highlighted a sharp increase in carbon emissions from the construction industry, leading to increasingly severe environmental issues that impede the sustainable development of both society and the construction sector. Consequently, there's a pressing need to ensure the industry's development through the formulation of scientifically sound and effective energy-saving and emission-reduction policies, along with the rational control of carbon emissions. Given the recent carbon emissions trends in the United States construction industry and projections for future emissions, it's imperative to develop appropriate environmental measures. Amidst soaring energy demand and escalating human activities, carbon dioxide emissions encounter a multitude of complexities, involving various factors and uncertainties. Consequently, policymakers necessitate more dependable prediction and assessment tools to thoroughly analyze and forecast carbon dioxide emissions [5]. Eventually, the purpose of this study is to establish a mathematical model with high prediction accuracy to obtain the predicted value of future carbon emissions of the construction industry.

## 2. Literature Review

The Random Forest method, a machine learning technique, was employed to build a predictive model aimed at establishing the correlation between carbon emissions during a building's construction phase and its design parameters. This predictive model, based on random forest methodology, aided designers in understanding the connection between the design characteristics of a building and the expected carbon emissions during its construction phase [6]. Some researchers argued that while econometric methods were commonly employed for regression analysis of carbon emissions and their influencing factors, they might not have fully captured the impact of these factors on carbon emissions. Therefore, these researchers proposed establishing a multi-factor prediction model using the system dynamics method to enhance the accuracy of predicting direct residential carbon emissions [7].

In addition, there were three main prediction methods for building carbon emissions at that time: combining the STIRPAT model with scenario analysis to test the prediction of carbon emissions [8], processing the GDP value to obtain the prediction value based on the cointegration relationship between GDP and carbon emissions [9], and using grey system theory, carbon emission prediction was obtained after analysis [5]. However, the accuracy of the predicted values produced by these methods was not high. Combining the grey system theory with Markov principle, this study gave full play to the advantages of the two methods, used Markov theory to determine the state transition probability, found out the characteristics of the carbon emission of the construction industry, and formed a suitable grey Markov prediction model. By calculating and improving the prediction accuracy, the future carbon emission prediction value of the building was obtained.

## 3. Grey prediction model GM (1,1)

The Grey prediction mode is a predictive method used for constructing mathematical models and making predictions when there's limited information available. It involves analyzing general differential equations to define Frey derivatives and Frey differential equations, and then using discrete data series to approximate dynamic models of these equations. The GM (1,1) model is built in the following steps [10]:

Let the input number of the original data be listed as $x^{(0)}$:

$$x^{(0)} = (x^{(0)}(1), x^{(0)}(2), \dots, x^{(0)}(n)) \tag{1}$$

The original series is accumulated to form a cumulative series $x^{(1)}$:

$$x^{(1)} = (x^{(0)}(1), x^{(1)}(2), \dots, x^{(1)}(n)) \tag{2}$$

$$x^{(1)}(k) = \sum_{i=1}^{k} x^{(0)}(i), k = 1, 2, \dots, n \tag{3}$$

Then the whitening differential of the original grey GM(1,1) prediction model is as follows:

$$\frac{dt^{(1)}(t)}{dt} + ax^{(1)}(t) = b \tag{4}$$

Where: a is the development coefficient and b is the gray action coefficient. By $\frac{dt^{(1)}(t)}{dt} + ax^{(1)}(t) = b$, and for equally spaced data $\Delta t=1$. So:

$$\Delta x^{(1)}(t) = x^{(1)}(t) - x^{(1)}(t-1) \tag{5}$$

The difference equation can be obtained:

$$x^{(0)}(t) + ax^{(1)}(t) = b \tag{6}$$

To make the result more reasonable, x 1 is modified to mean generating sequence z 1 .

$$z^{(1)}(k) = \frac{[x^{(1)}(k) + x^{(1)}(k+1)]}{2}, k = 2,3,\dots,n \tag{7}$$

The difference equation is changed to:

$$x^{(0)}(t) + az^{(1)}(t) = b \tag{8}$$

Substituting the schedule (t=2, 3,…,n) into the formula (8) is:

$$\begin{cases} x^{(0)}(2) + az^{(1)}(2) = b \\ x^{(0)}(3) + az^{(1)}(3) = b \\ \quad\vdots \\ x^{(0)}(n) + az^{(1)}(n) = b \end{cases} \tag{9}$$

