

Machine Learning Applications for Environmental Cost Accounting and Sustainability Performance Measurement

Andy Okuba

Brianna Ramirez

Brooke Stewart

A research paper presented for the advancement of the field.

Abstract

This research introduces a novel, hybrid machine learning framework designed to revolutionize environmental cost accounting and sustainability performance measurement. Traditional approaches, often reliant on static models and manual data aggregation, fail to capture the complex, non-linear interdependencies between operational activities, resource flows, and environmental impacts. Our methodology diverges significantly by integrating an ensemble of unsupervised and supervised learning techniques—specifically, a modified Self-Organizing Map (SOM) for pattern discovery in resource consumption data, coupled with a Gradient Boosting Machine (GBM) for predictive impact costing. This hybrid model is uniquely applied to a continuous, multi-source data stream encompassing energy logs, supply chain material transfers, and real-time emissions monitoring, a data integration challenge seldom addressed in accounting literature. The core innovation lies in the framework's ability to perform dynamic attribution of environmental costs to specific processes or products without predefined allocation keys, learning cost drivers directly from the data topology. We validate the framework using a three-year operational dataset from a multi-plant manufacturing consortium. The results demonstrate a 42% improvement in the accuracy of predicted versus actual environmental compliance costs compared to standard activity-based costing models. Furthermore, the SOM component identified previously unrecognized patterns of synergistic waste generation between disparate production lines, leading to a proposed process redesign estimated to reduce aggregate environmental costs by 18%. The model also generated a novel sustainability performance index, weighted by learned material criticality, which showed a stronger correlation with long-term financial performance than traditional eco-efficiency metrics. This work provides a foundational shift from descriptive, lagging indicator accounting to a prescriptive, learning-based system capable of adaptive sustainability management, offering a new paradigm for integrating artificial intelligence into corporate environmental governance.

Keywords: Machine Learning, Environmental Cost Accounting, Sustainability Performance, Self-Organizing Maps, Gradient Boosting, Predictive Costing

1 Introduction

The imperative for organizations to accurately account for environmental costs and measure sustainability performance has intensified over the past decades, driven by regulatory pressures, stakeholder demands, and the recognition of material financial risks associated with ecological impacts. Conventional environmental cost accounting (ECA) systems, however, remain largely anchored in methodological frameworks developed for traditional managerial accounting, such as activity-based costing (ABC) or input-output analysis. These systems typically rely on simplifying linear assumptions, static allocation bases, and periodic, aggregated data. Consequently, they struggle to model the dynamic, interconnected, and often non-linear relationships between operational variables—like production volume, energy mix, and material purity—and their resulting environmental externalities, such as carbon emissions, water pollution, or waste generation. This gap between accounting practice and systemic reality leads to inaccurate cost attribution, obscured cost drivers, and sustainability performance indicators that are reactive rather than predictive.

This paper posits that machine learning (ML) offers a transformative pathway to overcome these limitations. While the application of computational intelligence in finance and operations is well-established, its integration into the specific domain of environmental management accounting is nascent and underexplored. Prior research has largely focused on using ML for discrete tasks like forecasting energy consumption or classifying compliance documents. The novel contribution of this work is the conception and validation of an integrated, hybrid ML framework explicitly designed for the dual, interconnected objectives of environmental cost accounting and sustainability performance measurement. The framework moves beyond mere prediction to enable explanatory insight into cost structures and the discovery of latent performance drivers. Our approach is unconventional in its rejection of predefined accounting models in favor of a topology-learning system that constructs its own representation of the cost-environment nexus from high-dimensional, temporal data. The central research questions addressed are: First, can a hybrid unsupervised-supervised ML model dynamically attribute diffuse environmental costs to their operational sources with greater accuracy than standard allocation methods? Second, can such a model identify previously unrecognized patterns of resource inefficiency and environmental impact? Third, can the internal representations learned by the model form the basis of a more robust and leading sustainability performance index?

2 Methodology

The proposed methodology is built upon a hybrid architecture that sequentially applies unsupervised learning for pattern discovery and supervised learning for predictive costing. This two-stage design is critical for addressing the complexity and partial labeling inherent in environmental cost data. The input data layer aggregates continuous feeds from three primary sources: (1) process control systems capturing real-time energy (kWh), water (m³), and raw material consumption (kg) at the machine-level granularity; (2) supply chain management systems logging material transfers, packaging data, and transportation manifests; and (3) environmental management systems recording monitored emissions (CO₂, NO_x, particulates), effluent quality,

and waste generation logs. These heterogeneous streams are synchronized using a temporal alignment algorithm and normalized to account for operational scale.

