RESEARCH ARTICLE

# Network science and explainable AI-based life cycle management of sustainability models

Ádám Ipkovich[1☉], Tímea Czvetkó[1☉], Lilibeth A. Acosta[2☉], Sanga Lee[2☉], Innocent Nzimenyera[2☉], Viktor Sebestyén[1,3☉], János Abonyi[1☉]*

1 HUN-REN-PE Complex Systems Monitoring Research Group, University of Pannonia, Veszprém, Hungary, 2 Climate Action and Inclusive Development (CAID) Unit, Global Green Growth Institute, Jung-gu, Seoul, Republic of Korea, 3 Sustainability Solutions Research Lab, Faculty of Engineering, University of Pannonia, Veszprém, Hungary

☉ These authors contributed equally to this work.
* janos@abonyilab.com

## Abstract

Model-based assessment of the potential impacts of variables on the Sustainable Development Goals (SDGs) can bring great additional information about possible policy intervention points. In the context of sustainability planning, machine learning techniques can provide data-driven solutions throughout the modeling life cycle. In a changing environment, existing models must be continuously reviewed and developed for effective decision support. Thus, we propose to use the Machine Learning Operations (MLOps) life cycle framework. A novel approach for model identification and development is introduced, which involves utilizing the Shapley value to determine the individual direct and indirect contributions of each variable towards the output, as well as network analysis to identify key drivers and support the identification and validation of possible policy intervention points. The applicability of the methods is demonstrated through a case study of the Hungarian water model developed by the Global Green Growth Institute. Based on the model exploration of the case of water efficiency and water stress (in the examined period for the SDG 6.4.1 & 6.4.2) SDG indicators, water reuse and water circularity offer a more effective intervention option than pricing and the use of internal or external renewable water resources.

## Introduction

To achieve Sustainable Development Goals (SDGs) and their targets, regular assessments are needed to track the progress of countries and identify areas where more effort is necessary [1]. These assessments may require a range of data and information, including indicators on social, economic, and environmental issues [2]. Evaluations can also involve consultation with various stakeholders, including governments, civil society organizations, and the private sector [3]. The main objective is to provide a complete picture of progress toward the SDGs, identify challenges and opportunities [4], and inform policies and actions that can help accelerate progress [1].

The complex relationships between policies and social, economic, and environmental issues require models [5] and analyses that consider each issues and predict the performance of sustainable development in countries [6]. Model-based assessment of sustainability planning provides decision-makers with a powerful tool to understand the complex interdependencies between social, economic, and environmental systems and design policies and strategies that are likely to lead to more sustainable outcomes. It involves developing quantitative structural equation models that represent the systems that are being considered and using these models to explore different scenarios and assess the likely impacts of different policy choices. In assessing aggregated SDG indicators, effective decision support requires models that are suitable for identifying potential intervention points and for establishing knowledge in political strategy creation.

As SDG indicators are aggregated, effective decision support requires models that are suitable to identify potential intervention points and establish knowledge in the creation of political strategies. Country-specific factors such as unique development phases, policies, and databases can make it difficult to develop accurate and reliable system dynamics models. Therefore, it is important to pay attention to the life cycle of models and the difficulties that may arise during each stage, such as initial contextualization, data utilization, model creation, analysis, implementation, and monitoring.

A common challenge in developing reliable complex systems is model identification to accurately represent the workings of the system. Regarding sustainability, this challenge can arise due to the detailed structure of SDGs, data and information limitations, and modeling techniques. The development of structured models requires broad experience in modeling techniques and validation methods and relies heavily on data and information. An additional requirement is that the relationships between variables, and their change in behavior over time must be known [7]. In this regard, life cycle-based assessment encompasses the various stages of model development, including initial contextualization and data utilization for design, the creation of a complex model, the derivation of analysis, the implementation of the model, and the monitoring and revision of the model to improve its accuracy.

The primary objective of this paper is to highlight approaches that support structural equation model-based assessment of sustainability planning, identify the contribution of variables to the SDG indicators based on historical data, and promote evidence-based policy development. We propose methods based on network and data analysis to support modeler work and obtain more accurate, automatic, and efficient model development procedures during the entire life cycle of models. We introduce the potential for Shapley value utilization to identify the contribution of each variable in the model to the output SDG indicators. Furthermore, we highlight the opportunities for life cycle-based modeling as well as the links between life cycle phases with the help of network and data science technologies. It is important to highlight that we do not mean the identification of environmental impacts by the expression 'life cycle', we consider the modeling process as a life cycle, which draws attention to the fact that during the development of expert systems (in their life cycle) different steps can be followed to develop the concept, for which the tools proposed in this research can be used.

Structural equation models can be represented as networks, and sustainable development-related problems can be evaluated using network science tools [8]. Network analysis can be useful in understanding the complex networks of stakeholders, institutions, and processes that influence sustainability outcomes. Furthermore, it can be used to understand the relationships between different SDG goals and targets. For example, network tools can help identify the interdependencies between different goals and understand how progress in one area can impact others [9]. The important variables of the model can be identified that can serve as possible intervention points. Additionally, data-based methods can be used to support the

continuous development and analysis of structural equation models to assess progress toward the SDGs. Data-driven methods can reveal the relationship between different variables, support understanding of the underlying mechanisms that drive these relationships, identify patterns and trends, and the potential impact of different interventions on the system [10].

The related models are usually static in a way that the developed models are not further improved, validated, and maintained throughout the modeling life-cycle. The modeling structure is usually not prepared to integrate and track changes, whereas we are and will be witnessing more and more dramatic changes because of the radical steps taken in the direction of polycentrism and green transformation. Therefore, there is an emerging need for the model life-cycle-based development.

By combining network and data analysis throughout the entire model life cycle, stakeholders can gain a more complete understanding of the complex interactions and relationships that influence the achievement of the SDGs. This life cycle-based assessment of models can be used to ensure that models are developed, deployed, and maintained in a reliable, efficient, and effective manner. In this study, we consider the intertwining of MLOps and CRISP-ML (Q) (Cross-Industry Standard Process for the development of Machine Learning applications with Quality assurance methodology) life cycles. MLOps is an end-to-end framework of the machine learning (ML) development process and supports automation principles and includes three phases: design, model development, and operation [11]. This framework is highly intertwined with CRISP-ML(Q), which consists of six iterative phases with varying sequences: business and data understanding, data preparation, modeling, evaluation, deployment and monitoring, and maintenance [12].

The system model of the Sustainable Development Goals can take into account changes over time in the relationships between the SDG goals and the targets [13], so there is a significant need to maintain the models, as the behaviour of the system is dependent on time and maturity with a special focus on responses to political challenges [14], including critical and success factors in economic sectors and national infrastructures [15]. Action-focused approaches increase their utility in decision-making, so the scenario derivation supports the formulations of reasonable actions [16].

Additional variables could be added to the explanation of 'unexplained' interactions in the SDG models [13], which can be supported by machine learning [15]. Modeling efforts can be exploited with a feedback-rich structure of models [17].

Therefore, the potential contributions and the structure of this article are as follows:

- The life cycle-based model development and possible application of data and network science tools in each life cycle phase are discussed in the *Sustainability focused model life cycle management section.*

- *The methodology of using network and data science tools in model identification and development is explored in the Development of models based on network science and explainable AI tools section. This section emphasizes the challenges inherent in model development and how the utilization of network analysis and the Shapley value can mitigate these difficulties. They serve as sensitivity analyses by promoting the understanding of which variables have the greatest direct or indirect impact on changes in SDG indicators through the model network based on historical data.*

- *A case study is presented in the Application of network and data analysis for model development section to demonstrate the effectiveness of the proposed methods. The case study is conducted on a submodel developed by the Global Green Growth Institution for Hungary.*

It is imperative to underscore that the evaluation of models is contingent upon the data they reflect, meaning that the results are primarily influenced by the data. Consequently, insufficient quality of the data can lead to incorrect inferences about the model's quality and the effects of variables. Furthermore, we must consider that different economies or environmental backgrounds may exhibit different relationships. Therefore, a homogeneous sample is essential, requiring a clear definition of the modeled area. Naturally, the structure of models can also be flawed, and identifying structural errors may not be straightforward. MLOps primarily focuses on retraining the same model to address this issue. However, the methodology enables to creation of a new model by incorporating new variables and insights.

Nevertheless, we believe that the network and data analysis approach to model-based sustainability planning offers valuable information throughout the entire model life cycle. Furthermore, by identifying the direct and indirect contributions of variables to changes in the SDGs, this data-driven approach can serve as a fundamental basis for evaluating the efficacy of policy intervention. The relevancy of the proposed method is proved and validated through the Global Green Growth Institute (GGGI) who applied the method in real-life sustainability models. The interpreted application study relies on a real-life example that is developed by the Global Green Growth Institute, which is a treaty-based international, inter-governmental organization dedicated to supporting and promoting strong, inclusive, and sustainable economic growth in developing countries and emerging economies. The developed method will be an integrated tool that will be used for their work, enabling policymakers and researchers to assess the environmental and economic impacts of various policy interventions, and fostering informed decision-making for sustainable development initiatives worldwide.

## Materials and methods

### Sustainability focused model life cycle management

Environmental systems exhibit a high degree of complexity and dynamics, characterized by diverse interacting variables that may vary over time. Models that aim to capture such complexities are often susceptible to uncertainties attributable to the evolving nature of environmental systems, the lack of data, and the adaptability of the model to local conditions. The localization of model structures necessitates the customization of generalized models to reflect unique characteristics of countries, regions, or cities, including but not limited to social, economic, and environmental developmental levels, as well as diverse policy goals and interventions aimed at enhancing sustainable development.

The systematic development of these models can be supported with machine learning techniques to incorporate data-driven insights and evidence-based policy assessment. MLOps and CRISP-ML are systematic frameworks that consider the whole life cycle of data-driven models from design and business understanding until the operation phase and model maintenance [18]. This life cycle-based model evaluation can be used to ensure that models are developed, deployed, and maintained in a reliable, efficient, and effective manner. These concepts are fundamentally iterative and exploratory, so depending on the results from the later phases, the reexamination of earlier steps may be needed.

In Fig 1, the framework for applying data and network analysis-based life cycle management for the development of sustainability models is shown. The blue block represents model development, where model building is based on expert knowledge and time series data and requires the use of modeling and analyzing techniques, as well as country-specific targets and policy implications for accuracy. The Global Green Growth Institute built models for assessing and predicting green growth by linking the energy, agriculture, forestry and other land use (AFOLU), and water and waste models [19]. The orange block indicates the difficulties of
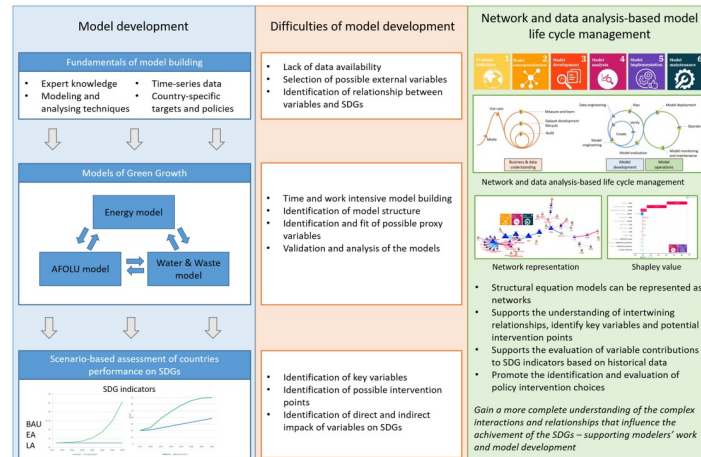
**Fig 1. The framework of applying network and data analysis-based life cycle management of sustainability model development.**

model development as mentioned above, such as lack of data availability, identification of possible external variables and their relationship, model validation, assessment of intervention points, analysis of the effect of variables on SDGs, *etc*. The green block represents a model life cycle-based solution for model development utilizing the tools of network and data sciences.

By following the CRISP-ML methodology, organizations can ensure that models are developed, evaluated, and deployed in a transparent and structured manner, which can help to improve decision making and support long-term planning and management of environmental systems. Additionally, the CRISP-ML provides a clear roadmap for monitoring and maintaining the model, allowing continuous improvement and adaptation to new information and scenarios.

