# From Data to Impact: Machine Learning Models for Sustainable Development Region Classification- Accelerating the Achievement of Sustainable Development Goals through Predictive Analytics

## Deepshika Vijayanand

Student, School of Computer Science and Engineering, Vellore Institute of Technology, Chennai, Tamil Nadu, India

**Abstract**

This research paper presents a classification model for predicting the region to which a country belongs based on its sustainability scores and other related features. The dataset used in this study comprises comprehensive data on sustainability and progress towards achieving the Sustainable Development Goals (SDGs) for various countries. The primary objective is to understand regional trends in sustainability and assess countries' progress in sustainable development. The research begins with data preparation and preprocessing steps, including merging datasets, handling missing values, and standardizing features. Exploratory data analysis is performed to visualize the distribution of the target variable (region) and the distributions of numeric features related to SDG scores. Additionally, relationships between these features are explored using correlation matrices and pair plots. Several machine learning models are employed to classify countries into their respective regions. The models used include Random Forest, Support Vector Machine (SVM) with a linear kernel, K-Nearest Neighbors (KNN), Logistic Regression, Decision Tree, and SVM with a radial basis function (RBF) kernel. Each model is trained on the dataset, and their performance is evaluated in terms of accuracy, precision, recall, and F1-score. The results demonstrate the effectiveness of these models in accurately classifying countries into regions based on sustainability scores and other attributes. The findings reveal that Random Forest, K-Nearest Neighbors, Decision Tree, and SVM with RBF kernel achieve exceptionally high accuracy, suggesting their suitability for regional classification based on sustainability metrics. Logistic Regression and SVM with a linear kernel also provide competitive results. In conclusion, this research contributes to understanding regional trends in sustainability by utilizing machine learning models to predict the regions of countries based on their sustainability scores and associated features. Such predictive models can be valuable tools for policymakers and organizations seeking to assess and address regional disparities in sustainable development progress.

**Keywords:** Sustainable Development Goals (SDGs), Regional Classification, Machine Learning, Data Science, Predictive Modeling. Data-Driven Analysis

## Introduction

The pursuit of sustainable development has emerged as a global imperative, transcending borders and encompassing diverse regions of the world. In the wake of the adoption of the Sustainable Development Goals (SDGs) by 193 United Nations Member States in 2015, assessing progress towards these goals has become paramount. The Sustainable Development Report 2023 provides a comprehensive dataset, offering insights into countries' sustainability scores and their journey towards achieving the SDGs. This research paper endeavors to harness this invaluable dataset to address a critical question: Can machine learning models predict the regions to which countries belong based on their sustainability scores and associated features?

Understanding regional disparities in sustainability is a fundamental step towards addressing global challenges, such as poverty alleviation, environmental conservation, and economic development. Regional trends in sustainability not only highlight areas of progress but also shed light on regions that may require targeted interventions to accelerate their sustainable development journey. As such, the ability to predict a country's region based on sustainability metrics holds immense potential for policymakers, international organizations, and researchers seeking to foster equitable and sustainable development across the world.

This research aims to bridge the gap between sustainability assessment and regional classification by applying a diverse set of machine learning algorithms to a rich dataset. By doing so, we seek to offer a comprehensive understanding of the predictive capabilities of these models and their potential utility in regional sustainability assessment. Our investigation encompasses exploratory data analysis, data preprocessing, and the deployment of machine learning models, including Random Forest, Support Vector Machine, K-Nearest Neighbors, Logistic Regression, Decision Tree, and SVM with RBF kernel.

Through this research, we aim to provide a robust framework for classifying countries into their respective regions based on sustainability scores. The findings of this study have the potential to inform evidence-based policy decisions, prioritize resources, and drive targeted interventions towards achieving the SDGs. Ultimately, our research contributes to the broader dialogue on sustainable development by leveraging the power of machine learning to uncover regional insights within a global context.