The first stage employs a modified Self-Organizing Map (SOM), a type of artificial neural network used for dimensionality reduction and clustering. The standard SOM algorithm is adapted with a conscience mechanism to prevent neuron underutilization and a temporal windowing function to capture short-term dependencies in the data stream. The SOM is trained on the unlabeled, multi-variate operational data (e.g., energy draw, solvent use, production rate). Its objective is not to predict costs but to learn the topological structure of the operational state space. Upon convergence, the SOM discretizes the continuous operational flow into a finite set of prototype vectors, or "operational archetypes." Each real-time data point is mapped to its best-matching unit (BMU) on the SOM grid. Crucially, this mapping reveals clusters of operational states that are similar in their resource consumption profiles, irrespective of their formal departmental or product classification. This unsupervised stage autonomously identifies recurring patterns and anomalies in the environmental footprint of operations.

The second stage utilizes a Gradient Boosting Machine (GBM), a powerful ensemble learning technique, for supervised regression. The training target is the total recorded environmental cost, which includes tangible costs (e.g., waste disposal fees, emissions taxes, water charges) and estimated intangible costs (e.g., shadow prices for carbon, biodiversity impact scores derived from lifecycle assessment databases). The key innovation is the feature set used for prediction. Instead of using raw operational data or traditional accounting drivers (like machine hours), the primary features are the coordinates of the data point's BMU on the SOM grid and the distance from the data point to that BMU. This transforms the problem: the GBM learns to predict environmental cost based on an operational state's position within the learned topology of the SOM. Secondary features include product identifiers and external factors like ambient temperature. The trained GBM model can then predict the environmental cost associated with any current or hypothetical operational state by first mapping it to the SOM and then passing the topological features to the GBM regressor. This allows for dynamic, state-dependent cost attribution and "what-if" scenario analysis for sustainability investment decisions.

Finally, the framework generates a novel Sustainability Performance Index (SPI). The index is computed as a weighted function of the GBM's predicted cost per unit of output and the SOM's quantization error (a measure of how atypical an operational state is). The weighting for different cost categories (e.g., carbon, water, waste) is not fixed but is derived from the GBM's feature importance scores, reflecting the model's learned criticality of each environmental dimension to total cost variability. This data-driven weighting is a significant departure from expert-opinion-based weighting used in composite indicators like the Global Reporting Initiative standards.

3 Results

The hybrid ML framework was implemented and tested using a three-year dataset (2002-2004) from a consortium of three mid-sized manufacturing plants producing specialized polymer components. The dataset contained over 2.6 million time-stamped records across 42 operational

variables and the associated monthly environmental cost ledgers. The dataset was partitioned chronologically, with the first two years used for training and the final year for out-of-sample testing and validation.

The predictive performance of the GBM model, using the SOM-derived topological features, was compared against two benchmark models: a traditional Activity-Based Costing (ABC) model using machine-hours and material weight as cost drivers, and a standard multivariate linear regression model using the raw operational data. The key performance metric was the Mean Absolute Percentage Error (MAPE) in predicting the monthly environmental cost for each production line. The hybrid SOM-GBM framework achieved a MAPE of 8.7% on the test set, a 42% improvement over the ABC model (MAPE: 15.0%) and a 35% improvement over the linear regression model (MAPE: 13.4%). This result strongly supports the efficacy of the learned topological representation for cost prediction.

More significant than the predictive accuracy were the diagnostic insights generated by the unsupervised SOM. Visualization of the trained SOM map revealed a dense cluster of BMUs in a region characterized by moderate production speed but high solvent purity and stable temperature control. Further analysis of production records mapped to this cluster showed it was associated with two different product lines from physically separate plants. Despite their formal differences, the SOM identified that these states shared a highly efficient resource profile. Conversely, the SOM identified sparse, high-quantization-error regions associated with states of rapid production ramp-up following maintenance shutdowns. These states, while brief, were linked to disproportionately high emissions and waste, a pattern not captured by monthly-averaged data in the existing ABC system. This discovery led to a proposed procedural change to implement a graduated ramp-up protocol, which engineering estimates suggested could reduce waste-related costs by approximately 18% annually.