The schematic workflow and connections between MLOps and CRISP-ML are illustrated in Fig 2, where the line colors represent the different stages (orange—business and data understanding, blue—model development, green—model operation) and the arrow colors identify where network and data science tools can be applied (network science—grey, data science—yellow).

In Table 1, the intersection of MLOps, CRISP-ML life cycle, and the potential applicability and implementation of network and data science techniques are examined within the framework of model-based assessment of Sustainable Development Goals.

According to Studer *et al*. [12] the *business and data understanding phase* includes defining objectives, collecting and verifying data quality, and assessing the project. This phase involves
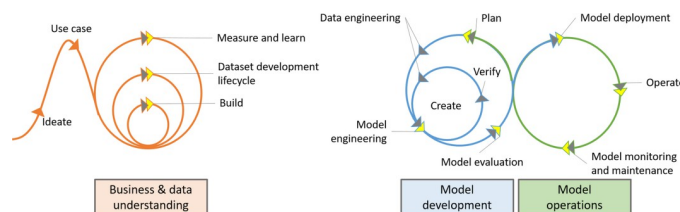


**Fig 2. Schematic representation of the connections between MLOps and CRISP-ML (based on [18]).** The line colors represent the different stages (orange—business and data understanding, blue—model development, green—model operation) and the arrow colors identify where network and data science tools can be applied (network science—grey, data science—yellow).

**Table 1. Data and network science methods supporting sustainability model development in the MLOps and standard modeling life cycle.**

| MLOps | CRISP-ML dimensions | Model relevance | Data science relevance | Data science methods | Network science relevance | SDG relevance | References |
|---|---|---|---|---|---|---|---|
| Design | Business understanding | S1 Fig Problem definition | Understand the problem, goals and objectives of the project, by analyzing historical data, identifying trends, and patterns and understanding the relationships between different variables. | Exploratory Data Analysis (EDA) [20], Descriptive statistics, Data mining [21], Data visualization, Causality analysis [24] | Identify the key variables and relationships that are important for understanding the behavior of the system, by creating a conceptual model of the system. | Identify the interlinkages between SDG goals [25], targets [36] and indicators [9]. Identify the connection between model elements and SDGs. Identify dependencies among risks associated with SDGs [26]. | [9, 20, 21, 23–27, 36–42] |
| | Data understanding | S2 Fig Model conceptualization | Collecting and verifying the data quality regarding the structured issues and, finally deciding upon whether the project should be committed. | Data cleaning, Data exploration [20], Correlation analysis [27], Causality analysis [23], Data visualization, Shapley value [42] | | | |
| Model development | Data preparation | S3 Fig Model development | Representing sustainability-related concepts and structured knowledge of the relationship between model elements. Producing a data set for the modeling design phase by gathering and linking data from diverse data sources. | Feature engineering, Data imputation, Feature scaling, Data normalization, Data selection, Data cleaning, Ontology modeling, Knowledge graphs, Open linked data | Defining knowledge graphs representing the relationship between model elements. Gather and analyze the data that is needed to validate the model, by identifying patterns and structures in the data that can be used to inform the development of the model. | Represent the SDGs and their interrelationships with other ontologies to facilitate the integration and analysis of data applying to the SDGs. Provides transparent and reproducible analysis, establishing a shared vocabulary and framework for the description and sharing of data and models associated with the SDGs. | [43–48] |
| | Modeling | | Creating models that satisfy the given constraints and requirements. | Regression analysis, Time series analysis [28], Clustering, Decision trees, Random forest, Neural networks [32], Bayesian methods [49], and other machine learning algorithms [29], Digital Twins [35], Composite indicators Cross-validation, Hyperparameter tuning, Sensitivity analysis, Classification, Clustering, Regression analysis, Frequent itemset/pattern mining, Visualization | Build and test the model, by using network metrics such as centrality to identify key elements in the system and to understand the overall connectivity of the system. | Predictions [50], simulations, structural model development, risk assessment [26] | [26, 28–35, 49–52] |
| | Evaluation | S4 Fig Model analysis | Design and analysis are intertwined and as new change options are proposed, they are analyzed using data analysis techniques. The performance, robustness and explainability must be evaluated. | | Analyze the results of the simulation and to identify the best course of action for achieving a desired outcome, by identifying patterns and structures in the results that can be used to inform decision-making. | Identify and categorize countries, regions based on their performance on SDGs. Identify patterns and key variables which can be possible intervention points to succeed on SDGs. | [5, 25, 29, 38, 40, 52] |

(*Continued*)

**Table 1.** (Continued)

| MLOps | CRISP-ML dimensions | Model relevance | Data science relevance | Data science methods | Network science relevance | SDG relevance | References |
|---|---|---|---|---|---|---|---|
| Operations | Deployment | S5 Fig Model implementation | Before rolling out a model to all, it is best practice to deploy it first to a small subset and evaluate its behaviour. | Model deployment, Simulation models, Digital Twins [35], Model interpretability, Monitoring [53], Web page development, Dashboard development, Model maintenance | Identify the most critical elements in the network that need to be targeted for the implementation of policies or strategies that have been identified as the best course of action. | Identify transition trajectories, different scenarios based on diverse strategies. Identify the possible performance on SDGs based on the diverse strategies. | [5, 30, 35, 51, 52] |
| | Monitoring and Maintenance | S6 Fig Model maintenance | Monitoring and maintenance processes to assure quality performance and identify corrective actions. The model has to adapt to changes in the environment. | | | Continuous monitoring [54], improvement and validation of the models. Monitor the transition trajectories. Identify possible external variables and mitigation actions for better accuracy. | [53–55] |

identifying the project's scope, risks, and success criteria such as measurable features, systems alignment, Key Performance Indicators (KPIs), feasibility, and data availability, which are critical components of model development. Regarding data science methods, exploratory data analysis (EDA) provides insight and understanding of databases, visualization of potential relationships between variables, detection of outliers and anomalies, development of simple models (predictive or exploratory), the precondition of data [20]. Data mining and machine learning methods can be valuable tools for discovering knowledge from databases [21]. Integration of data mining and system dynamics supports evidence-based decision-making and a better understanding of the dynamics and complexity of a system [22]. The relationships between, *e.g.*, sustainability pillars [23], the interconnectedness of SDGs [24], and their patterns can be identified through causality analysis that is important from a policy point of view. The causal relationships [24] and complex interactions of the SDGs can be mapped through network analysis [25]. For example, a probabilistic network model can be used to explore the dependencies among the risks associated with the SDGs [26].

During the phase of model development, the data is preprocessed and delivered in a format that is suitable for integration with the model. This process, also known as data preparation, includes the selection, cleaning, and standardization of the data. The modeling activity involves selecting appropriate models, incorporating domain knowledge, evaluating, validating, and documenting the models. The data preparation stage is supported by techniques such as feature engineering, scaling, normalization, and data imputation [27]. There are various methods that can be applied to the building and analysis of models, such as time series analysis and forecasting [28], regression analysis, clustering, classification, dimension reduction techniques [29], Monte Carlo simulation [30], system dynamics modeling [31], network analysis [25], neural networks [32], Bayesian networks [33], and composite indicators [34]. These methods can be integrated to provide a more comprehensive understanding of the system under study. It is worth noting that, while these methods provide valuable information, the identification of other factors such as dynamics, uncertainties, and external events contribute to a more complete understanding of the system.

This will be the basis for *model operations*, where the models are continuously deployed, evaluated, and monitored, as well as fine-tuned if necessary. It is considered good practice to

first deploy a model to a small subset and evaluate its behavior prior to extensive implementation *e.g.* transition trajectories, and digital twin simulation based on scenarios may help improve understanding of a phenomenon [35].

The literature on utilizing a network and data-driven tools for the assessment of SDGs has been growing, but their integration and contribution throughout the entire model life cycle remains an area of underexplored research, even though they are used in the steps. Developing a framework that systematically evaluates the potential and contribution of network and data science techniques throughout the model life cycle could prove to be a valuable asset in supporting the work of modelers. In this way, the integration of data-driven techniques and network-based analysis of the system can improve the accuracy, robustness, interpretability, and efficiency of environmental system dynamics models, by allowing the selection of the most relevant variables to the problem, and by identifying key drivers and relationships, it helps organizations and decision-makers to better understand the dynamics of environmental systems and make more informed decisions.

## Development of models based on network science and explainable AI tools

The MLOps approach introduces much-needed structure and efficiency to the development of machine learning and network science models. By combining MLOps practices with explainable AI tools, organizations can leverage the power of these models across different domains, ensuring transparency and interpretability, which in turn fosters greater trust in AI-driven decision-making. In this context, we present a novel Shapley value-based method in Section, which complements the MLOps life cycle by evaluating the contributions of variables to the model. This approach promotes a better understanding of the model, as well as the data, and facilitates model conceptualization, development, and analysis. Furthermore, in Section we highlight the significance of combining the Shapley value and network science in identifying key drivers of model behavior.

**Shapley value-based evaluation of variable contribution.** Complex models are often represented in the form of structural equation models (SEMs), whose equations describe the relationship between variables [56]. These structured equation models can be represented and analyzed as networks. Combining the measures of the structural network with the Shapley value [57] has a great potential to achieve a more complete understanding of the system. As the Shapley value defines the contribution of the variables, in this regard, these values can be applied as weights of the network edges. Thus, a deeper understanding of the interactions of the models can be achieved, which can lead to valuable insights for model development and analysis. The combination of Shapley value and networks is emerging, and has been utilized for example to discover influential nodes in social network [58], identify centrality in weighted and unweighted networks [59], identify individuals' performance in group influence within a real-world social network [60], used as an extension to betweenness centrality and define a new metric called stress centrality [61]. Furthermore, the Shapley value-based interpretation of feature contribution to model predictions have been applied in a variety of studies, such as predicting antifungal peptides [62], anti-tubercular peptides [63], or anti-inflammatory peptides [64].

Consider a set of $V$ that contains $n$ number of variables $[\hat{x}_1^{(t)}, \ldots, \hat{x}_l^{(t)}, x_{l+1}^{(t)}, \ldots, x_n^{(t)}]$ with $l$ number of observed input variables ($\hat{x}$) and $n - l$ number of derived variables. The *kth* structural model predicts the *kth* derived variable $x_k^{(t)}$ at a given time $t = 1, \ldots, T$. The model $f_k$ requires the set of inputs ($\chi_k^{(t)} \subseteq V$) and parameters $\theta_k$ to predict a variable:

$$x_k^{(t)} = f_k(\chi_k^{(t)}, \theta_k) \tag{1}$$

where $x_k^{(t)}$ denotes the predicted variable, $f_k$ defines the structural equation, $\chi_k^{(t)}$ is considered as the set of input variables, while $\theta_k$ stands for the parameters of the equation.

A structural equation model is built from several equations; the model structure can be hierarchical, the output of an equation may be realized as an input of another function.

$$\chi_k^{(t)} = \left\{ x_j^{(t)} \middle| a_{j,k} \neq 0 \right\}, j = 1, \ldots, n; t = 1, \ldots, T; x_j^{(t)} \in V \tag{2}$$

where $\chi_k^{(t)}$ is the set of variables at time $t$ for function $f_k$. The $jth$ variable is denoted by $x_j^{(t)}$. Set $V$ contains all variables, $a_{j,k}$ stands for the one-way connection between the $jth$ input and the $kth$ output variable.

Moreover, data may differ by context, so the models must be parameterized properly. In some cases, parameters may require re-identification as time goes on; new technological and cultural changes may trigger a need for adjustment in the model behavior. One possible solution is to minimize a cost function, *e.g.* the mean squared error of the prediction of the model to an observed output.

$$\min_{\theta_k} \frac{1}{T} \sum_{t=1}^{T} \left( x_k^{(t)} - f_k \left( \chi_k^{(t)}, \theta_k \right) \right)^2 \tag{3}$$

The inclusion of extraneous variables is a potential issue that may arise during the process of developing a model. Therefore, it is important to determine which variable may be relevant to the dependent variable, as it may result in a less convoluted and more accurate model. Establishing the contribution of a variable to the model may enable us to filter out extraneous ones, while also determining the highest contributor, known as biases. However, the impact of one variable should include all contributions to the model, including the added value of cooperation with a group. As such, the average marginal contribution is required, which is often determined by the Shapley value [57].