Introduce matrix vector notation:

$$Y = \begin{bmatrix} x^{(0)}(2) \\ x^{(0)}(3) \\ \vdots \\ x^{(0)}(n) \end{bmatrix} \quad B = \begin{bmatrix} -z^{(1)}(2)1 \\ -z^{(1)}(3)1 \\ \vdots \quad \vdots \\ -z^{(1)}(n)1 \end{bmatrix} \quad u = \begin{bmatrix} a \\ b \end{bmatrix} \tag{10}$$

So the GM(1,1) model can be expressed as Y=Bu. The least square method is used to obtain the estimates of a and b.

$$\hat{u} = \begin{bmatrix} \hat{a} \\ \hat{b} \end{bmatrix} = (B^T B)^{-1} B^T Y \tag{11}$$

According to the initial condition x 1 1 = x 0 1 , the time response of the grey GM(1,1) model can be obtained by solving the differential equation:

$$\hat{x}^{(1)}(k+1) = \left(x^{(0)}(1) - \frac{b}{a}\right)e^{-ak} + \frac{b}{a}, k = 1,2,3,\dots,n \tag{12}$$

Then the predicted value obtained by this model is:

$$\hat{x}^{(0)}(k+1) = (x^{(0)}(1) - \frac{b}{a})(1 - e^a)e^{-ak} \tag{13}$$

k is the quantity related to the time series.

### 3.1. Case application
The data of total annual carbon emissions from buildings in the United States from 2002 to 2022 were selected for analysis. After calculation, the carbon emissions of buildings from 2002 to 2022 are shown in Table 1.

**Table 1.** 2002-2022 US carbon emissions of buildings

| Year | Carbon emission(t) | Year | Carbon emission(t) | Year | Carbon emission(t) |
|------|-------------------|------|-------------------|------|-------------------|
| 2002 | 633,204,033 | 2009 | 601,427,211 | 2016 | 542,719,551 |
| 2003 | 650,829,350 | 2010 | 596,394,077 | 2017 | 543,363,166 |
| 2004 | 646,306,527 | 2011 | 570,156,529 | 2018 | 596,926,784 |
| 2005 | 626,896,744 | 2012 | 506,638,815 | 2019 | 589,779,399 |
| 2006 | 577,605,473 | 2013 | 578,159,792 | 2020 | 547,222,571 |
| 2007 | 604,796,709 | 2014 | 593,271,035 | 2021 | 558,939,991 |
| 2008 | 609,816,806 | 2015 | 560,663,759 | 2022 | 584,672,52 |

The grey GM(1,1) model is established using the total $CO_2$ emission data from buildings in Table 1 and subsequent prediction calculation is carried out. Then the original input sequence $x^{(0)}$ is:

$x^{(0)}$=(633204033, 650829350, 646306527, 626896744, 577605473, 604796709, 609816806, 601427211, 596394077, 570156529, 506638815, 578159792, 593271035, 560663759, 542719551, 543363166, 596926784, 589779399, 547222571, 558939991, 584672528)

Then the sum once is listed as:

x(1)=(633204033, 1284033383, 1930339910, 2557236654, 3134842127, 3739638836, 4349455642, 4950882853, 5547276930, 6117433459, 6624072274, 7202232066, 7795503101, 8356166860, 8898886411, 9442249577, 10039176361, 10628955760, 11176178331, 11735118322, 12319790850)

By establishing the grey GM(1,1) model and solving it:

$$\hat{u} = [a, b]^T = [0.0065, 627036471.6717]^T \tag{14}$$

Then the GM(1,1) model predicted value of the total carbon dioxide emitted by buildings in the original data is:

$$\hat{x}^{(0)}(k+1) = (633204033 - \frac{627036471.6717}{0.0065})(1 - e^{0.0065})e^{-0.0065k} \tag{15}$$

The comparison between the predicted value of the grey model and the real value after the final calculation is shown in Table 2.