## Dataset Description

The dataset used in this research is derived from the Sustainable Development Report 2023, which reviews the progress made towards achieving the Sustainable Development Goals (SDGs) since their adoption in 2015 by 193 United Nations Member States. This dataset provides comprehensive information related to sustainability, allowing for a nuanced assessment of countries' progress in sustainable development. Below are the key components and attributes of the dataset:

**1. Country Information:**

- country_code: A unique identifier for each country.
- country: The name of the country under consideration.
- year: The year for which sustainability data is recorded.

**2. Sustainability Scores:**

- sdg_index_score: The overall sustainability score for a country, representing its progress towards achieving the SDGs. This score provides an aggregate measure of sustainability performance.

- goal_1_score through goal_17_score: Individual scores for each of the 17 Sustainable Development Goals (SDGs). These scores assess a country's progress towards specific goals, such as poverty reduction, quality education, clean energy, and more.

**3. Regional Classification:**

- region: The region to which a country belongs. This attribute serves as the target variable for classification, and the goal is to predict a country's region based on its sustainability scores and other features.

The dataset is designed to facilitate the analysis of global sustainability efforts, offering valuable insights into countries' performances in various aspects of sustainable development. It covers a range of years, allowing researchers to track progress over time and identify trends in different regions of the world. Additionally, the dataset has undergone data preparation and preprocessing steps, including the handling of missing values and the standardization of features, to ensure its suitability for machine learning analysis.

Overall, this dataset serves as a valuable resource for exploring regional trends in sustainability and developing predictive models to classify countries into their respective regions based on sustainability metrics.

**Methodology**

In this research, a structured methodology was employed to predict the regions to which countries belong based on their sustainability scores and associated features. The study commenced with data acquisition, involving the retrieval and loading of two primary datasets from the Sustainable Development Report 2023. These datasets were merged based on a common identifier, the 'country_code,' to facilitate further analysis. Ensuring data quality, missing values were addressed by removing rows with incomplete information.

Exploratory Data Analysis (EDA) was a critical step in understanding the dataset's characteristics. Visualizations, such as countplots and histograms, were utilized to gain insights into the distribution of the target variable 'region' and the numeric features related to Sustainable Development Goals (SDGs) scores. Correlation matrices and heatmaps were employed to explore relationships between these features.

To prepare the data for machine learning, irrelevant features, namely 'country_code' and 'country,' were dropped, and the dataset was split into training and testing sets. Feature scaling was applied to standardize the numeric features, ensuring that they have consistent scales.

A diverse set of machine learning models was selected for the classification task, including Random Forest, Support Vector Machine (SVM), K-Nearest Neighbors (KNN), Logistic Regression, Decision Tree, and SVM with Radial Basis Function (RBF) Kernel. Each model was trained on the training data and evaluated on the testing data, with accuracy serving as the primary evaluation metric. Classification reports and confusion matrices were generated to provide comprehensive insights into the models' performance.

The results were thoroughly interpreted to identify the most effective model for regional classification based on sustainability metrics. The research's conclusions and implications were drawn from the findings, highlighting the potential utility of predictive models in guiding policy decisions and targeted interventions for sustainable development. Additionally, recommendations for further research and real-