The marginal contribution requires marginalizing over ("averaging out") the variables that are excluded from the evaluated selected group of variables. Consider $x_k^{(t)}$ as the dependent variable, $\chi_k^{(t)}$ the set of independent variables. If marginal contributions of a group variable are required, then a contribution evaluation function $\phi(\cdot)$ is required for a subset. The excluded variables are marginalized over (averaged out), and, therefore, the contribution of the selected variables is returned. The contribution should be understood as a difference from an expected value (total average of the predictions) [65]:

$$\phi \left( \chi_k^{(t)} / \{ x_j^{(t)} \} \right) = \int f_k \left( \chi_k^{(t)}, \theta_k \right) dx_j \quad - E[f_k]; \tag{4}$$

where $\phi(\chi_k^{(t)} / x_j^{(t)})$ defines the marginal contribution of variable $x_j^{(t)}$. Here, the expected value is approximated as the average of the model predictions: $E[f_k] \approx \frac{1}{T} \sum_{t=1}^{T} f_k \left( \chi_k^{(t)}, \theta_k \right)$.

The classical way to approach the Shapley value (calculation of the average marginal contribution) is by subtracting the contribution of a subset of variables with and without the selected variable ($x_j^{(t)}$), and taking their weighted average [57], summed over all possible subsets:

$$S_{j,k}^{(t)} = \frac{1}{|\chi_k^{(t)}|!} \sum_{P \subseteq \chi_k^{(t)} / \left\{ x_j^{(t)} \right\}} \phi \left( P \cup \{ x_j^{(t)} \} \right) - \phi(P) \tag{5}$$

where $S_{j,k}$ denotes the Shapley value, $P$ denotes a subset of variables $\chi_k^{(t)}$. $\phi(P)$ denotes the

marginal contribution of the set without, $\phi(P \cup x_j^{(t)})$ defines the marginal contribution of the subset with the examined variable.

As the Shapley value is computationally intensive, it is often approximated by Monte Carlo simulation. For a set with $n$ variables, if all possible subset orders are considered, the operation time would be factorial ($O(n)!$). As such, variable subsets are randomly sampled with a $M$ number of permutations rather than calculating the contribution of each subset. The calculation of the Monte Carlo-based Shapley value requires a set for averaging out the combination of the variables (*e.g.* substituting the expected value of the variable to the set), and the other is for the original values. As such, the computational time can be reduced to $O(nM)$. The formalization of the approximation of the contribution of the variable can be considered as [66]:

$$
\begin{aligned}
\hat{\phi}_j(\hat{P}_m) &= f_k\left(\chi_{k\ [xi=E(x_i^{(t)})]}^{(t)}, \theta_k\right) - f_k\left(\chi_{k\ [x_i=E(x_i^{(t)}),x_j^{(t)}=E(x_j^{(t)})]}^{(t)}, \theta_k\right), \\
i, j &\notin m, i \neq j \\
\bar{S}_{j,k}^{(t)} &= \frac{1}{|M|}\sum_{m \in M}\hat{\phi}_j(\hat{P}_m)
\end{aligned}
\tag{6}
$$

where $M$ is a population of a randomly sampled set of feature combinations, $m$ is a feature combination in $M$ with variables that are not marginalized, $\bar{S}_{j,k}^{(t)}$ denotes the approximation of the Shapley value $S_{j,k}^{(t)}$. $[x_i = E(x_i^{(t)})]$ suffix defines the set of variables, where the specific variables (here, the *ith*) is fixed to an expected value.

The Shapley value aims to explain the change in the value function compared to the expected value, therefore, the sum of the Shapley values for a model must return the difference between the function value at point $t$ and the expected value:

$$
\sum \bar{S}_{j,k}^{(t)} = f_k(\chi_k^{(t)}, \theta_k) - E[f_k]
\tag{7}
$$

Eq 7 defines the efficiency property of the Shapley value, which is to be used to scale individual values between [0, 1]. The marginal contribution of a variable can be defined for an individual sample (*e.g.* contributions of variables for year 2000), or for each to provide information of the trends over the year. For the *jth* input of the *kth* model, the individual and mean Shapley is as follows:

$$
w_{j,k} = \begin{cases} \dfrac{\bar{S}_{j,k}^{(t)}}{\sum_{i \in \chi_k^{(t)}}|\bar{S}_{i,k}^{(t)}|} & \text{if } d_{ind} = 1 \\[4ex] \dfrac{\frac{1}{T}\sum_{t=1}^{T}\bar{S}_{j,k}^{(t)}}{\sum_{i \in \chi_k^{(t)}}\frac{1}{T}\sum_{t=1}^{T}|\bar{S}_{i,k}^{(t)}|} & \text{if } d_{ind} = 0 \end{cases}
\tag{8}
$$

where $d_{ind}$ defines a dummy variable as to whether the weight should be the individual or the average Shapley value. The first member defines the normalized individual contribution, while the second denotes the normalized mean contribution. For a proper weight matrix (**W**), all $d_{ind}$ should take the same value. If the weight type is selected, then for each model, a vector can be defined where relevant inputs to the models have contributions, *i.e.* the Shapley value is designated as a zero for diagonals (so that there is no cycle) and observed inputs, who have no models.

As shown above, the Shapley value provides a perfect measure to establish the direct connection between variables, making it particularly valuable in the context of MLOps. In the development and deployment of hierarchical models within the MLOps framework, the

Shapley value can be used to assess the individual contributions of variables, aiding in the interpretation, monitoring, and optimization of these models. In the following subsection, the models will be defined as a directed acyclic graph of the variables that uses the Shapley value as weight. By incorporating the Shapley value-based method and network analysis into the MLOps life cycle, organizations can enhance their understanding of model behavior, promote transparency, and optimize the performance of hierarchical machine learning models.

**Network analysis for identifying key drivers of the model.** Network analysis is a technique that can be used in MLOps to analyze the relationships between different variables or features in a machine learning model. The application of specific network analysis tools may help select key variables and intervention points in the graph that affect the model the most. If the contribution of a variable is added as weight to the network, it may alter the result of the analysis. By building a network, one can a) employ network science tools to understand the relationship of variables, c) determine the key nodes of the model, and b) select intervention points that significantly influence the output of the model.

The Shapley-weighted degree centrality may support the identification of the central position of a variable for their contribution to change *i.e.* has various connections throughout the model, and can be manipulated by several inputs, or can change the value of several outputs [67]. It is a simple measure of the number of neighbors of a node; however, it provides information on the degree one node influences or being influenced by other nodes. Degree centrality is defined as the number of edges incident on a node, and also considering the weight and directions:

$$C_d(k) = \frac{\sum_{j=1}^{n}(w_{j,k} + w_{k,j})}{2(n-1)}, k \neq j \tag{9}$$

where $C_d$ defines the degree centrality measure for directed graphs, and $w_{j,k}$ defines the weight of the edges (*e.g.* individual Shapley values) for the directed edge between node $j$ and $k$, respectively.

Closeness centrality focuses on the average shortest path connections to other (significant) nodes, so with consideration, nodes with high closeness may be able to manipulate various key nodes at once [68]. It is defined as the reciprocal of the sum of the shortest path distances between the node and all other nodes in the network. For a node $x_k$, the formula is as follows:

$$C_c(k) = \frac{1}{\sum_{j \neq k} d(x_j, x_k)} \tag{10}$$

where $C_c(k)$ defines the closeness centrality, $x_k$ is the node in question, $x_j$ is a node in the network, $d(x_j, x_k)$ is the shortest path distance between node $x_k$ and node $x_j$.

Betweenness centrality may improve the interpretation of the structure of the variables, as it measures the amount of shortest paths going through the node [69]. The shortest path may be influenced by the weights, so the role of this centrality is to find the measure of how dominant the role of the node is.

$$C_b(k) = \sum_{x_j \neq x_i \neq x_k} \frac{\delta_{x_j,x_i}(k)}{\delta_{x_j,x_i}} \tag{11}$$

where $C_b(k)$ defines the betweenness centrality, $x_j$ and $x_i$ denote nodes in the network, $\delta_{x_j,x_i}(k)$ stands for the number of shortest paths from node $x_j$ to node $x_i$ that pass through node $x_k$, and $\delta_{x_j,x_i}$ is the total number of shortest paths from node $x_j$ to node $x_i$.

Other network analysis tools such as community detection (*e.g.* Louvian [70]) and shortest paths algorithms (*e.g.* Dijkstra's algorithm [71]) may provide additional information on the relationship of the variables and significance of intervention points.

Network tools and the incorporation of the Shapley value allow us to identify the most influential variables, supporting the work of data scientists and modelers to focus on optimizing and fine-tuning the key drivers for improved model performance. With respect to the MLOps framework, network analysis offers a more holistic and systematic approach to understanding the underlying dynamics of complex machine learning models.

## Results

### Application of network and data analysis for model development

This section presents the processes for improving a system of environmental simulation based on structural equation models through a Hungarian case study provided by the Global Green Growth Institute [19]. The Green Growth Model of Hungary consists of three connected structural equation submodels, namely: energy, agriculture, forestry and other land use (AFOLU), water and waste models. To demonstrate of the efficiency of the data-driven Shapley value and network analysis for hierarchical models, we apply the methods on a submodel of the Hungarian water model. The introduced GGGI models rely on literature-based structural information and the model connections are validated with subject-matter experts. The model itself is deployed and maintained as a Python and Dash-based application. The practical and accurate development, monitoring, and maintenance of the ecosystem of the environmental models requires the principles of MLOps to apply.

As Hungary has progressed the least in SDG6 [5] since 2015, the interpretation of the water-related submodel (SDG6 indicators) could be of great help in explaining the required planning aspects of effective political interventions. Due to the fact that the modeling of the SDG indicators is a complex problem, the time series data (2000–2019) were compiled from sources such as FAO [72–76], NASA [77, 78], World Bank [79], United Nations [80], WHO & UNICEF [81], and the Hungary Ministry of Innovation and Technology [82] by GGGI. Some variables were imputed with the 2017 values due to unknown or missing data. The abbreviations and short description of the variables handled in the water model are shown in Table 2.

The structural equations of the water model is indicated in Table 3. The function indices are named after the variable indices. It is also indicated which final output the derived variable contributes to, as well as the equation of each function is specified as well as the literature-based evidence of the model relationships are added for each model equation.

In Fig 3, the schematic representation of the Shapley value-based interpretation of the structural equation model and the contributions of the variables is shown. The structured equation and the Shapley value method are formulated as follows.

In the following sections, we focus on the validation of the hierarchical models to show that such models provide correct information and knowledge to transfer to decision-makers. As the MLOps principles focus not only on the development but also on validation and maintenance, they contain the necessary toolset to confirm the possible use of the models. Therefore, first, we validate the model by comparing it with linear regression and k-th nearest neighbors models to ensure the validity of the connection between the independent and the dependent variables. Then the roles of the variables are discussed with the help of Shapley-value-based analysis. Lastly, we introduce network analysis to select the most important input.