**Table 2.** GM(1,1) Comparison of model predicted values and true values

| Year | Actual value (t) | Prediction value (t) | Residual error | Relative error (%) |
|------|-----------------|---------------------|---------------|-------------------|
| 2002 | 633,204,033 | 633,204,033 | 0 | 0 |
| 2003 | 650,829,350 | 620,927,401 | 29901949 | 4.59 |
| 2004 | 646,306,527 | 616,922,072 | 29384455 | 4.55 |
| 2005 | 626,896,744 | 612,942,579 | 13954165 | 2.23 |
| 2006 | 577,605,473 | 608,988,755 | -31383282 | 5.43 |
| 2007 | 604,796,709 | 605,060,437 | -263728 | 0.04 |
| 2008 | 609,816,806 | 601,157,458 | 8659348 | 1.42 |
| 2009 | 601,427,211 | 597,279,655 | 4147556 | 0.69 |
| 2010 | 596,394,077 | 593,426,867 | 2967210 | 0.5 |
| 2011 | 570,156,529 | 589,598,931 | -19442402 | 3.41 |
| 2012 | 506,638,815 | 585,795,688 | -79156873 | 15.62 |
| 2013 | 578,159,792 | 582,016,977 | -3857185 | 0.67 |
| 2014 | 593,271,035 | 578,262,642 | 15008393 | 2.53 |
| 2015 | 560,663,759 | 574,532,523 | -13868764 | 2.47 |
| 2016 | 542,719,551 | 570,826,467 | -28106916 | 5.18 |
| 2017 | 543,363,166 | 567,144,316 | -23781150 | 4.38 |

**Table 2.** (continued).

| 2018 | 596,926,784 | 563,485,918 | 33440866 | 5.6 |
| 2019 | 547,222,571 | 559,851,118 | 29928281 | 5.07 |
| 2020 | 558,939,991 | 556,239,764 | -9017193 | 1.65 |
| 2021 | 584,672,528 | 552,651,706 | 6288285 | 1.13 |
| 2022 | 545,544,875 | 549,086,793 | 35585735 | 6.09 |
| 2023 | 542,025,805 | | | |
| 2024 | 589,779,399 | | | |
| 2025 | 538,529,435 | | | |
| 2026 | 535,055,618 | | | |
| 2027 | 531,604,210 | | | |
| Average relative error | | | | 3.49 |

### 3.2. Accuracy check

The accuracy of grey GM(1,1) model is tested by residual test and posterior test.

Average relative error:

$$\overline{\varphi} = \frac{1}{n}\sum\nolimits_{i=1}^{n} \varphi_i , \varphi_i = \frac{\Delta^{(0)}(i)}{x^{(0)}(i)} , \Delta^{(0)}(i) = \left|x^{(0)}(i) - \overline{x}^{(0)}(i)\right| \tag{16}$$

Mean square error ratio:

$$C = \frac{S_1}{S_2} , S_1 = \sqrt{\frac{1}{n}\sum\nolimits_{i=1}^{n}\left[x^{(0)}(i) - \overline{x}^{(0)}(i)\right]^2} , \overline{x}^{(0)} = \frac{1}{n}\sum\nolimits_{i=1}^{n} x^{(0)}(i) \tag{17}$$

$$S_2 = \sqrt{\frac{1}{n}\sum_{i=1}^{n}[\Delta^{(0)}(i) - \overline{\Delta}^{(0)}(i)]^2} , \overline{\Delta}^{(0)} = \frac{1}{n}\sum_{i=1}^{n}\Delta^{(0)}(i) \tag{18}$$

Small probability error:

$$P = P\left\{\left|\Delta^{(0)}(i) - \overline{\Delta}^{(0)}\right| < 0.06745S_1\right\} \tag{19}$$

According to the above formula, the accuracy test of the model shows that the average relative error of the model is only 0.0349, the mean square error ratio is 0.7500, and the small probability error is 0.6191.

**Table 3.** Grey model prediction accuracy test level

| Level | α | c |
|---|---|---|
| Level 1 (good) | >0.95 | <0.35 |
| Level 2 (qualified) | >0.80 | <0.45 |
| Level 3(simply qualified) | >0.70 | <0.50 |
| Level 4 (unqualified) | ≤0.70 | ≥0.65 |

After comparing Table 3, it is found that the accuracy level of relative error is level 2, the mean square error is, and the precision level of small probability error is level 4. Therefore, it is proved that the predicted value calculated by the grey GM(1,1) model is poorly fitted to the real value, and the accuracy of the prediction model is unqualified. Therefore, it is necessary to modify the residual error of the grey GM(1,1) model and construct a grey Markov prediction model with higher accuracy by constructing Markov chain to predict the future numerical changes.