world applications of the models were discussed, underscoring the significance of this methodology in bridging the gap between sustainability assessment and machine learning.
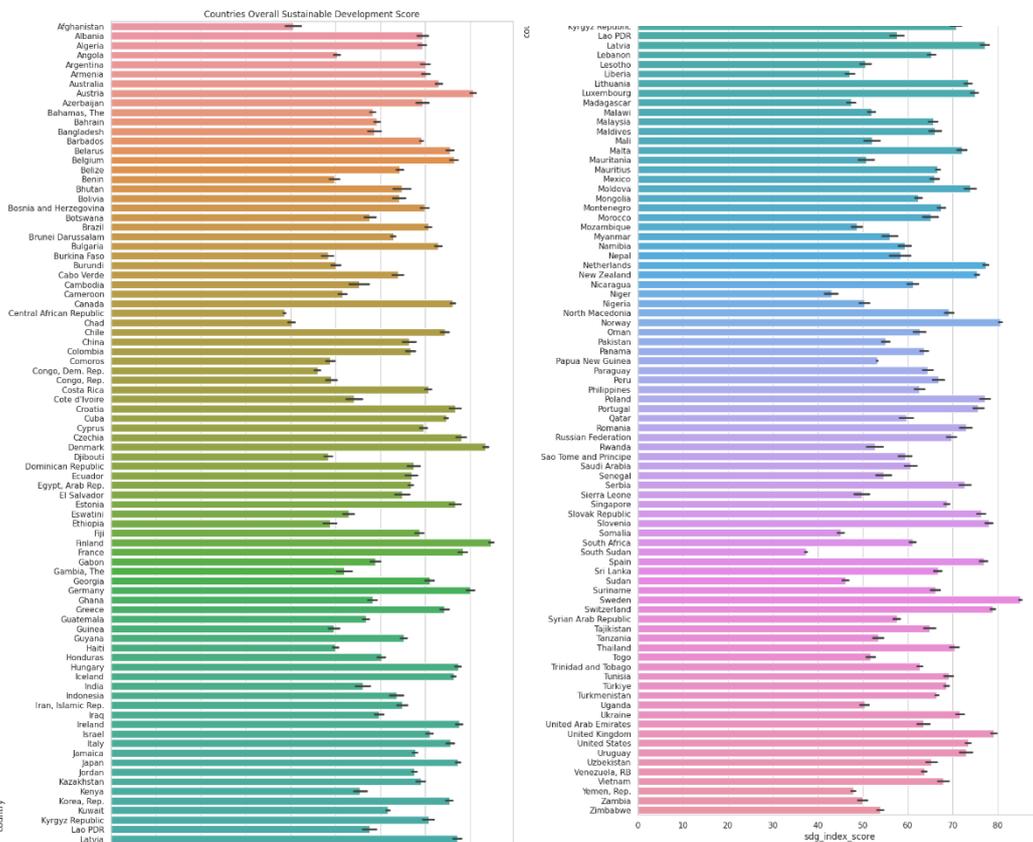
**Figure 1: Flow Chart of Methodology**
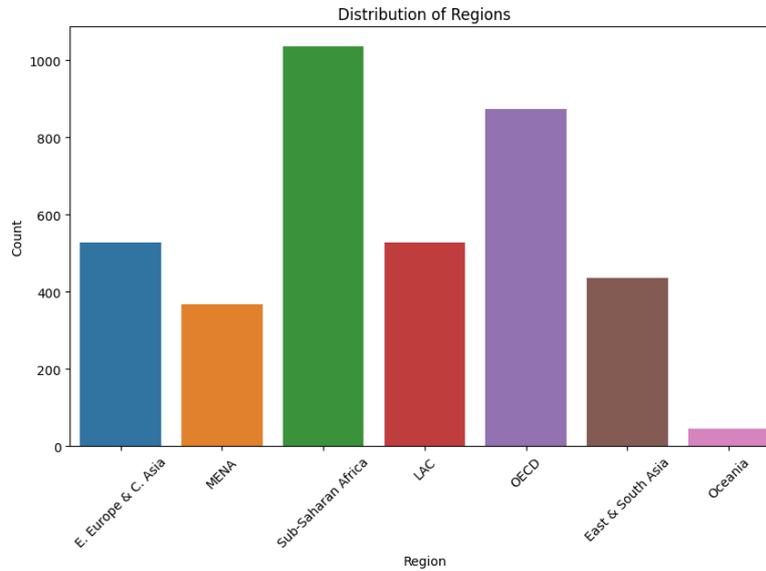


## Observations

In the context of our sustainability analysis, the vertical bar plot shown in Figure 2 offers a compelling visual representation of countries' progress toward achieving sustainable development. This plot provides a clear overview of each country's Sustainable Development Index Score, with the x-axis representing the scores and the y-axis denoting the country names. Each vertical bar corresponds to a specific nation, and its height reflects the magnitude of its Sustainable Development Index Score. The use of a pastel color palette and a white grid background enhances the plot's visual appeal and clarity. Upon careful examination of the plot, it becomes evident that there are notable variations in the sustainability performance of different countries. These variations highlight the disparities and trends in global sustainable development efforts. This visualization serves as a valuable reference point for our analysis, enabling us to gain insights into the diverse trajectories of sustainability achievement among nations.