**MLOps life cycle in environmental model development.** The three stages of the MLOps life cycle include the business and data understanding, model development, and model operation phases. Continuous revision, development, and evaluation of environmental models are

**Table 2. Variables for the water efficiency model.** The inputs and outputs contain data from 2000–2019. It is important to note that various inputs are imputed. (PI) denotes policy intervention points. Parameter (I) represents identifiable parameters. Variables without a specific unit are symbolized with a one in the Unit column of the table.

| Notation | Abbreviation | Name | Unit | Type |
|---|---|---|---|---|
| $\hat{x}_1$ | AGVA | Agricultural Gross Value Added | $ | Input |
| $\hat{x}_2$ | AIR | Agriculture area actually irrigated | 1000 ha | Input |
| $\hat{x}_3$ | Arice | Area of Rice Paddy Irrigation | 1000 ha | Input (PI) |
| $\hat{x}_4$ | CL | Cropland | 1000 ha | Input |
| $\hat{x}_5$ | DW | Desalination Water | $m^3$/year | Input |
| $\hat{x}_6$ | ERWR | External Renewable Water Resources | $m^3$/year | Input |
| $\hat{x}_7$ | ETa | Actual Evapotranspiration | mm/year | Input |
| $\hat{x}_8$ | ETo | Evapotranspiration | mm/year | Input |
| $\hat{x}_9$ | GDPC | GDP per capita | $ | Input |
| $\hat{x}_{10}$ | IGVA | Industrial Gross Value Added | $ | Input |
| $\hat{\mathbf{x}}_{11}$ | IRRTECHi | Irrigation technology proportion | 1 | Input (PI) |
| $\hat{x}_{12}$ | IRWR | Internal Renewable Water Resources | $m^3$/year | Input |
| $\hat{x}_{13}$ | IWU | Industrial Water Withdrawal | $10^9 m^3$/year | Input |
| $\hat{x}_{14}$ | Pop | Population | capita | Input |
| $\hat{x}_{15}$ | SGVA | Service Sector Gross Value Added Resources | $ | Input |
| $\hat{x}_{16}$ | TW | Treated Wastewater | $/15$m^3$ | Input |
| $\hat{x}_{17}$ | WP | Water Price | $/15$m^3$ | Input |
| $\mathbf{x}_{18}$ | AIRi | Irrigated area per irrigation technology type | 1000 ha | Variable |
| $x_{19}$ | AWU | Agricultural Water Withdrawal | $10^9 m^3$/year | Variable |
| $\mathbf{x}_{20}$ | CI | Cropping Intensity | 1 | Variable |
| $x_{21}$ | Cr | Corrective coefficient | 1 | Variable |
| $x_{22}$ | ETc | Potential Crop Evaporation Vector | mm/year | Variable |
| $x_{23}$ | ICU | Irrigation Consumptive Use | mm/year | Variable |
| $x_{24}$ | IWR | Irrigation Water Requirement | $10^9 m^3$/year | Variable |
| $\mathbf{x}_{25}$ | IWRi | Irrigation Water Requirement per irrigation | 1e9 $m^3$/year | Variable |
| $x_{26}$ | IWW | Irrigation Water Withdrawal | $10^9 m^3$/year | Variable |
| $x_{27}$ | MWU | Municipal Water Withdrawal | $10^9 m^3$/year | Variable |
| $x_{28}$ | PAIR | Proportion of Irrigated Cropland | 1 | Variable |
| $x_{29}$ | TFA | Total Freshwater Available | $m^3$/year | Variable |
| $x_{30}$ | TNCW | Total Non Conventional Water | $m^3$/year | Variable |
| $x_{31}$ | TRF | Total Renewable Freshwater | $m^3$/year | Variable |
| $x_{32}$ | TWW | Total Water Withdrawal | $10^9 m^3$/year | Variable |
| $y_1$ | SDG 6.4.1 | Water Use Efficiency | $/($m^3$/year) | Output |
| $y_2$ | SDG 6.4.2 | Share of Freshwater Withdrawal to Freshwater Availability | % | Output |
| $\theta_{20}$ | ICA | Cropland area actually irrigated (per crop type) | 1000 ha | Parameter |
| $\theta_{22}$ | Kc | Crop Factor (per crop type) | 1 | Parameter |
| $\theta_{26}, \theta_{28}$ | IRRTECHEFFi | Irrigation efficiency per irrigation technology | % | Parameter |
| $\theta_{27_{1,2,3,4}}$ | $\beta_{MWU}$ | Municipal Water Withdrawal coefficient vector | 1 | Parameter (I) |
| $\theta_{y_2}$ | EFR | Environmental Flow Requirement | $m^3$/year | Parameter |

required to ensure adaptation to spatial and temporal changes. Reviewing and analyzing the model can help ensure that it is aligned with the problem statement and the data requirements. This includes reviewing the policy needs and requirements, identifying country specifications, data sources, and quality, as well as changes over time to ensure that the model is designed to meet the desired SDG outcomes. Therefore, the modeling development phase includes

**Table 3. Structural model equations.** The function indices are named after the variable indices. $\hat{\theta}$ defines parameters that ensure the matching of dimensions or any physical conversions or attributes. The source of model equation is also indicated in the table.

| Function | Derived variable name | Relevant to | Equation | Reference |
|---|---|---|---|---|
| $f_{18}$ | Irrigated area per irrigation technology type (AIRi) | SDG 6.4.1, SDG 6.4.2 | $x_{18,i} = \hat{x}_{11,i}\hat{x}_2$ | [83] |
| $f_{19}$ | Agricultural Water Withdrawal (AWU) | SDG 6.4.1, SDG 6.4.2 | $x_{19} = x_{26}$ | [84] |
| $f_{20}$ | Cropping intensity (CI) | SDG 6.4.1, SDG 6.4.2 | $x_{20,i} = \frac{\theta_{20,i}}{\hat{x}_2}; i = 1, \ldots, \|\theta_{20}\|$ | [85] |
| $f_{21}$ | Corrective coefficient (Cr) | SDG 6.4.1 | $x_{21} = 1/(1 + (x_{28}/(1 - x_{28}) * \theta_{21}))$ | [86] |
| $f_{22}$ | Potential Crop Evaporation Vector (ETc) | SDG 6.4.1, SDG 6.4.2 | $x_{22} = \theta_{22} \cdot \mathbf{x}_{20} \cdot \hat{x}_8$ | [85] |
| $f_{23}$ | Irrigation Consumptive Use (ICU) | SDG 6.4.1, SDG 6.4.2 | $x_{23} = \|x_{22} - \hat{x}_7\|$ | [85] |
| $f_{25}$ | Irrigation Water Requirement per irrigation (IWRi) | SDG 6.4.1, SDG 6.4.2 | $x_{25,i} = \hat{\theta}x_{23}x_{18,i}$ | [85] |
| $f_{26}$ | Irrigation Water Withdrawal (IWW) | SDG 6.4.1, SDG 6.4.2 | $x_{26} = \sum_{i=1}^{\|x_{25}\|}(x_{25,i}/\theta_{26,i}) + \hat{x}_3\hat{\theta}_{\text{rice\_height}}$ | [84] |
| $f_{27}$ | Municipal Water Withdrawal (MWU) | SDG 6.4.1, SDG 6.4.2 | $x_{27} = e^{\theta_{27_1}}\hat{x}_{17}^{\theta_{27_2}}\hat{x}_9^{\theta_{27_3}}\hat{x}_{14}^{\theta_{27_4}}10^{-9}$ | [87] |
| $f_{28}$ | Proportion of Irrigated Cropland (PAIR) | SDG 6.4.1 | $x_{28} = (\sum_{i=1}^{\|\theta_{28}\|}\theta_{28,i})/\hat{x}_4$ | [86] |
| $f_{29}$ | Total Freshwater Available (TFA) | SDG 6.4.2 | $x_{29} = x_{30} + x_{31}$ | [88] |
| $f_{30}$ | Total Non-Conventional Water (TNCW) | SDG 6.4.2 | $x_{30} = \hat{x}_5 + \hat{x}_{16}$ | [88] |
| $f_{31}$ | Total Renewable Freshwater (TRF) | SDG 6.4.2 | $x_{31} = \hat{x}_6 + \hat{x}_{12}$ | [88] |
| $f_{32}$ | Total Water Withdrawal (TWW) | SDG 6.4.1, SDG 6.4.2 | $x_{32} = x_{27} + x_{19} + \hat{x}_{13}$ | [88] |
| $f_{y_1}$ | Water Use Efficiency (SDG 6.4.1) | SDG 6.4.1 | $y_1 = (\hat{x}_1(1 - x_{21}) + \hat{x}_{10} + \hat{x}_{15})/(x_{32}10^9)$ | [86] |
| $f_{y_2}$ | Share of Freshwater Withdrawal to Freshwater Availability (SDG 6.4.2) | SDG 6.4.2 | $y_2 = x_{32}/(x_{29} - \theta_{y_2})10^2$ | [88] |

https://doi.org/10.1371/journal.pone.0300531.t003

optimization of the model parameters, validation of the model against new data, and integration of new knowledge and insights into the model.

The GGSim model aims to predict several sustainable development goals. The water sub-models include the water use efficiency (SDG 6.4.1) and share of freshwater withdrawal to freshwater availability (SDG 6.4.2) as outputs. As the Shapley value evaluates the contribution of variables to the model output, first, the accuracy of the models should be evaluated. The predictions are evaluated for the GGSim model, two machine learning techniques, namely linear regression and k-th nearest neighbors, and the parameter optimization.

Fig 4 presents the results of different models for the efficiency of water use (SDG 6.4.1) and the share of freshwater withdrawal to freshwater availability (SDG 6.4.2) SDG indicators. The red line represents the observed data, whereas the blue dashed line shows the output of the GGSim model, which was built by experts. It seems that the error is somewhat systematic;
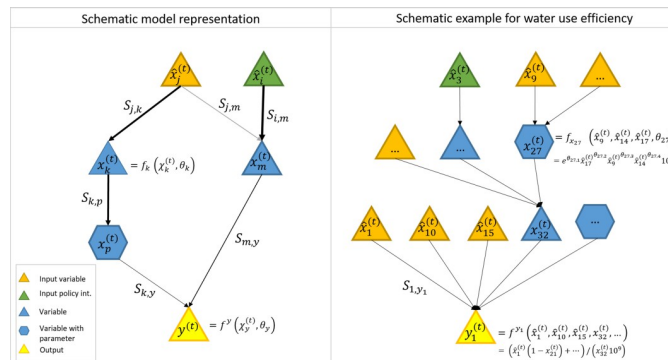


**Fig 3. Schematic representation of the Shapley value-based network interpretation of structural equation model and variable contributions.**
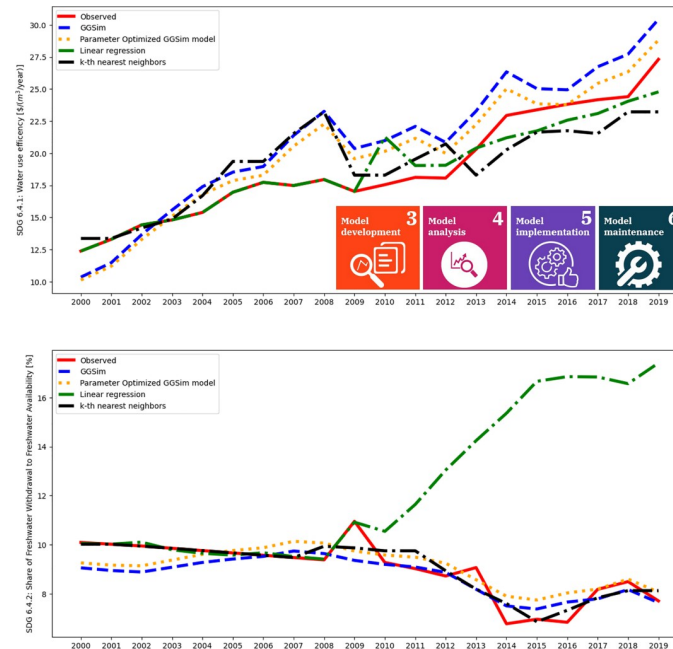
https://doi.org/10.1371/journal.pone.0300531.g003

**Fig 4. Prediction of the GGSim model against observed data.** The expert-built model (SDG 6.4.1 $r^2$: 0.519, SDG 6.4.2 $r^2$: 0.639) is also compared against a linear regression (SDG 6.4.1 $r^2$: 0.907, SDG 6.4.2 $r^2$: -20.46) and a k-th nearest neighbor algorithm (SDG 6.4.1 $r^2$: 0.91, SDG 6.4.2 $r^2$: 0.69). The k-th nearest neighbor shows the most promise, however, one cannot easily explain the inner structure of the black-box model. The parameter optimization of the municipal water withdrawal equation improved the model (SDG 6.4.1 $r^2$: 0.744, SDG 6.4.2 $r^2$: 0.65).

https://doi.org/10.1371/journal.pone.0300531.g004

however, some variables have no variance due to data imputation. Thus, the systematic error may be associated with a lack of adequate data. The $r^2$ was also calculated: for SDG 6.4.1 $r^2$: 0.519, SDG 6.4.2 $r^2$: 0.639. Additionally, parameter optimization was performed on the GGSim model for municipal water use. The SDG 6.4.1 model improved significantly, with a $r^2$ of 0.744, while SDG 6.4.2 only had an insignificant change in fitness ($r^2$ = 0.65).