## 4. Grey Markov prediction model

### 4.1. Establishment of grey Markov model

The prediction accuracy of the grey GM(1,1) model based on the above total $CO_2$ emission of buildings is poor. Next, the residual correction of the model is carried out, and the grey Markov prediction model is further established by constructing Markov chain.

The absolute value sequence of the residuals in Table 2 is the absolute value sequence of residuals $\varepsilon^{(0)}(k)$.

$$\varepsilon^{(0)}(k) = \left|x^{(0)}(k) - \bar{x}^{(0)}(k)\right| = \left\{\varepsilon^{(0)}(1), \varepsilon^{(0)}(2), \dots, \varepsilon^{(0)}(n)\right\} \tag{20}$$

The one-time cumulative sequence $\varepsilon^{(1)}(k)$ is:

$$\varepsilon^{(1)}(k) = \left\{\varepsilon^{(1)}(1), \varepsilon^{(1)}(2), \dots, \varepsilon^{(1)}(n)\right\} \tag{21}$$

By establishing GM(1,1) model for $\varepsilon^{(1)}(k)$, the differential equation is as follows:

$$\frac{d\varepsilon^{(1)}(k)}{dt} + a_1\varepsilon^{(1)}(k) = b_1 \tag{22}$$

Finally, the GM(1,1) model predicted value of the absolute residual data is as follows:

$$\hat{\varepsilon}^{(0)}(k + 1) = (\varepsilon^{(0)}(1) - \frac{b_1}{a_1})(1 - e^{a_1})e^{-a_1 k} \tag{23}$$

Then the grey GM(1,1) prediction model after residual correction is:

$$\hat{x}^{(0)}(k + 1) = \left(x^{(0)}(1) - \frac{b}{a}\right)(1 - e^a)e^{-ak} + sng(k+1)(\varepsilon^{(0)}(1) - \frac{b_1}{a_1})(1 - e^{a_1})e^{-a_1 k}$$

$$sng(k + 1) = \begin{cases} -1, x^{(0)}(k) - \hat{x}^{(0)}(k) < 0 \\ 0, x^{(0)}(k) - \hat{x}^{(0)}(k) = 0 \\ 1, x^{(0)}(k) - \hat{x}^{(0)}(k) > 0 \end{cases} \tag{24}$$

The specific situation of sgn(k+1) is determined by the residual difference between the true value and the predicted value.

### 4.2. Case calculation

The comparison between the predicted value of the grey GM(1,1) model with residual correction and the true value is shown in Table 4.

**Table 4.** Comparison of the predicted value with the true value after the residual correction

| Year | Actual value (t) | Prediction value (t) | Residual error | Relative error (%) |
|------|------------------|----------------------|----------------|--------------------|
| 2002 | 633,204,033 | 633,204,033 | 0 | 0.00 |
| 2003 | 650,829,350 | 639,984,492 | -10844858 | 1.69 |
| 2004 | 646,306,527 | 636,162,596 | -10143931 | 1.59 |
| 2005 | 626,896,744 | 632,368,302 | 5471558 | 0.87 |
| 2006 | 577,605,473 | 589,376,051 | 11770578 | 2.00 |
| 2007 | 604,796,709 | 585,258,951 | -19537758 | 3.34 |
| 2008 | 609,816,806 | 621,149,542 | 11332736 | 1.82 |
| 2009 | 601,427,211 | 617,464,172 | 16036961 | 2.60 |
| 2010 | 596,394,077 | 613,805,669 | 17411592 | 2.84 |
| 2011 | 570,156,529 | 569,023,974 | -1132555 | 0.20 |
| 2012 | 506,638,815 | 565,022,687 | 58383872 | 10.33 |

**Table 4.** (continued).

| 2013 | 578,159,792 | 561,044,026 | -17115766 | 3.05 |
| 2014 | 593,271,035 | 599,437,467 | 6166432 | 1.03 |
| 2015 | 560,663,759 | 553,153,881 | -7509878 | 1.36 |
| 2016 | 542,719,551 | 549,242,046 | 6522495 | 1.19 |
| 2017 | 543,363,166 | 545,352,135 | 1988969 | 0.36 |
| 2018 | 596,926,784 | 585,487,859 | -11438925 | 1.95 |
| 2019 | 589,779,399 | 582,064,838 | -7714561 | 1.33 |
| 2020 | 547,222,571 | 533,812,227 | -13410344 | 2.51 |
| 2021 | 558,939,991 | 575,295,118 | 16355127 | 2.84 |
| 2022 | 584,672,528 | 571,948,159 | -12724369 | 2.22 |
| Average relative error | | | | 2.15 |

The accuracy test of the established grey Markov model shows that the average relative error of the model is only 0.0215, the mean square ratio is 0.4314, and the small probability error is 0.9524.