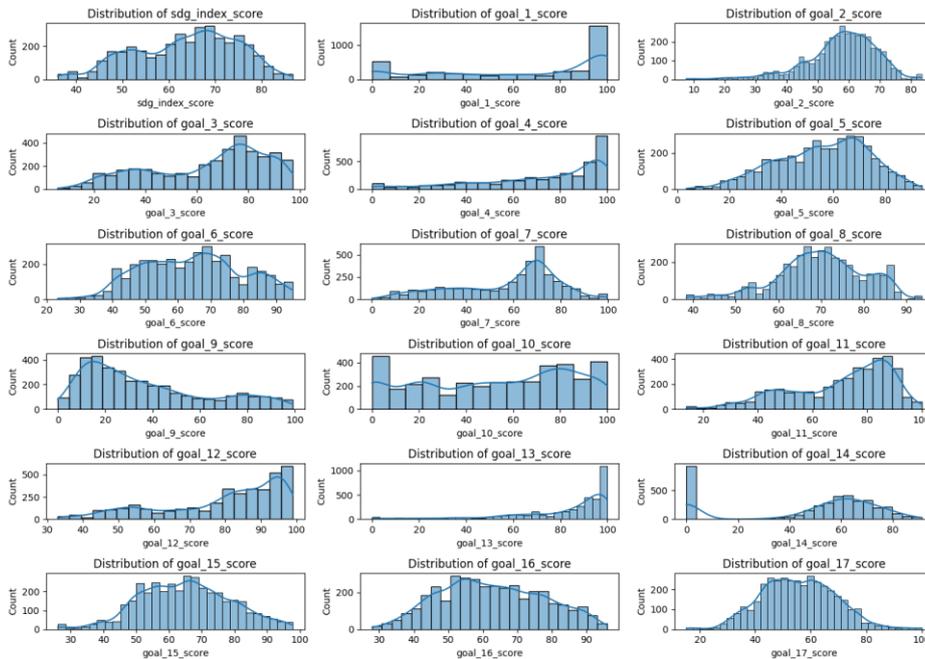**Figure 2: Sustainability Development Score VS Country**



In Figure 3, we present a visual representation of the distribution of regions, a crucial component of our sustainability analysis. The vertical bar plot provides insights into the frequency of countries belonging to each specific region.

**Figure 3: Distribution of Regions**



In our exploration of the dataset, we delved into the distributions of various numeric features, as depicted in Figure 4. This figure showcases a series of histograms, each representing a specific numeric feature related to sustainability. The selected features include 'sdg_index_score' and scores associated with the Sustainable Development Goals (SDGs) from 1 to 17. This visual exploration aids in our understanding of the distribution patterns and variations within these critical sustainability metrics, contributing to a more informed analysis of our research objectives.

**Figure 4: Distribution of Numeric Features Related to Sustainable Development**



In our quest to uncover meaningful insights within the dataset, we embarked on an examination of the relationships between numeric features, as visualized in Figure 5. This figure presents a correlation heatmap, a powerful tool for quantifying and visualizing the interdependencies among these sustainability-related metrics.