Two machine learning models have been applied to the input data to predict the output while leaving out intermediate (derived) variables. Linear regression and the k-th nearest neighbor [89] models were trained every second year, starting from 2000 to 2009. It seems that linear regression cannot capture the complexity of the computation for SDG 6.4.2, and its $r^2$ is -20.46, while the SDG 6.4.1 model provides an acceptable fit (0.907). Although the coefficients are known, the white-box nature cannot be used due to the failing prediction of SDG 6.4.2. However, the k-th nearest neighbors method was able to fit a well-performing model on the observed output variable (SDG 6.4.1 $r^2$: 0.91, SDG 6.4.2 $r^2$: 0.69). Due to the black-box nature of the model, one may not be able to interpret why it provides predictions as such, but the improvement in accuracy reassured the connection between the inputs and outputs. If the model is hierarchical, it may also be hard to interpret; in some cases, the conversions and calculations are not trivial; therefore, it requires a method that is capable of describing the contributions between the variables.

This example has revealed the necessity for continuous model development and maintenance through data-driven machine learning techniques to ensure the quality of model performance. As is evident from the application of linear regression and k-th nearest-neighbor machine learning models in a given scenario, the model accuracy can vary significantly. Models should be trained over time to adapt to changes. It includes updating model parameters
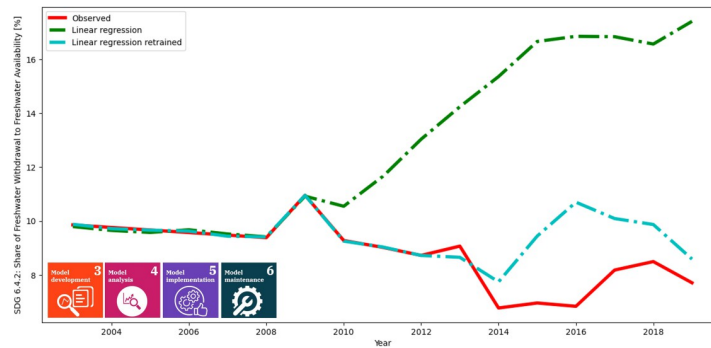
**Fig 5. Retrained linear regression for SDG 6.4.2.** The red line depicts the observed variable, the green one represents the linear regression from Fig 4, and the cyan colored one represent the retrained linear regression. The years before 2003 were dropped, and the model was retrained on years 2003–2012. The "new" model seems to fit better on the observed model, proving the necessity of continuous maintenance for machine learning models.

https://doi.org/10.1371/journal.pone.0300531.g005

and incorporating new data, and variables, thereby improving the predictive capability of the model. The application of the MLOps-based continuous development approach is a must to ensure that the models capture the complexity of the sustainability problems and iteratively enhance accuracy and robustness and provide more reliable predictions. Therefore, we retrained linear regression for SDG 6.4.2, which is presented in Fig 5. Here, the red line represents the observed variable, the green line represents the previous linear regression, and the cyan line represents the retrained linear regression. During retraining, we dropped the years from 2000–2002, and trained on from 2003 to 2012. The regression predicts the whole time interval. The new prediction seems to be closer to the observed data than the previous linear regression.

Emphasizing the significance of validated model structures is essential for performing in-depth analysis and model development. By focusing on validated models, data scientists and modelers can confidently investigate the underlying dynamics without concerns about flawed or unreliable outcomes. The reliable model structure facilitates the analysis of variable contributions to the model output and the identification of potential policy intervention points.

**Direct and indirect contribution of the variables to the model output.** Understanding the direct contribution of variables provides insight into which variables have the most significant impact on the output of the model. This information is crucial to identify the key drivers of the system being modeled, which can be used to inform decision-making and policy development. The indirect contribution of a variable refers to the impact that a variable has on the output of the model through its interactions with other variables in the network, which can allow us to identify complex relationships. During model validation, it is imperative to ensure that only the relevant variables remain part of the model. The Shapley value heavily relies on the role of the variable in the model and the variance of the data, and so ensuring variance may define the importance of a variable.

Figs 6 and 7 illustrate the indirect contribution of variables to water use efficiency and water stress level in 2017. It is advantageous to visualize how one variable changes the expected value; the impact of one particular year may help in deciding what policies should be implemented next year. The cumulative effects on the predictions can be measured, which provides information on how different variables may interact with each other. In this representation, the negative contribution values should not be associated with negative correlations to the output. The water use efficiency model is biased towards industrial water withdrawal (IWU), and service sectoral gross value added resources (SGVA), the data for 2017 indicate that most

**Fig 6. Indirect contribution of variables for SDG 6.4.1 in year 2017.** The dataset contains data imputations that were made with the 2017 data, therefore, this year is the best candidate for Shapley analysis. The X-axis shows the function value of SDG 6.4.1 in 2017, while the y axis presents how one variable changes the expected value of the function.

variables have no proper impact on the output of the model, except for what is required for the calculation of the expected value and conversions. For water efficiency (SDG 6.4.1), industrial water withdrawal, service sector gross value added resources, industrial gross value added (IGVA) and GDP per capita (GDPC) can be considered high-impact variables (Fig 6). For the level of water stress (EW2), industrial water withdrawal and GDP per capita are the drivers of contribution from the point of view of the SDG output indicator (Fig 7).

In the case of the water efficiency and water stress submodels, the decisive input in 2017 was the withdrawal of industrial water in Hungary; therefore, in order to improve the



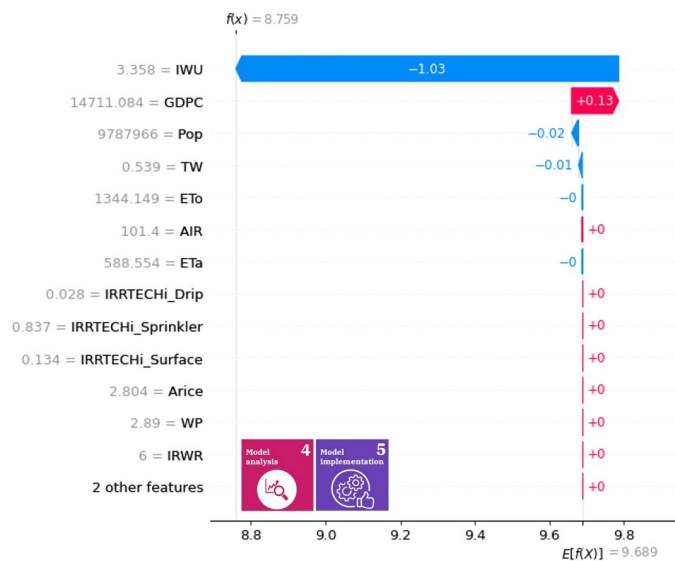**Fig 7. Indirect contribution of variables for SDG 6.4.2 in year 2017.** The X-axis shows the function value of SDG 6.4.2 in 2017, while the y axis presents how one variable changes the expected value of the function.

**Fig 8. Indirect mean Shapley values of variables for SDG 6.4.1.** The X-axis shows the average of caused (relative) change in the SDG 6.4.1 from 2000 to 2019, while the y axis presents the average impact of a variable.

https://doi.org/10.1371/journal.pone.0300531.g008

achievement of the 2030 Agenda, it is necessary to identify political interventions that reduce the value of this input, since the effect of IWU strongly affects the performance of SDG 6.4.1 & 6.4.2 indicators. However, due to the complexity of the SDG indicators, disaggregated inputs [90] are easier to handle. The reduction of water footprint [91], the circularity of water [92], and the nature-based solutions [93] can be effective tools to support SDG6.

By applying the framework we propose, it is possible to analyze the entire model structure, which also laid the foundation for the validation of the elements of the model. It is important to emphasize that all models included in this manuscript have undergone expert validation.

It is important to note that the effectiveness of individual state variables (policy intervention points) can vary depending on the current value of the other inputs, so policy monitoring is also an important task, which can be supported by our proposed Shapley value-based management. If it is not the given annual contributions but the tendentious driving forces that need to be identified in the SDG framework, the indirect mean Shapley values can be called upon.

The indirect mean Shapley values of the variables for the efficiency of water use are illustrated in Fig 8 and for the level of water stress in Fig 9. The mean contribution determines the
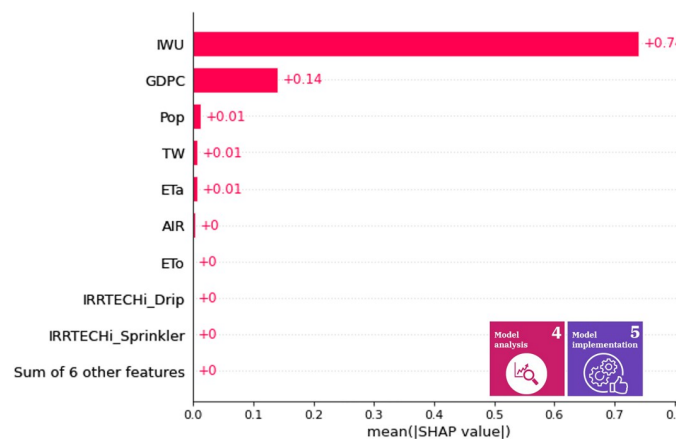


**Fig 9. Indirect mean Shapley values of variables for SDG 6.4.2 the level of water stress.** The X-axis shows the average of caused (relative) change in the SDG 6.4.2 from 2000 to 2019, while the y axis presents the average impact of a variable.

https://doi.org/10.1371/journal.pone.0300531.g009

significance of a variable over the years, so these mean contributions can be used as a general rule of thumb during the selection of appropriate policies. The high impact variables are similar compared to the 2017 data. Therefore, we can deduce that the roles do not differ significantly over time. In decision support, this representation may help in selecting features with high impact over time, model interpretation (by defining the high-impact variables), and validation. For example, if one variable relevance is a scientific fact and the mean absolute Shapley value is insignificant, then the model may have flaws that must be corrected before the actual information visualization is provided to the decision makers.

Fig 8 shows that the average contribution of the service sector to water efficiency is the most significant. This indicator indirectly characterizes tourist activity, the operation of hotels, restaurants, laundries, *etc*., thus affecting a wide range of water uses, and in SDG12 for the co-benefits that appear through food waste [94]. This high Shapley value draws attention to the fact that the potential of the service sector for sustainable watershed management in Hungary is worth considering as an effective point of political intervention.

The indirect mean Shapley values of the variables for the level of water stress are illustrated in Fig 9. The decisive role of industrial water withdrawals in water stress is outstanding when examined over the entire period. In this case, the importance of productivity developments comes to the fore, which already show a decreasing trend in water stress in several developed countries [95]. Water stress is an excellent example of the fact that SDG indicators are difficult to directly regulate, so the identification of intervention points for which effective political measures [96] can be formulated is essential for the implementation of the 2030 Agenda.

The direct contribution of the variables to the changes in the output variables is shown in Figs 10 and 11. Nodes are differentiated by their color and shape, with policy intervention points represented by green triangles ($\hat{x}_3, \hat{x}_{11}$), input variables by orange triangles ($\hat{x}_k, k = 1, \ldots, 17$), derived variables by blue triangles ($x_k, k = 18, \ldots, 32$), variables with (possibly identifiable) parameters by blue hexagons ($x_{21}, x_{27}$), and output SDG indicators by yellow triangles ($y_1, y_2$). The thickness (the greater the better) of the arrows represents the strength of direct variable contribution (the width is proportional to the individual to total contribution ratio). Here, the percentage of the total average contribution is provided for each arrow with respect to all inputs of a model. The gray arrows depict the zero contribution.

Figs 10 and 11 show that there are aggregated and raw data input sources to describe the aggregated SDG indicators, which can be used to identify optimal intervention points (policies) based on their systematic exploration and impact of the contribution on the SDG output, and this systematic exploration can also be used to check the effectiveness of existing policies.



**Fig 10. Direct mean contribution of variables to the change in the output variables.** The types of nodes are represented with different colors and shapes (green triangle—policy intervention points; orange triangle—input variables; blue triangle—variables; blue hexagon—variables with parameters; yellow triangle—output). The thickness of the arrows represents the strength of direct variable contribution.

