

Predictive Modeling of Environmental Fines Using Historical Accounting Data

Elijah Rivera

Ella Adams

Emily Coleman

An original research paper submitted for peer review.

Abstract

This paper introduces a novel methodological framework for predicting corporate environmental fines by leveraging historical accounting data, a previously underexplored data source in environmental compliance forecasting. Traditional approaches to environmental risk assessment have relied heavily on direct environmental metrics, regulatory history, and industry-specific factors, often overlooking the rich predictive signals embedded in financial statements. We propose that patterns in accounting data—including expense allocations, capital expenditure trends, depreciation methods, and footnote disclosures—contain latent indicators of environmental management priorities and potential compliance vulnerabilities. Our research formulates the prediction of environmental fines as a time-series classification problem, employing a hybrid model that combines autoregressive integrated moving average (ARIMA) components for capturing temporal financial trends with a multilayer perceptron for capturing non-linear interactions between accounting variables. We construct a unique dataset spanning 1998 to 2004, linking the financial statements of manufacturing firms from Compustat with environmental penalty records from the EPA’s Enforcement and Compliance History Online (ECHO) database. The model identifies several key accounting predictors, most notably the ratio of repair and maintenance expenses to capital expenditures, the volatility of cost of goods sold, and specific linguistic cues in management discussion and analysis (MDA) sections related to environmental contingencies. Results demonstrate a predictive accuracy of 82.7% in classifying firms that will incur a significant environmental fine within the next fiscal year, a substantial improvement over baseline models using only environmental performance indicators. This work establishes a new, cross-disciplinary research direction at the intersection of environmental informatics, forensic accounting, and predictive analytics, offering regulators and investors a proactive tool for risk assessment derived from routinely published financial information.

Keywords: environmental fines, predictive modeling, accounting data, forensic analytics, regulatory compliance, hybrid ARIMA-MLP model

1 Introduction

The prediction of corporate environmental misconduct and subsequent regulatory penalties represents a significant challenge for stakeholders including investors, regulators, and the firms themselves. Conventional predictive models in this domain have predominantly relied on historical environmental performance data, such as past violations, emissions reports, and permit compliance records. While informative, this approach suffers from a fundamental limitation: it is inherently reactive, relying on data generated by the very regulatory system it seeks to predict. This paper proposes a paradigm shift by investigating the predictive power of historical accounting data—a rich, standardized, and publicly available information source—for forecasting environmental fines. The core hypothesis is that a firm’s financial decisions and reporting practices embed subtle signals regarding its operational priorities, risk management ethos, and, by extension, its propensity for environmental compliance failures.

Accounting data offers a unique lens. Choices in capitalizing versus expensing environmental upgrades, the level of detail in contingent liability disclosures, trends in maintenance expenditures, and the linguistic tone of management discussions regarding environmental matters may all serve as proxies for a firm’s environmental stewardship culture. For instance, a firm that consistently expenses minor environmental repairs rather than capitalizing them might be signaling a short-term operational focus that could correlate with higher compliance risks. Similarly, increasing volatility in production costs might indicate process instability, a potential precursor to environmental incidents. This research seeks to formalize these intuitions into a quantifiable, predictive model.

Our work is distinguished by its cross-disciplinary novelty, bridging the fields of environmental science, accounting, and machine learning. We move beyond correlation studies that link poor environmental performance to financial outcomes and instead invert the causal inquiry: can financial data predict environmental outcomes? The primary research questions guiding this study are: (1) Which specific accounting variables and ratios demonstrate sta-

tistically significant predictive power for future environmental fines? (2) Can a hybrid time-series and pattern recognition model effectively integrate these financial signals to achieve robust classification accuracy? (3) Does the predictive power of accounting data vary across industry sub-sectors within the manufacturing domain? By answering these questions, this paper contributes a novel methodological framework and provides empirical evidence for a previously untapped source of predictive intelligence in environmental governance.

2 Methodology

The methodological framework of this study is built upon a hybrid modeling architecture designed to capture both the temporal dynamics of financial time-series and the complex, non-linear interactions between disparate accounting variables. The core innovation lies in the feature engineering process, which translates raw accounting data into predictive signals, and the model design, which synergistically combines statistical and connectionist approaches.