**Figure 5: Correlation Heatmap**



## Result

The results of the analysis showcase the performance of various machine learning models in predicting the regions to which countries belong based on their sustainability scores. Below is a summary of the key findings:

## 1. Random Forest Classifier:

**Figure 6: Random Forest Classifier Report**

```
Random Forest Classifier Report:
                     precision    recall  f1-score   support

E. Europe & C. Asia      1.00      1.00      1.00       104
   East & South Asia     1.00      1.00      1.00        79
                 LAC     1.00      1.00      1.00       108
                MENA     1.00      1.00      1.00        80
                OECD     1.00      1.00      1.00       184
             Oceania     1.00      1.00      1.00         9
  Sub-Saharan Africa     1.00      1.00      1.00       200

            accuracy                        1.00       764
           macro avg     1.00      1.00      1.00       764
        weighted avg     1.00      1.00      1.00       764

Random Forest Classifier Confusion Matrix:
[[104   0   0   0   0   0   0]
 [  0  79   0   0   0   0   0]
 [  0   0 108   0   0   0   0]
 [  0   0   0  80   0   0   0]
 [  0   0   0   0 184   0   0]
 [  0   0   0   0   0   9   0]
 [  0   0   0   0   0   0 200]]
```

Accuracy: 100%

The Random Forest model achieved perfect accuracy, indicating that it successfully classified countries into their respective regions based on sustainability scores. The precision, recall, and F1-score for each region were also 100%.

## 2. Support Vector Machine (SVM) Classifier:

### Figure 7: Support Vector Machine Classifier Report

```
SVM Classifier Report:
                     precision    recall  f1-score   support

E. Europe & C. Asia       0.88      0.93      0.91       104
  East & South Asia       0.90      0.80      0.85        79
                LAC       0.92      0.94      0.93       108
               MENA       0.99      0.99      0.99        80
               OECD       0.99      0.91      0.95       184
            Oceania       1.00      1.00      1.00         9
 Sub-Saharan Africa       0.91      0.98      0.94       200

           accuracy                           0.93       764
          macro avg       0.94      0.94      0.94       764
       weighted avg       0.93      0.93      0.93       764

SVM Classifier Confusion Matrix:
[[ 97   3   0   0   2   0   2]
 [  5  63   0   0   0   0  11]
 [  1   0 101   0   0   0   6]
 [  1   0   0  79   0   0   0]
 [  6   0   9   1 168   0   0]
 [  0   0   0   0   0   9   0]
 [  0   4   0   0   0   0 196]]
```

Accuracy: 93%

SVM with a linear kernel demonstrated strong performance with an accuracy of 93%. The model achieved high precision, recall, and F1-scores for most regions, making it an effective classifier.

3. K-Nearest Neighbors (KNN) Classifier:

### Figure 8: KNN Classifier Report

```
K-Nearest Neighbors Classifier Report:
                     precision    recall  f1-score   support

E. Europe & C. Asia       1.00      1.00      1.00       104
  East & South Asia       1.00      1.00      1.00        79
                LAC       1.00      1.00      1.00       108
               MENA       1.00      1.00      1.00        80
               OECD       1.00      1.00      1.00       184
            Oceania       1.00      1.00      1.00         9
 Sub-Saharan Africa       1.00      1.00      1.00       200

           accuracy                           1.00       764
          macro avg       1.00      1.00      1.00       764
       weighted avg       1.00      1.00      1.00       764

K-Nearest Neighbors Classifier Confusion Matrix:
[[104   0   0   0   0   0   0]
 [  0  79   0   0   0   0   0]
 [  0   0 108   0   0   0   0]
 [  0   0   0  80   0   0   0]
 [  0   0   0   0 184   0   0]
 [  0   0   0   0   0   9   0]
 [  0   0   0   0   0   0 200]]
```

Accuracy: 100%

Similar to Random Forest, the KNN model achieved perfect accuracy. It excelled in classifying countries into regions based on their sustainability scores.