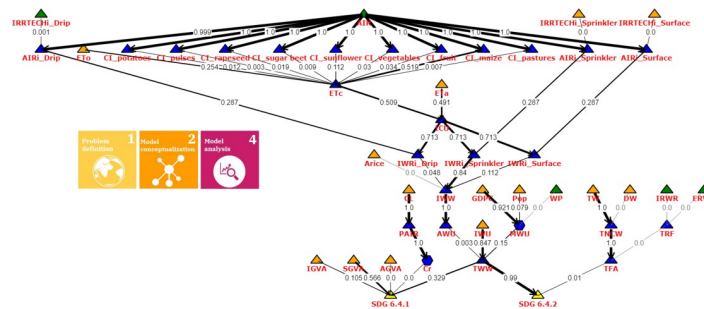https://doi.org/10.1371/journal.pone.0300531.g010

**Fig 11. Direct contribution of variables to the change in the output variables for 2017.** The type of nodes are represented with different colors and shapes (green triangle—policy intervention points; orange triangle—input variables; blue triangle—variables; blue hexagon—variables with parameters; yellow triangle—output). The thickness of the arrows represents the strength of direct variable contribution.
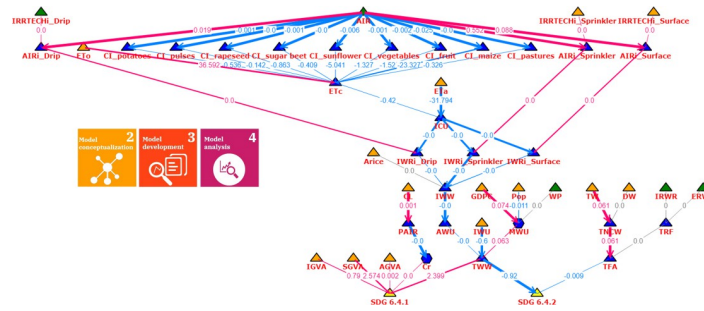
Fig 10 shows that irrigated agricultural areas (AIR) influence irrigation water use (IWW) through crop specifications (ETc, ICU), which in relation to total water withdrawal (TWW) affects the level of water stress (SDG 6.4.2) the most. It seems that the water price (WP), the internal and external renewable water resources (IRWR, ERWR) as political intervention points did not have a significant impact in the period examined based on Hungarian data, so it is more expedient to develop the service sector instead (SGVA), or improve water retention (ETa), more precisely in developed and scheduled irrigation solutions, such as the internet of things and wireless sensor network technologies [97]. Based on the model structure, the other option is to reduce industrial water use, which can be supported by moving toward water reuse and recycling [98] and technological solutions saving water [99]. The role of the variables can also be analyzed in a yearly breakdown, which lays the foundation for a better understanding of the dynamics of the SDG indicators and their input data sources. The annual contribution of the variables in 2017 can be seen in Fig 11.

The values on the edges show the contribution relative to the value of the output variable of the model, so if the change in the output is low, the relative contribution values on the edges are also low. However, the thickness of the edges symbolizes the degree of contribution of the given variable, so the absolute and relative effects can be read together for all SDG indicators included in the model. As can be seen in Fig 11, the effect of the progress achieved in industrial water use on both SDG6 indicators is even more pronounced than in the case of the mean time series contribution (Fig 10).

Conducting the analysis of the direct and indirect contribution of the variables to the model output reveals key information regarding the drivers of the models, which serves as validation of the developed model, a base for further improvement, and enabling informed decision-making and selection of key policy intervention points through understanding the underlying interactions within the variables and the model output.

**Selection and validation of policy intervention points.** The selection and validation of policy intervention points is an important step in the development of effective policies and interventions that can improve outcomes in the water model. We utilized the potential of sensitivity and network analysis to validate policy intervention points and to support the selection of possible ones. By testing the sensitivity of the model to changes in the value of intervention points, we can identify the points that have the greatest potential to improve outcomes. This information can be used to prioritize interventions and focus resources on the most effective intervention points. The tools of network science enable us to identify critical paths and nodes within the water model, based on which information we can focus on the most important
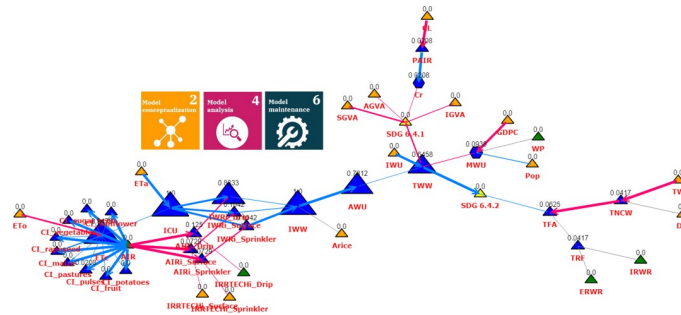
**Fig 12. Network representation of the betweenness centralities for the SDG 6.4.1 & SDG 6.4.2 indicators in water model.** The size of the nodes presents the importance of centralities. The types of nodes are represented with different colors and shapes (green triangle—policy intervention points; orange triangle—input variables; blue triangle—variables; blue hexagon—variables with parameters; yellow triangle—output). The thickness of the arrows represents the strength of the direct variable contribution.

https://doi.org/10.1371/journal.pone.0300531.g012

variables in the system, and develop interventions that target these variables. Fig 12 shows the betweenness centrality of network water efficiency (6.4.1) and water stress (6.4.2) SDG indicators.

In the network of Fig 12, the size of the nodes is the same as the degree of their intermediary role (betweenness centrality measure), so the larger nodes represent successful mediators, i.e. potential political interveners, while the influence of the smaller nodes is less significant in the state changes of the SDG system of the water model. Regarding the intermediary role, in the case of the SDG6 indicators, crop evaporation (Etc), irrigation (ICU, IWW), agricultural water use (AWU) and total water withdrawal (TWW) play a prominent role in the network, so a direct or indirect reduction of these node values may mean political implementation potential. Based on the example shown in Fig 12, it can also be seen that if the used state variables e.g. water withdrawals cannot be reduced directly, the model extension with additional politically perturbable variables can be performed based on the network science-based model analysis.

Utilizing network science tools facilitates the identification and validation of crucial policy intervention points, a key aspect of model validation and development. By employing these tools, researchers and decision-makers can rigorously verify the effectiveness and reliability of proposed policy changes before implementing them in real-world scenarios, ensuring a more robust and informed decision-making process.

## Discussion

The proper translation of information empowers decision-makers. The decision support often relies on the visualization of data and models. Dashboards with incorporated visualizations and hierarchical structures can help in identifying trends, relationships between variables, sensitivities of features *etc*. Knowledge, insight, and information transfer is required to improve the decision-making capacity of stakeholders [100]. The role of explainable AI is to provide knowledge transfer, when the actual knowledge is hidden [101]. Shapley-based networks have already been researched [102], however, we focus on their implications for decision-makers and their connection with MLOps practices. We must note that the Shapley values may not be interpretable if the data is incorrect, the model acts on incorrect assumptions, has inappropriate biases, or the incorrect set of variables is used as inputs [103], which is why checking and validating the prediction ability of the models is crucial before understanding their structure. The accuracy of model evaluation relies heavily on the quality of the underlying data, and poor data quality can result in inaccurate assessments of model performance and variable effects.

Additionally, the need for generally accepted data sources is emphasized, especially when dealing with diverse economies or environmental backgrounds. Furthermore, addressing structural flaws in models is a key focus of MLOps, which involves retraining existing models or creating new ones by incorporating new variables and insights.

However, the method that facilitates the identification of potential intervention points cannot be utilized to propose exact policies. Furthermore, the method was developed to support decision-makers rather than substitute them. Also, the proposed case study is an isolated sub-model of the whole SDG framework, and the change of one variable may seem to be very effective; it may negatively influence other areas and models, provided that the models are connected. Consider the following counterfactual: removing rice paddies may decrease water use and improve water use efficiency; however, the harvest will also be reduced, and so it will affect the food market. In other words, water use is improved, but the number of starving people can increase due to costly imports and reduced harvests. Therefore, the selection of the policy intervention is a great balancing act. The change in one variable may initiate a trade-off between several of the outputs, even if the models may not account for it. Furthermore, for one country in a global economy, it may seem reasonable to take actions, *e.g.* reduce rice production, if all countries do the same, which may affect the global food supply. Therefore, there is a need to evaluate policy-making in a system of systems settings.

Regarding the evaluation of sustainability models with Shapley-based network analysis, the results present a handful of irrelevant variables to the change in prediction. Namely, water price, rice paddy irrigation area, and internal and external renewable water resources lack variance, therefore they have no impact. The agricultural gross value added (AGVA) also has a minuscule contribution due to its small size compared to its industrial (IGVA) and service sectoral (SGVA) versions. Generally speaking, agriculture draws the most water worldwide (about 69%) [104], due to crop and animal needs, and therefore plays a significant role in feeding the population. The output shows a one-sided picture of water use. The description of water use efficiency can be understood as the economic value per cubic meter of water [105]. Therefore, the interpretation may have methodological flaws [106], as it may not consider other important factors. The SDG indicator may be related to behavioral factors that impact municipal water withdrawal [104]. Here, the water price indicator can be used as an incentive to control behavior, yet this variable does not have variance in the Hungarian data set provided by the GGGI. However, the price of water was selected as an intervention point in the original model.

The application and key scientific contribution of the MLOps methodology presented in the research are two-fold since the experts who create the models can represent the entire model structure, therefore it is possible to validate the model features. In addition, our proposed methodology also supports problem exploration and understanding through data. In other words, data-based qualitative validation and network science, as well as Shapley value-based qualitative validation, can be implemented together, thereby effectively supporting modeling processes.

## Conclusion

This paper emphasizes the importance of developing life cycle-based models that incorporate network and data analysis practices in sustainability planning, as demonstrated through an application study of the Hungarian water model developed by the Global Green Growth Institution. To accurately track progress toward achieving the Sustainable Development Goals (SDGs), it is necessary to use data-driven models that are country-specific and aligned with policy interventions. However, modeling such complex relationships is often difficult due to

the diverse scale and quality of data, policy initiatives, and economic and social developments across countries. Continuous life cycle-based revision, analysis, development, and evaluation of environmental models are needed to ensure model adaptation to spatial and temporal changes.

Therefore, we propose the utilization of tools based on network and data science throughout the modeling life cycle to ensure more accurate models to assess the SDGs and evidence-based policy-making. The potential of Shapley value has been introduced to identify key drivers of the system being modeled and to support decision-making and validate policy interventions, which facilitates understanding of the direct and indirect contribution of variables. This enables greater insight into which variables have the most significant impact on the output of the model (SDG indicators). Furthermore, we suggest using a Shapley value-driven sensitivity analysis of the changes in intervention point values, and network analysis to identify critical paths and nodes in the selecting and validating of policy intervention points. Ensuring the accuracy of model evaluation is imperative, hinging on the detailed handling of underlying data quality to prevent incorrect conclusions about the model's performance. It is crucial to establish homogeneous, generally accepted samples, particularly when considering diverse economic and environmental backgrounds, underscoring the need for a clear definition in the modeled area. Future development plans should prioritize advancements in model evaluation by refining methodologies for handling data, while also focusing on the continuous innovation of MLOps to proactively identify and rectify structural flaws in models through either retraining existing models or creating new ones with updated variables and insights.

The real-life application study of network and data analysis for the development of the Hungarian water model underlines the efficiency of the proposed aspects. We believe that the proposed data-driven life cycle management of sustainability models has great potential in real-life sustainability planning and decision-making at any administrative level including model identification, development and validation.

## Supporting information

**S1 Fig.**
(TIF)

**S2 Fig.**
(TIF)

**S3 Fig.**
(TIF)

**S4 Fig.**
(TIF)

**S5 Fig.**
(TIF)

**S6 Fig.**
(TIF)

## Author Contributions

**Conceptualization:** Lilibeth A. Acosta, Viktor Sebestyén, János Abonyi.

**Data curation:** Sanga Lee, Innocent Nzimenyera.

**Formal analysis:** Ádám Ipkovich, Tímea Czvetkó.

**Funding acquisition:** Ádám Ipkovich, Viktor Sebestyén, János Abonyi.

**Methodology:** Ádám Ipkovich, Tímea Czvetkó, Viktor Sebestyén, János Abonyi.

**Resources:** Lilibeth A. Acosta, Sanga Lee, Innocent Nzimenyera.

**Software:** Ádám Ipkovich.

**Supervision:** Lilibeth A. Acosta, Viktor Sebestyén, János Abonyi.

**Validation:** Lilibeth A. Acosta, Sanga Lee, Innocent Nzimenyera, Viktor Sebestyén, János
 Abonyi.

**Visualization:** Ádám Ipkovich.

**Writing – original draft:** Ádám Ipkovich, Tímea Czvetkó, Viktor Sebestyén.

# References

1. Nations U. Transforming our world: The 2030 agenda for sustainable development. New York: United Nations, Department of Economic and Social Affairs. 2015;.