2.1 Data Collection and Preprocessing

A longitudinal dataset was constructed for the period 1998 to 2004. Firm-level financial data was extracted from the Standard & Poor's Compustat database, focusing on North American Industrial (NI) and Industrial Supplemental (SI) files. Environmental fine data was sourced from the U.S. Environmental Protection Agency's Enforcement and Compliance History Online (ECHO) database, which records administrative and judicial penalties under major environmental statutes. The datasets were linked via firm name and location, resulting in a matched panel of 542 manufacturing firms (SIC codes 2000-3999). The target variable was binary: whether a firm incurred an environmental fine exceeding \$50,000 in a given fiscal year (t+1).

From the raw accounting data, 47 initial features were engineered. These included tradi-

tional ratios (e.g., current ratio, debt-to-equity), industry-specific metrics (e.g., cost of goods sold to sales), and novel constructs hypothesized to be environmentally relevant. Key novel features included: the Environmental Capex Ratio (capital expenditures tagged as environmental in footnotes divided by total capex), Maintenance Intensity (repair and maintenance expenses over property, plant, and equipment net), Contingency Disclosure Score (a text-analysis-derived measure of the specificity and volume of environmental liability discussions in MD&A and footnotes), and Earnings Smoothing Index (a measure of discretionary accruals). All financial variables were scaled by total assets to control for firm size and winsorized at the 1st and 99th percentiles to mitigate outlier effects.

2.2 Hybrid Model Architecture

The predictive model, termed the ARIMA-MLP Hybrid, consists of two integrated components. For each firm and each continuous accounting variable (e.g., maintenance intensity over time), a firm-specific ARIMA model is fitted to the time-series from years $t-4$ to t . The parameters of this ARIMA model (order p, d, q) and the one-step-ahead forecast error for year t become new input features. This process transforms temporal patterns into static, cross-sectional descriptors.

These derived time-series features are then concatenated with the other static accounting ratios and text-based scores for year t . This combined feature vector serves as the input to a Multilayer Perceptron (MLP). The MLP architecture comprised an input layer (size equal to the feature vector), two hidden layers with hyperbolic tangent activation functions (32 and 16 neurons, respectively), and a single output neuron with a sigmoid activation function for binary classification. The model was trained to minimize binary cross-entropy loss using backpropagation with a learning rate scheduler.

2.3 Training and Evaluation Protocol

The dataset was partitioned temporally: data from 1998-2002 was used for training and validation (using 5-fold cross-validation), and data from 2003-2004 was held out as a strict test set to evaluate real-world predictive performance. This ensures the model is evaluated on its ability to predict future fines based on past data, simulating a real deployment scenario. Performance was assessed using accuracy, precision, recall, F1-score, and the area under the Receiver Operating Characteristic curve (AUC-ROC). Baseline models, including a logistic regression on environmental history variables only and a standard MLP on accounting data without ARIMA features, were implemented for comparison.

3 Results

The experimental results provide strong support for the core hypothesis that historical accounting data contains significant predictive signals for future environmental fines. The hybrid ARIMA-MLP model achieved a test set accuracy of 82.7% and an AUC-ROC of 0.891, substantially outperforming both baseline models. The logistic regression model using only prior environmental violations and industry dummies achieved an accuracy of 68.2% (AUC-ROC: 0.723), while the standard MLP (without ARIMA features) achieved 76.4% accuracy (AUC-ROC: 0.832). This performance delta highlights the value added by both the novel accounting features and the hybrid modeling approach that captures temporal dynamics.

Feature importance analysis, conducted via permutation importance and analysis of the MLP’s first-layer weights, identified the most potent predictors. The top three features were: (1) the forecast error from the ARIMA model fitted to the Maintenance Intensity time-series (a large positive error, indicating an unexpected spike in maintenance spending, was predictive of fines), (2) the Contingency Disclosure Score (firms with more vague or minimal disclosures were at higher risk), and (3) the volatility of the Cost of Goods Sold ratio over the previous four years. Interestingly, traditional profitability measures like return

on assets showed negligible predictive power.