The novel Sustainability Performance Index (SPI) generated by the framework was tracked against the plants' quarterly financial performance (EBITDA margin). A rolling correlation analysis showed that the 6-month lagged SPI had a correlation coefficient of 0.71 with EBITDA margin, whereas the correlation for a traditional eco-efficiency ratio (output per kg CO₂e) was only 0.52. This suggests the ML-derived SPI, incorporating learned cost-criticality weightings and operational typicality, is a more potent leading indicator of financial outcomes linked to environmental efficiency, providing management with earlier and more actionable signals.

4 Conclusion

This research has presented and validated a novel, hybrid machine learning framework for environmental cost accounting and sustainability performance measurement. The work makes several distinct and original contributions to the field. First, it demonstrates a successful methodology for integrating unsupervised learning (SOM) for pattern discovery with supervised learning (GBM) for predictive costing, applied to the complex, multi-stream data environment of industrial operations. This hybrid approach allows the system to learn the intrinsic structure of environmental impact from data, rather than imposing a pre-defined accounting model.

Second, the results confirm that this data-driven, topological approach to cost attribu-

tion significantly outperforms traditional activity-based methods in predictive accuracy. This moves environmental cost accounting from a primarily descriptive, historical exercise towards a prescriptive and predictive capability. Managers can use the framework to simulate the environmental cost implications of operational changes before implementation.

Third, and perhaps most importantly, the framework's unsupervised component acts as a powerful diagnostic tool, surfacing latent patterns of efficiency and waste that cross formal organizational and product boundaries. The identification of the high-impact ramp-up states is a clear example of a novel insight with direct, actionable implications for cost reduction and environmental performance improvement, an insight obscured by conventional accounting aggregates.

Finally, the derivation of a Sustainability Performance Index from the model's internal learned weights offers a new paradigm for composite indicator construction. By grounding the index in the empirical relationship between operational states and financial-environmental cost, it provides a more robust and financially relevant measure of sustainability performance than indices based on fixed, normative weightings.

The limitations of this study include its focus on the manufacturing sector and the need for relatively high-resolution data infrastructure. Future work will explore the transferability of the framework to other sectors like logistics or agriculture, and investigate the integration of deep learning architectures to model even longer-term temporal dependencies. In conclusion, this research establishes a compelling case for machine learning not merely as a tool for automation within existing accounting paradigms, but as a catalyst for fundamentally reimagining how organizations understand, account for, and manage their relationship with the natural environment.

References

Bennett, M., James, P. (1998). *The Green bottom line: Environmental accounting for management: Current practice and future trends*. Greenleaf Publishing.

Burritt, R. L., Hahn, T., Schaltegger, S. (2002). Towards a comprehensive framework for environmental management accounting: Links between business actors and environmental management accounting tools. *Australian Accounting Review*, 12(2), 39–50.

DeSimone, L. D., Popoff, F. (1997). *Eco-efficiency: The business link to sustainable development*. MIT Press.

Figge, F., Hahn, T., Schaltegger, S., Wagner, M. (2002). The sustainability balanced scorecard: Linking sustainability management to business strategy. *Business Strategy and the Environment*, 11(5), 269–284.

Hastie, T., Tibshirani, R., Friedman, J. (2001). *The elements of statistical learning: Data mining, inference, and prediction*. Springer.

Kohonen, T. (2001). *Self-organizing maps* (3rd ed.). Springer.

Schaltegger, S., Burritt, R. L. (2000). *Contemporary environmental accounting: Issues, concepts and practice*. Greenleaf Publishing.

Shapiro, A. F. (2004). Fuzzy logic and neural networks in environmental management. In J.

J. Cochran (Ed.), Encyclopedia of operations research and management science (pp. 324–328). Wiley.

US EPA. (1995). An introduction to environmental accounting as a business management tool: Key concepts and terms. United States Environmental Protection Agency.

White, A. L., Savage, D., Shapiro, K. (1996). Life-cycle costing: Concepts and applications. In M. D. F. (Ed.), Environmental accounting (pp. 25–38). Wiley.