2. Purvis B, Mao Y, Robinson D. Three pillars of sustainability: in search of conceptual origins. Sustainability science. 2019; 14(3):681–695. https://doi.org/10.1007/s11625-018-0627-5

3. Yamane T, Kaneko S. Impact of raising awareness of Sustainable Development Goals: A survey experiment eliciting stakeholder preferences for corporate behavior. Journal of Cleaner Production. 2021; 285:125291. https://doi.org/10.1016/j.jclepro.2020.125291

4. Nastasi B, Markovska N, Puksec T, Duić N, Foley A.: Renewable and sustainable energy challenges to face for the achievement of Sustainable Development Goals. Elsevier.

5. Sebestyén V, Abonyi J. Data-driven comparative analysis of national adaptation pathways for Sustainable Development Goals. Journal of Cleaner Production. 2021; 319:128657. https://doi.org/10.1016/j.jclepro.2021.128657

6. Singh RK, Murty HR, Gupta SK, Dikshit AK. An overview of sustainability assessment methodologies. Ecological indicators. 2009; 9(2):189–212. https://doi.org/10.1016/j.ecolind.2008.05.011

7. Pejić-Bach M, Čerić V. Developing system dynamics models with "step-by-step" approach. Journal of information and organizational sciences. 2007; 31(1):171–185.

8. Honti G, Dörgő G, Abonyi J. Review and structural analysis of system dynamics models in sustainability science. Journal of Cleaner Production. 2019; 240:118015. https://doi.org/10.1016/j.jclepro.2019.118015

9. Sebestyén V, Bulla M, Rédey Á, Abonyi J. Network model-based analysis of the goals, targets and indicators of sustainable development for strategic environmental assessment. Journal of environmental management. 2019; 238:126–135. PMID: 30849597

10. Dhar V. Data science and prediction. Communications of the ACM. 2013; 56(12):64–73. https://doi.org/10.1145/2500499

11. Treveil M, Omont N, Stenac C, Lefevre K, Phan D, Zentici J, et al. Introducing MLOps. O'Reilly Media; 2020.

12. Studer S, Bui TB, Drescher C, Hanuschkin A, Winkler L, Peters S, et al. Towards CRISP-ML (Q): a machine learning process model with quality assurance methodology. Machine Learning and Knowledge Extraction. 2021; 3(2):392–413. https://doi.org/10.3390/make3020020

13. Anderson CC, Denich M, Warchold A, Kropp JP, Pradhan P. A systems model of SDG target influence on the 2030 Agenda for Sustainable Development. Sustainability science. 2022; 17(4):1459–1472. https://doi.org/10.1007/s11625-021-01040-8 PMID: 34659581

14. Bennich T, Weitz N, Carlsen H. Deciphering the scientific literature on SDG interactions: A review and reading guide. Science of the Total Environment. 2020; 728:138405. https://doi.org/10.1016/j.scitotenv.2020.138405 PMID: 32388023

15. Di Vaio A, Palladino R, Hassan R, Escobar O. Artificial intelligence and business models in the sustainable development goals perspective: A systematic literature review. Journal of Business Research. 2020; 121:283–314. https://doi.org/10.1016/j.jbusres.2020.08.019

16.  Scharlemann JP, Brock RC, Balfour N, Brown C, Burgess ND, Guth MK, et al. Towards understanding interactions between Sustainable Development Goals: The role of environment–human linkages. Sustainability science. 2020; 15(6):1573–1584. https://doi.org/10.1007/s11625-020-00799-6

17.  Moallemi EA, Bertone E, Eker S, Gao L, Szetey K, Taylor N, et al. A review of systems modelling for local sustainability. Environmental Research Letters. 2021; 16(11):113004. https://doi.org/10.1088/1748-9326/ac2f62

18.  MLOps. ml-ops org, editor.: CRISP-ML(Q). The ML Lifecycle Process.

19.  Acosta LA, Gerrard SP, Luchtenbelt HGH, Nazareth M, Ruben Sabado J, Eugenio JR, et al. Green Growth Simulation Tool Phase 1—Concept, Methods and Applications. GGGI Technical Report No 17. 2020;.

20.  Data MC, Komorowski M, Marshall DC, Salciccioli JD, Crutain Y. Exploratory data analysis. Secondary analysis of electronic health records. 2016;p. 185–203.

21.  Mikut R, Reischl M. Data mining tools. Wiley interdisciplinary reviews: data mining and knowledge discovery. 2011; 1(5):431–443.

22.  Zhang B, Tay FE. An Integrated Approach Using Data Mining and System Dynamics to Policy Design: Effects of Electric Vehicle Adoption on CO 2 Emissions in Singapore. In: Industrial Conference on Data Mining. Springer; 2017. p. 258–268.

23.  Hosseini HM, Kaneko S. Causality between pillars of sustainable development: Global stylized facts or regional phenomena? Ecological Indicators. 2012; 14(1):197–201. https://doi.org/10.1016/j.ecolind.2011.07.005

24.  Dörgő G, Sebestyén V, Abonyi J. Evaluating the interconnectedness of the sustainable development goals based on the causality analysis of sustainability indicators. Sustainability. 2018; 10(10):3766. https://doi.org/10.3390/su10103766

25.  Swain RB, Ranganathan S. Modeling interlinkages between sustainable development goals using network analysis. World Development. 2021; 138:105136. https://doi.org/10.1016/j.worlddev.2020.105136

26.  Qazi A, Angell LC, Daghfous A, Al-Mhdawi M. Network-based risk assessment of country-level sustainable development goals. Environmental Impact Assessment Review. 2023; 99:107014. https://doi.org/10.1016/j.eiar.2022.107014

27.  Nazabal A, Williams CK, Colavizza G, Smith CR, Williams A. Data engineering for data analytics: a classification of the issues, and case studies. arXiv preprint arXiv:200412929. 2020;.

28.  Alharbi Y, Arribas-Bel D, Coenen F. Sustainable development goals monitoring and forecasting using time series analysis. In: PROCEEDINGS OF THE 2ND INTERNATIONAL CONFERENCE ON DEEP LEARNING THEORY AND APPLICATIONS (DELTA). SCITEPRESS-Science and Technology Publications; 2021. p. 123–131.

29.  Holloway J, Mengersen K. Statistical machine learning methods and remote sensing for sustainable development goals: a review. Remote Sensing. 2018; 10(9):1365. https://doi.org/10.3390/rs10091365

30.  Fakhimi M, Stergioulas L, Mustafee N, Eldabi T. A review of literature in modeling approaches for sustainable development. In: 2013 Winter simulations conference (WSC). IEEE; 2013. p. 282–290.

31.  Hjorth P, Bagheri A. Navigating towards sustainable development: A system dynamics approach. Futures. 2006; 38(1):74–92. https://doi.org/10.1016/j.futures.2005.04.005

32.  Mirghaderi SH. Using an artificial neural network for estimating sustainable development goals index. Management of Environmental Quality: An International Journal. 2020;. https://doi.org/10.1108/MEQ-12-2019-0266

33.  Düspohl M, Frank S, Döll P. A review of Bayesian networks as a participatory modeling approach in support of sustainable environmental management. Journal of Sustainable Development. 2012; 5 (12):1–18.

34.  Kwatra S, Kumar A, Sharma P. A critical review of studies related to construction and computation of Sustainable Development Indices. Ecological Indicators. 2020; 112:106061. https://doi.org/10.1016/j.ecolind.2019.106061

35.  Hassani H, Huang X, MacFeely S. Enabling Digital Twins to Support the UN SDGs. Big Data and Cognitive Computing. 2022; 6(4):115. https://doi.org/10.3390/bdcc6040115

36.  Le Blanc D. Towards integration at last? The sustainable development goals as a network of targets. Sustainable Development. 2015; 23(3):176–187. https://doi.org/10.1002/sd.1582

37.  Dalampira ES, Nastis SA. Mapping sustainable development goals: A network analysis framework. Sustainable Development. 2020; 28(1):46–55. https://doi.org/10.1002/sd.1964

38.  Laumann F, von Kügelgen J, Uehara THK, Barahona M. Complex interlinkages, key objectives, and nexuses among the Sustainable Development Goals and climate change: a network analysis. The

Lancet Planetary Health. 2022; 6(5):e422–e430. https://doi.org/10.1016/S2542-5196(22)00070-5 PMID: 35550081

39. Sebestyén V, Bulla M, Rédey Á, Abonyi J. Data-driven multilayer complex networks of sustainable development goals. Data in brief. 2019; 25:104049. https://doi.org/10.1016/j.dib.2019.104049 PMID: 31194124

40. Sebestyén V, Domokos E, Abonyi J. Focal points for sustainable development strategies—Text mining-based comparative analysis of voluntary national reviews. Journal of Environmental Management. 2020; 263:110414. PMID: 32174539

41. Cosenz F, Rodrigues VP, Rosati F. Dynamic business modeling for sustainability: Exploring a system dynamics perspective to develop sustainable business models. Business Strategy and the Environment. 2020; 29(2):651–664. https://doi.org/10.1002/bse.2395

42. Narayanam R, Narahari Y. A Shapley Value-Based Approach to Discover Influential Nodes in Social Networks. IEEE Transactions on Automation Science and Engineering. 2011; 8(1):130–147. https://doi.org/10.1109/TASE.2010.2052042

43. Dorgo G, Honti G, Abonyi J. Automated analysis of the interactions between sustainable development goals extracted from models and texts of sustainability science. Chemical Engineering Transactions. 2018; 70:781–786.

44. Orellana DFP, Piedra N. Semantic Enrichment of Open Dataset related to sustainable Development Goals using Open Knowledge Graphs. In: 2021 XVI Latin American Conference on Learning Technologies (LACLO). IEEE; 2021. p. 470–473.

45. Eguiguren JE, Piedra N. Connecting Open Data and Sustainable Development Goals using a Semantic Knowledge Graph approach. 2019;.

46. Bonanni L, Ebner H, Hockenberry M, Sayan B, Zapico Lamela JL, Brandt N, et al. The Open Sustainability Project: A Linked Data Approach to LCA. LCA X, Bridging Science, Policy, and the Public 2-4 November 2010, Portland, Oregon. 2010;.

47. Perez A, Larrinaga F, Curry E. The role of linked data and semantic-technologies for sustainability idea management. In: Software Engineering and Formal Methods: SEFM 2013 Collocated Workshops: BEAT2, WS-FMDS, FM-RAIL-Bok, MoKMaSD, and OpenCert, Madrid, Spain, September 23-24, 2013, Revised Selected Papers 11. Springer; 2014. p. 306–312.

48. Serra F, Delgado T. DW2RDF4SDG–Ontology modeling from multi-dimensional cubes for Sustainable Development Goals. Sistemas & Telemática. 2018; 16(44):9–24.

49. Requejo-Castro D, Giné-Garriga R, Pérez-Foguet A. Data-driven Bayesian network modelling to explore the relationships between SDG 6 and the 2030 Agenda. Science of the total environment. 2020; 710:136014. https://doi.org/10.1016/j.scitotenv.2019.136014 PMID: 32050357

50. Ospina-Forero L, Castañeda G, Guerrero OA. Estimating networks of sustainable development goals. Information & Management. 2022; 59(5):103342. https://doi.org/10.1016/j.im.2020.103342

51. Zelinka D, Amadei B. A systems approach for modeling interactions among the Sustainable Development Goals Part 2: System dynamics. International Journal of System Dynamics Applications (IJSDA). 2019; 8(1):41–59. https://doi.org/10.4018/IJSDA.2019010103

52. Lemaire GG, Carnohan SA, Grand S, Mazel V, Bjerg PL, McKnight US. Data-Driven System Dynamics Model for Simulating Water Quantity and Quality in Peri-Urban Streams. Water. 2021; 13(21): 3002. https://doi.org/10.3390/w13213002

53. Yeh C, Meng C, Wang S, Driscoll A, Rozi E, Liu P, et al. SustainBench: Benchmarks for Monitoring the Sustainable Development Goals with Machine Learning. arXiv preprint arXiv:211104724. 2021;.