The model’s precision (the proportion of predicted fines that were correct) was 79.5%, and its recall (the proportion of actual fines that were predicted) was 73.8%. This indicates a slightly conservative bias, erring on the side of missing some fines rather than generating excessive false alarms—a desirable characteristic for a regulatory screening tool. Performance was consistent across the chemical (SIC 28), primary metal (SIC 33), and fabricated metal (SIC 34) industries, but slightly weaker for machinery (SIC 35), suggesting the financial signals of environmental risk may be somewhat industry-contextual.

A case study analysis of several firms that incurred large fines in 2004 revealed compelling narratives. One chemical manufacturer showed a three-year declining trend in its Environmental Capex Ratio alongside a sharp, ARIMA-model-predicted error in its maintenance expenses in the year prior to the fine, coinciding with a major wastewater treatment violation. The model successfully flagged this firm as high-risk. These results empirically validate the proposed link between specific financial decision-making patterns and subsequent environmental compliance failures.

4 Conclusion

This research has demonstrated the feasibility and efficacy of a novel approach to predicting environmental regulatory fines using historical accounting data. By conceptualizing financial statements as a latent data source for environmental risk assessment, we have opened a new cross-disciplinary research avenue. The hybrid ARIMA-MLP model successfully extracted and synthesized predictive signals from the temporal evolution of accounting ratios and the qualitative content of financial disclosures, achieving robust classification performance.

The primary original contributions of this work are threefold. First, it establishes a new predictive paradigm that is proactive rather than reactive, leveraging data generated independently of the regulatory enforcement cycle. Second, it introduces and validates a set of

novel accounting-based features, such as the Maintenance Intensity trend and the Contingency Disclosure Score, which serve as quantifiable proxies for environmental management quality. Third, it provides a scalable methodological framework—the hybrid ARIMA-MLP architecture—for integrating time-series analysis with pattern recognition in the context of socio-economic prediction problems.

The implications are significant for multiple stakeholders. Regulators could employ such a model to prioritize inspections and allocate monitoring resources more efficiently. Investors and financial analysts could integrate environmental fine risk into their valuation and due diligence models. Firms themselves could use this as an internal diagnostic tool to identify financial patterns that may signal underlying operational risks. Future work will focus on expanding the feature set to include textual analysis of annual report narratives beyond the MD&A, incorporating intra-industry network effects, and testing the model’s generalizability to other regulatory domains, such as workplace safety or financial misconduct. This study underscores the profound insights that can be gleaned from re-examining conventional data sources through an unconventional, interdisciplinary lens.

References

Beaver, W. H. (1966). Financial ratios as predictors of failure. *Journal of Accounting Research*, 4, 71–111.

Blacconiere, W. G., Patten, D. M. (1994). Environmental disclosures, regulatory costs, and changes in firm value. *Journal of Accounting and Economics*, 18(3), 357–377.

Box, G. E. P., Jenkins, G. M., Reinsel, G. C. (1994). *Time series analysis: Forecasting and control* (3rd ed.). Prentice Hall.

Cormier, D., Magnan, M. (1997). Investors’ assessment of implicit environmental liabilities: An empirical investigation. *Journal of Accounting and Public Policy*, 16(2), 215–241.

Dechow, P. M., Sloan, R. G., Sweeney, A. P. (1995). Detecting earnings management.

The Accounting Review, 70(2), 193–225.

Hamilton, J. T. (1995). Pollution as news: Media and stock market reactions to the Toxics Release Inventory data. *Journal of Environmental Economics and Management*, 28(1), 98–113.

Haykin, S. (1999). *Neural networks: A comprehensive foundation* (2nd ed.). Prentice Hall.

Konar, S., Cohen, M. A. (2001). Does the market value environmental performance? *The Review of Economics and Statistics*, 83(2), 281–289.

Rumelhart, D. E., Hinton, G. E., Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088), 533–536.

Shane, P. B., Spicer, B. H. (1983). Market response to environmental information produced outside the firm. *The Accounting Review*, 58(3), 521–538.