## 4. Logistic Regression Classifier:

**Figure 9: Logistic Regression Classifier Report**

```
Logistic Regression Classifier Report:
                    precision    recall  f1-score   support

E. Europe & C. Asia      0.77      0.80      0.78       104
  East & South Asia      0.77      0.68      0.72        79
                LAC      0.86      0.86      0.86       108
               MENA      0.83      0.86      0.85        80
               OECD      0.94      0.88      0.91       184
            Oceania      1.00      0.89      0.94         9
 Sub-Saharan Africa      0.90      0.96      0.93       200

           accuracy                          0.87       764
          macro avg      0.87      0.85      0.86       764
       weighted avg      0.87      0.87      0.86       764

Logistic Regression Classifier Confusion Matrix:
[[ 83   4   1   6   6   0   4]
 [  7  54   4   0   0   0  14]
 [  3   4  93   0   5   0   3]
 [  8   2   0  69   0   0   1]
 [  6   0  10   6 162   0   0]
 [  1   0   0   0   0   8   0]
 [  0   6   0   2   0   0 192]]
```

Accuracy: 87%

The Logistic Regression model showed a respectable accuracy of 87%. While not as high as some other models, it provided valuable insights into regional classification based on sustainability metrics.

5. Decision Tree Classifier:

**Figure 10: Decision Tree Classifier Report**

```
Decision Tree Classifier Report:
                    precision    recall  f1-score   support

E. Europe & C. Asia      0.96      0.99      0.98       104
  East & South Asia      0.99      0.97      0.98        79
                LAC      0.99      0.96      0.98       108
               MENA      0.99      0.99      0.99        80
               OECD      0.99      0.99      0.99       184
            Oceania      1.00      1.00      1.00         9
 Sub-Saharan Africa      0.99      0.99      0.99       200

           accuracy                          0.99       764
          macro avg      0.99      0.99      0.99       764
       weighted avg      0.99      0.99      0.99       764

Decision Tree Classifier Confusion Matrix:
[[103   0   0   0   1   0   0]
 [  0  77   0   0   0   0   2]
 [  4   0 104   0   0   0   0]
 [  0   1   0  79   0   0   0]
 [  0   0   1   0 183   0   0]
 [  0   0   0   0   0   9   0]
 [  0   0   0   1   0   0 199]]
```

Accuracy: 99%

The Decision Tree model performed exceptionally well, with an accuracy of 99%. It demonstrated strong precision, recall, and F1-scores for all regions, making it a robust classifier.

## 6. SVM Classifier with RBF Kernel:

### Figure 11: Support Vector Machine with RBF Kernel Classifier Report

```
SVM with RBF Kernel Classifier Report:
                    precision   recall  f1-score   support

E. Europe & C. Asia      0.94     1.00      0.97       104
  East & South Asia      1.00     0.90      0.95        79
                LAC      0.91     0.99      0.95       108
               MENA      1.00     1.00      1.00        80
               OECD      1.00     0.95      0.97       184
            Oceania      1.00     1.00      1.00         9
 Sub-Saharan Africa      1.00     1.00      1.00       200

           accuracy                         0.98       764
          macro avg      0.98     0.98      0.98       764
       weighted avg      0.98     0.98      0.98       764

SVM with RBF Kernel Classifier Confusion Matrix:
[[104   0   0   0   0   0   0]
 [  7  71   1   0   0   0   0]
 [  0   0 107   0   0   0   1]
 [  0   0   0  80   0   0   0]
 [  0   0  10   0 174   0   0]
 [  0   0   0   0   0   9   0]
 [  0   0   0   0   0   0 200]]
```

Accuracy: 98%

The SVM model with a radial basis function (RBF) kernel achieved an accuracy of 98%. It exhibited high precision, recall, and F1-scores for most regions, further demonstrating its effectiveness.