54. Czvetko T, Honti G, Sebestyen V, Abonyi J. The intertwining of world news with Sustainable Development Goals: An effective monitoring tool. Heliyon. 2021; 7(2):e06174. https://doi.org/10.1016/j.heliyon.2021.e06174 PMID: 33598579

55. Marcovecchio I, Thinyane M, Estevez E, Fillottrani P. Capability maturity models towards improved quality of the sustainable development goals indicators data. In: 2017 ITU Kaleidoscope: Challenges for a Data-Driven Society (ITU K). IEEE; 2017. p. 1–8.

56. fan Y, Chen J, Shirkey G, John R, Wu R, Park H, et al. Applications of structural equation modeling (SEM) in ecological research: An updated review. Ecological Processes. 2016 10;5. https://doi.org/10.1186/s13717-016-0063-3

57. Shapley LS. Notes on the N-Person Game–II: The Value of an N-Person Game. Santa Monica, CA: RAND Corporation; 1951.

58. Narayanam Ramasuri and Narahari Yadati A shapley value-based approach to discover influential nodes in social networks. IEEE transactions on automation science and engineering, 8, 1, 130–147. 2010. https://doi.org/10.1109/TASE.2010.2052042

**59.** Aadithya, Karthik V and Ravindran, Balaraman and Michalak, Tomasz P and Jennings, Nicholas R Efficient computation of the shapley value for centrality in networks. Internet and Network Economics: 6th International Workshop, WINE 2010, Stanford, CA, USA, December 13-17, 2010. Proceedings 6 Springer, 1-13. 2010.

**60.** Chen, Wei and Teng, Shang-Hua Interplay between social influence and network centrality: a comparative study on shapley centrality and single-node-influence centrality Proceedings of the 26th international conference on world wide web, 967. 2017.

**61.** Szczepanski, Piotr L and Michalak, Tomasz and Rahwan, Talal A new approach to betweenness centrality based on the shapley value 2012

**62.** Ahmad Ashfaq and Akbar Shahid and Tahir Muhammad and Hayat Maqsood and Ali Farman. iAFPs-EnC-GA: identifying antifungal peptides using sequential and evolutionary descriptors based multi-information fusion and ensemble learning approach. Chemometrics and Intelligent Laboratory Systems. Elsevier. 2022. 222. 104516. https://doi.org/10.1016/j.chemolab.2022.104516

**63.** Akbar, Shahid and Raza, Ali and Al Shloul, Tamara and Ahmad, Ashfaq and Saeed, Aamir and Ghadi, Yazeed Yasin et al. pAtbP-EnC: identifying anti-tubercular peptides using multi-feature representation and genetic algorithm based deep ensemble model. IEEE,Access. IEEE. 2023.

**64.** Raza Ali and Uddin Jamal and Almuhaimeed Abdullah and Akbar Shahid and Zou Quan and Ahmad Ashfaq. AIPs-SnTCN: Predicting Anti-Inflammatory Peptides Using fastText and Transformer Encoder-Based Hybrid Word Embedding with Self-Normalized Temporal Convolutional Networks. Journal of Chemical Information and Modeling. ACS Publications. 2023. 63. 21. 6537–6554. https://doi.org/10.1021/acs.jcim.3c01563 PMID: 37905969

**65.** Molnar C. Interpretable Machine Learning; 2019. https://christophm.github.io/interpretable-ml-book/.

**66.** Štrumbelj E, Kononenko I. Explaining prediction models and individual predictions with feature contributions. Knowledge and Information Systems. 2013 12;41:647–665. https://doi.org/10.1007/s10115-013-0679-x

**67.** Barrat A, Barthélemy M, Pastor-Satorras R, Vespignani A. The architecture of complex weighted networks. Proceedings of the National Academy of Sciences. 2004; 101(11):3747–3752. https://doi.org/10.1073/pnas.0400087101 PMID: 15007165

**68.** Sabidussi G. The centrality index of a graph. Psychometrika. 1966; 31:581–603. https://doi.org/10.1007/BF02289527 PMID: 5232444

**69.** White DR, Borgatti SP. Betweenness centrality measures for directed graphs. Social Networks. 1994; 16(4):335–346. https://doi.org/10.1016/0378-8733(94)90015-9

**70.** De Meo P, Ferrara E, Fiumara G, Provetti A. Generalized Louvain method for community detection in large networks. In: 2011 11th International Conference on Intelligent Systems Design and Applications; 2011. p. 88–93.

**71.** Rosita YD, Rosyida EE, Rudiyanto MA. Implementation of Dijkstra Algorithm and Multi-Criteria Decision-Making for Optimal Route Distribution. Procedia Computer Science. 2019;161:378–385. The Fifth Information Systems International Conference, 23-24 July 2019, Surabaya, Indonesia.

**72.** Food and Agriculture Organization of the United Nations.: Irrigated crop calendars. Available from: https://www.fao.org/aquastat/en/databases/crop-calendar.

**73.** Food and Agriculture Organization of the United Nations.: Core Database. Available from: https://www.fao.org/aquastat/en/databases/maindatabase.

**74.** Brouwer C, Prins K, Heibloem M. Irrigation Water Management: Irrigation Scheduling. Training Manual (No 4). 1989;.

**75.** Food and Agriculture Organization of the United Nations.: FAOSTAT Statistical Database. Available from: https://www.fao.org/faostat/en/#data.

**76.** Global Perspectives Studies (GPS) Team.: Projections of future total renewable water resources (TWR) by country for different climate change scenarios available based on ISI-MIP data. Available from: https://www.fao.org/global-perspectives-studies/resources/detail/en/c/1157059/.

**77.** Running S, Mu Q, Zhao M, Moreno A.: MODIS Global Terrestrial Evapotranspiration (ET) Product (NASA MOD16A2/A3) Algorithm Theoretical Basis Document, Collection 5. Available from: https://lpdaac.usgs.gov/documents/93/MOD16_ATBD.pdf.

**78.** Running S, Mu Q, Zhao M, Moreno A.: MOD16A3GF MODIS/Terra Net Evapotranspiration Gap-Filled Yearly L4 Global 500 m SIN Grid V006 [Data set].

**79.** World Bank.: Water Efficiency. Available from: https://data.worldbank.org/.

**80.** United Nations Development Programme.: Human Development Index. Human Development Report 2021-22. Available from: http://hdr.undp.org/en/composite/HDI.

**81.** WHO and UNICEF.: JMP Data. Available from: https://washdata.org/data.

**82.** Hungary Ministry of innovation and technology.: National Clean Development Strategy 2020-2050.

**83.** Food and Agriculture Organization for the United Nation AQUASTAT—FAO's Global Information System on Water and Agriculture 2024 https://www.fao.org/aquastat/en/databases/maindatabase/;

**84.** Brouwer, C and Prins, Kees and Heibloem, Marjan Irrigation water management: Training manual no. 4: Irrigation scheduling Rome, Italy: FAO, 1985.

**85.** Allen, Richard G and Pereira, Luis S and Raes, Dirk and Smith, Martin and others Crop evapotranspiration-Guidelines for computing crop water requirements-FAO Irrigation and drainage paper 56. Fao, Rome. 300,9. 1998

**86.** UNSTATS Indicator 6.4.1: Change in water-use efficiency over time. United Nations Statistics Division, 1-28. 2018.

**87.** Hejazi Mohamad and Edmonds James and Chaturvedi Vaibhav and Davies Evan and Eom Jiyong Scenarios of global municipal water-use demand projections over the 21st century Hydrological Sciences Journal, Taylor & Francis, 58,3, 519–538, 2013 https://doi.org/10.1080/02626667.2013.772301

**88.** UNSTATS STEP-BY-STEP MONITORING METHODOLOGY FOR INDICATOR 6.4.2 Integrated Monitoring Guide for SDG 6, Step-by-step monitoring methodology for indicator 6.4.2 on water stress United Nations Statistics Division, 1-11. 2017.

**89.** Azadkia M.: Optimal choice of $k$ for $k$-nearest neighbor regression.

**90.** Hellegers P, van Halsema G. SDG indicator 6.4. 1 "change in water use efficiency over time": Methodological flaws and suggestions for improvement. Science of the Total Environment. 2021; 801:149431. https://doi.org/10.1016/j.scitotenv.2021.149431 PMID: 34411791

**91.** Weerasooriya R, Liyanage L, Rathnappriya R, Bandara W, Perera T, Gunarathna M, et al. Industrial water conservation by water footprint and sustainable development goals: a review. Environment, Development and Sustainability. 2021;p. 1–49.

**92.** Kakwani NS, Kalbar PP. Measuring urban water circularity: Development and implementation of a Water Circularity Indicator. Sustainable Production and Consumption. 2022; 31:723–735. https://doi.org/10.1016/j.spc.2022.03.029

**93.** Nika C, Gusmaroli L, Ghafourian M, Atanasova N, Buttiglieri G, Katsou E. Nature-based solutions as enablers of circularity in water systems: A review on assessment methodologies, tools and indicators. Water research. 2020; 183:115988. https://doi.org/10.1016/j.watres.2020.115988 PMID: 32683049

**94.** Beretta C, Hellweg S. Potential environmental benefits from food waste prevention in the food service sector. Resources, Conservation and Recycling. 2019; 147:169–178. https://doi.org/10.1016/j.resconrec.2019.03.023

**95.** Doeffinger T, Hall JW. Water stress and productivity: an empirical analysis of trends and drivers. Water Resources Research. 2020; 56(3):e2019WR025925. https://doi.org/10.1029/2019WR025925

**96.** Glass LM, Newig J. Governance for achieving the Sustainable Development Goals: How important are participation, policy coherence, reflexivity, adaptation and democratic institutions? Earth System Governance. 2019; 2:100031. https://doi.org/10.1016/j.esg.2019.100031

**97.** García L, Parra L, Jimenez JM, Lloret J, Lorenz P. IoT-based smart irrigation systems: An overview on the recent trends on sensors and IoT systems for irrigation in precision agriculture. Sensors. 2020; 20(4):1042. https://doi.org/10.3390/s20041042 PMID: 32075172

**98.** Klemeš JJ. Industrial water recycle/reuse. Current opinion in chemical engineering. 2012; 1(3):238–245. https://doi.org/10.1016/j.coche.2012.03.010

**99.** Flörke M, Kynast E, Bärlund I, Eisner S, Wimmer F, Alcamo J. Domestic and industrial water uses of the past 60 years as a mirror of socio-economic development: A global simulation study. Global Environmental Change. 2013; 23(1):144–156. https://doi.org/10.1016/j.gloenvcha.2012.10.018

**100.** Burkhard RA. Learning from architects: the difference between knowledge visualization and information visualization. In: Proceedings. Eighth International Conference on Information Visualisation, 2004. IV 2004.; 2004. p. 519–524.

**101.** Confalonieri R, Coba L, Wagner B, Besold TR. A historical perspective of explainable Artificial Intelligence. WIREs Data Mining and Knowledge Discovery. 2021; 11(1):e1391. https://wires.onlinelibrary.wiley.com/doi/pdf/10.1002/widm.1391.

**102.** Wang J, Wiens J, Lundberg S.: Shapley Flow: A Graph-based Approach to Interpreting Model Predictions.

**103.** Kumar IE, Venkatasubramanian S, Scheidegger C, Friedler S. Problems with Shapley-value-based explanations as feature importance measures. In: III HD, Singh A, editors. Proceedings of the 37th International Conference on Machine Learning. vol. 119 of Proceedings of Machine Learning Research. PMLR; 2020. p. 5491–5500. Available from: https://proceedings.mlr.press/v119/kumar20e.html.

**104.** Callejas Moncaleano DC, Pande S, Rietveld L. Water Use Efficiency: A Review of Contextual and Behavioral Factors. Frontiers in Water. 2021; 3. https://doi.org/10.3389/frwa.2021.685650

**105.** Hoekstra AY, Chapagain AK, Van Oel PR. Advancing Water Footprint Assessment Research: Challenges in Monitoring Progress towards Sustainable Development Goal 6. Water. 2017; 9(6). https://doi.org/10.3390/w9060438

**106.** Hellegers P, van Halsema G. SDG indicator 6.4.1 "change in water use efficiency over time": Methodological flaws and suggestions for improvement. Science of The Total Environment. 2021; 801:149431. https://doi.org/10.1016/j.scitotenv.2021.149431 PMID: 34411791