Through comparison, it is found that the accuracy level of the average relative error and the mean square error ratio is two levels, and the accuracy level of the small probability error is one level, so it is proved that the established grey Markov model is feasible and effective, and the prediction accuracy is higher than that of the grey GM(1,1) model. Then the predicted value of the next five years is obtained by establishing the state transition matrix of Markov chain.

### 4.3. Establishment of state

At k<n, the value of sgn(k) can be obtained from the difference between the real value and the predicted value. At k>n, the value of sgn(k) needs to be calculated by a Maldivian chain formed by a state transition probability structure with a model, where the state transition matrix P is:

$$P = \begin{bmatrix} P_{11}P_{12}\cdots P_{1n} \\ P_{21}P_{22}\cdots P_{2n} \\ \cdots \cdots \cdots \cdots \\ P_{n1}P_{n2}\cdots P_{nn} \end{bmatrix} \tag{25}$$

The calculation of probability can make the frequency value approximate to the probability value:

$$P_{ij} = \frac{M_{ij}}{M_i} \tag{26}$$

$M_{ij}$ is the number of times that state $E_i$ transfers to state $E_j$; $M_i$ is the total number of occurrences of state $E_i$.

According to the variation of residuals in Table 2, set $S_1$ to represent the state of positive residuals, $S_2$ to represent the state of zero residuals, and $S_3$ to represent the state of negative residuals. Get from 2002 to 2022 of the state is: $[S_2,S_1,S_1,S_1,S_3,S_3,S_1,S_1,S_1,S_3,S_3,S_3,S_1,S_3,S_3,S_3,S_1,S_1,S_3,S_1,S_1]$, so obtained:

$$P_{11} = \frac{M_{11}}{M_1} = \frac{3}{5}, P_{12} = \frac{M_{12}}{M_1} = 0, P_{13} = \frac{M_{13}}{M_1} = \frac{2}{5} \tag{27}$$

$$P_{21} = \frac{M_{21}}{M_2} = 1, P_{22} = \frac{M_{22}}{M_2} = 0, P_{23} = \frac{M_{23}}{M_2} = 0 \tag{28}$$

$$P_{31} = \frac{M_{31}}{M_3} = \frac{4}{9}, P_{32} = \frac{M_{32}}{M_3} = 0, P_{33} = \frac{M_{33}}{M_3} = \frac{5}{9} \tag{29}$$

Then the state transition matrix is:

$$\begin{bmatrix} 3/5 & 0 & 2/5 \\ 1 & 0 & 0 \\ 4/9 & 0 & 5/9 \end{bmatrix}$$

By using the established transfer matrix and the data in Table 2 and Table 4, the grey Markov prediction model is established to predict the total $CO_2$ emissions from buildings in 2023-2027. In addition, since the state of total carbon oxide in 2012 is $S_1$, we can observe the first row of the transition state matrix P, and the maximum $P_{ijMAX}=P_{11}$ is 3/5, then the next state is most likely to change from $S_1$ to $S_1$. Therefore, the final 2023-2027 forecast value is the gray forecast value +2/3× residual forecast value. The results are shown in Table 5.

**Table 5.** The subsequent predicted value of the grey Markov model

| Year | Prediction value (t) |
| --- | --- |
| 2023 | 559,393,725 |
| 2024 | 556,007,957 |
| 2025 | 552,646,171 |
| 2026 | 549,308,234 |
| 2027 | 545,994,014 |