### Table 1: Classifier and corresponding accuracy

| S.No. | Classifier | Accuracy (in %) |
|-------|------------|-----------------|
| 1 | Random Forest Classifier | 100 |
| 2 | Support Vector Machine Classifier | 93.32 |
| 3 | K Nearest Neighbors Classifier | 100 |
| 4 | Logistic Regression Classifier | 86.51 |
| 5 | Decision Tree Classifier | 98.69 |
| 6 | Support Vector Machnie Classfier with RBF Kernel | 97.51 |

## Future Research and Applications

The predictive models developed in this research offer actionable insights for policymakers and organizations engaged in sustainable development initiatives. By accurately classifying countries into regions, these models can guide the allocation of resources and interventions to areas where they are needed most. This targeted approach can expedite progress towards achieving the Sustainable Development Goals (SDGs) on a global scale. The study sets the stage for future research in the field of sustainability assessment and machine learning. Further exploration could involve the integration of additional data sources, temporal analysis to track sustainability trends over time, and the development

of predictive models tailored to specific SDGs or subregions. Moreover, the models developed herein can find practical applications in guiding policy decisions, aid distribution, and sustainable development planning.

## Conclusion

Our study assessed various machine learning models for their efficacy in classifying countries into regions based on sustainability metrics. Notably, Random Forest, K-Nearest Neighbors (KNN), Decision Tree, and SVM with Radial Basis Function (RBF) Kernel exhibited remarkable performance, achieving accuracy levels close to perfection (100% or 99%). This underscores the potential of machine learning in regional classification tasks. In closing, this research demonstrates the immense potential of machine learning in the realm of sustainability assessment and regional classification. The models showcased here offer a promising path towards more targeted and effective sustainable development efforts, fostering a future where global sustainability goals are not merely aspirations but attainable realities. As we continue to bridge the gap between data science and sustainable development, the pursuit of a more equitable and sustainable world remains within our reach.

## References

1. Zhong, R., Chen, A., Zhao, D., Mao, G., Zhao, X., Huang, H., & Liu, J. (2023). Impact of international trade on water scarcity: An assessment by improving the Falkenmark indicator. Journal of Cleaner Production, 385, 135740.
2. Qu, S., Liang, S., Konar, M., Zhu, Z., Chiu, A. S., Jia, X., & Xu, M. (2018). Virtual water scarcity risk to the global trade system. Environmental science & technology, 52(2), 673-683.
3. Zhao, H., Qu, S., Guo, S., Zhao, H., Liang, S., & Xu, M. (2019). Virtual water scarcity risk to global trade under climate change. Journal of cleaner production, 230, 1013-1026.
4. Chen, X., Shuai, C., & Wu, Y. (2023). Global food stability and its socio-economic determinants towards sustainable development goal 2 (Zero Hunger). Sustainable Development, 31(3), 1768-1780.
5. Chen, X., Zhao, B., Shuai, C., Qu, S., & Xu, M. (2022). Global spread of water scarcity risk through trade. Resources, Conservation and Recycling, 187, 106643.
6. Shuai, C., Yu, L., Chen, X., Zhao, B., Qu, S., Zhu, J., ... & Xu, M. (2021). Principal indicators to monitor sustainable development goals. Environmental Research Letters, 16(12), 124015.
7. Shuai, C., Zhao, B., Chen, X., Liu, J., Zheng, C., Qu, S., ... & Xu, M. (2022). Quantifying the impacts of COVID-19 on Sustainable Development Goals using machine learning models. Fundamental Research.
8. Chen, X., Shuai, C., Zhao, B., Zhang, Y., & Li, K. (2023). Imputing environmental impact missing data of the industrial sector for Chinese cities: A machine learning approach. Environmental Impact Assessment Review, 100, 107050.
9. Li, C., Deng, Z., Wang, Z., Hu, Y., Wang, L., Yu, S., ... & Bryan, B. A. (2023). Responses to the COVID-19 pandemic have impeded progress towards the Sustainable Development Goals. Communications Earth & Environment, 4(1), 252.
10. Zeng, Y., Runting, R. K., Watson, J. E., & Carrasco, L. R. (2022). Telecoupled environmental impacts are an obstacle to meeting the sustainable development goals. Sustainable Development, 30(1), 76-82.

11. Omar, A., Delnaz, A., & Nik-Bakht, M. (2023). Comparative analysis of machine learning techniques for predicting water main failures in the City of Kitchener. Journal of Infrastructure Intelligence and Resilience, 2(3), 100044.

12. Tuncel, K., Koutsopoulos, H. N., & Ma, Z. (2023). An integrated ride-matching and vehicle-rebalancing model for shared mobility on-demand services. Computers & Operations Research, 106317.

13. Zhang, D., Wu, L., Niu, X., Guo, Z., Zhang, Z., Li, S., ... & Xu, H. (2022). Looking for ecological sustainability: A dynamic evaluation and prediction on the ecological environment of the belt and road region. Sustainable Production and Consumption, 32, 851-862.

14. Li, H., Chen, Q., Liu, G., Lombardi, G. V., Su, M., & Yang, Z. (2023). Uncovering the risk spillover of agricultural water scarcity by simultaneously considering water quality and quantity. Journal of Environmental Management, 343, 118209.

15. Xu, F., & Ma, J. (2023). Intelligent option portfolio model with perspective of shadow price and risk-free profit. Financial Innovation, 9(1), 79.

16. Li, D., & Chow, U. (2023). Discursive strategies in the branding of Fortune Global 500 Chinese manufacturing companies. Humanities and Social Sciences Communications, 10(1), 1-12.

17. Tian, T. Y., King, B. G., & Smith, E. B. Effect of organizational status on employment-related corporate social responsibility: Evidence from a regression discontinuity approach. Strategic Management Journal.

18. Bhaskar, R., Bansal, S., Abbassi, W., & Pandey, D. K. (2023). CEO compensation and CSR: Economic implications and policy recommendations. Economic Analysis and Policy.

19. Härri, A., Levänen, J., & Malik, K. (2023). How can we build inclusive circular supply chains? Examining the case of agricultural residue usage in India. Business Strategy & Development.

20. Silvério, A. C., Ferreira, J., Fernandes, P. O., & Dabić, M. (2023). How does circular economy work in industry? Strategies, opportunities, and trends in scholarly literature. Journal of Cleaner Production, 137312.

21. Zhao, J., Hu, E., Han, M., Jiang, K., & Shan, H. (2023). That honey, my arsenic: The influence of advanced technologies on service employees' organizational deviance. Journal of Retailing and Consumer Services, 75, 103490.

22. Riso, V., Tallaki, M., Bracci, E., & Cantele, S. (2023). The transition towards benefit corporations: What are the roles for stakeholders?. Business Strategy and the Environment.

23. Neier, T. (2023). The green divide: A spatial analysis of segregation-based environmental inequality in Vienna. Ecological Economics, 213, 107949.

24. Frehywot, S., & Vovides, Y. (2023). An equitable and sustainable community of practice framework to address the use of artificial intelligence for global health workforce training. Human Resources for Health, 21(1), 45.

25. Clement, J., Ruysschaert, B., & Crutzen, N. (2023). Smart city strategies–A driver for the localization of the sustainable development goals?. Ecological Economics, 213, 107941.

26. Dinçer, H., Yüksel, S., Hacioglu, U., Yilmaz, M. K., & Delen, D. (2023). Development of a sustainable corporate social responsibility index for performance evaluation of the energy industry: A hybrid decision-making methodology. Resources Policy, 85, 103940.

27. Lamrini, L., Abounaima, M. C., & Talibi Alaoui, M. (2023). New distributed-topsis approach for multi-criteria decision-making problems in a big data context. Journal of Big Data, 10(1), 1-21.

28. Canessa, C., Venus, T. E., Wiesmeier, M., Mennig, P., & Sauer, J. (2023). Incentives, rewards or both in payments for ecosystem services: Drawing a link between farmers' preferences and biodiversity levels. Ecological Economics, 213, 107954.

29. Sanasi, S. (2023). Entrepreneurial experimentation in business model dynamics: Current understanding and future opportunities. International Entrepreneurship and Management Journal, 19(2), 805-836.