The actual value is compared with the predicted value of grey GM(1,1) and the predicted value of grey Markov model, and the image is drawn by using MATLAB software, as shown in Figure 1.
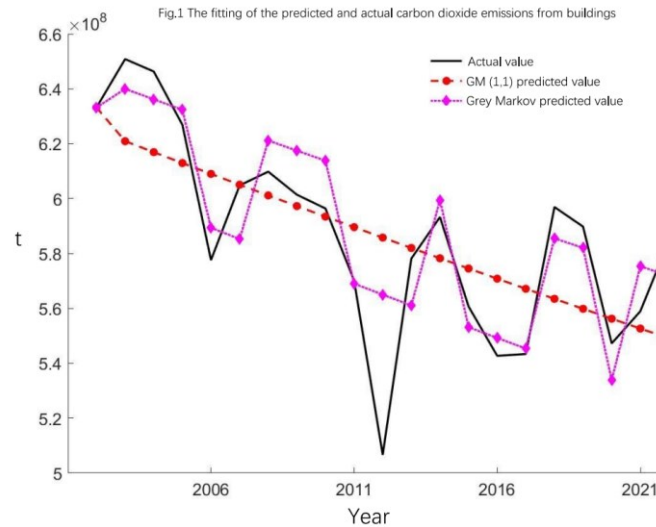


**Figure 1.** The fitting of the predicted and actual carbon dioxide emissions from buildings

## 5. Conclusion

Based on the $CO_2$ emission data of buildings in the United States from 2002 to 2022 published on the official website of the European Union, this paper compares the prediction value of the GM(1,1) model with that of the GM(1,1) model under Markov optimization, and finds that the predicted value of the GM(1,1) model under Markov optimization is closer to the actual value. On this basis, the data of $CO_2$ emissions from buildings in the United States from 2023 to 2027 are predicted, which in 2023 was 559,393,725 tons, in 2024 was 556,007,957 tons, in 2025 was 552,646,171 tons, in 2026 was 549,308,234 tons, in 2027 was 545,994,014 tons. The forecast results show that the carbon dioxide emissions of buildings in the United States will continue to decline in the next few years, which is basically in line with the long-term development trend of energy conservation and emission reduction in the building industry in the United States. The results also confirm that the model proposed in the article has a certain degree of feasibility, and can improve the accuracy of carbon dioxide emissions detection.

However, there are some shortcomings in the experimental process of this paper, such as only choosing the United States as a reference, not carefully categorizing and analyzing buildings, and the relatively small amount of data. The depth of analysis and the testing of the model will be further optimized and upgraded based on the in-depth study in this direction.

**References**

[1] Javanmard, M. E., Ghaderi, S. F., & Hoseinzadeh, M. (2021). Data mining with 12 machine learning algorithms for predict costs and carbon dioxide emission in integrated energy-water optimization model in buildings. Energy Conversion and Management, 238, 114153.

[2] Ahmed Ali, K., Ahmad, M. I., & Yusup, Y. (2020). Issues, impacts, and mitigations of carbon dioxide emissions in the building sector. Sustainability, 12(18), 7427.

[3] Atmaca, A., & Atmaca, N. (2016). Comparative life cycle energy and cost analysis of post-disaster temporary housings. Applied energy, 171, 429-443.

[4] Atmaca, N., Atmaca, A., & Özçetin, A. İ. (2021). The impacts of restoration and reconstruction of a heritage building on life cycle energy consumption and related carbon dioxide emissions. Energy and Buildings, 253, 111507.

[5] Zhou, W., Zeng, B., Wang, J., Luo, X., & Liu, X. (2021). Forecasting Chinese carbon emissions using a novel grey rolling prediction model. Chaos, Solitons & Fractals, 147, 110968.

[6] Fang, Y., Lu, X., & Li, H. (2021). A random forest-based model for the prediction of construction-stage carbon emissions at the early design stage. Journal of Cleaner Production, 328, 129657.

[7] Xie. Z., Gao. X., Yuan. W., Fang. J., and Jiang. Z, (2020). Decomposition and prediction of direct residential carbon emission indicators in Guangdong Province of China. Ecological Indicators, 115, 106344.

[8] Wu, R., Wang, J., Wang, S., & Feng, K. (2021). The drivers of declining CO2 emissions trends in developed nations using an extended STIRPAT model: A historical and prospective analysis. Renewable and Sustainable Energy Reviews, 149, 111328..

[9] Kumar, S., Shukla, A. K., & Muhuri, P. K. (2021). Anomaly based novel multi-source unsupervised transfer learning approach for carbon emission centric GDP prediction. Computers in Industry, 126, 103396.

[10] Jia, Z. Q., Zhou, Z. F., Zhang, H. J., Li, B., & Zhang, Y. X. (2020). Forecast of coal consumption in Gansu Province based on Grey-Markov chain model. Energy, 199, 117444.