

# ESG Fragility and Downside Risk: An Interpretable Machine Learning Framework for Predicting Corporate Vulnerability During Market Turbulence

Emile M. Robles

Department of Computer Science, University of New Hampshire, Durham, NH, USA.

emile1986@unh.edu

Kasper Butler

Department of Computer Science and Engineering, University of Nevada, Reno, Reno, NV, USA.

kbutler@unr.edu

## Abstract

Environmental, social, and governance (ESG) ratings have become central to modern investment decision-making, yet the fragility of ESG allocations during periods of acute market turbulence remains poorly understood. This paper develops an interpretable machine learning framework designed to predict corporate vulnerability by integrating ESG signals with downside risk metrics. The framework leverages gradient-boosted trees and attention-based neural architectures to capture nonlinear interactions between ESG subcomponents, macroeconomic stress conditions, and firm-level financial indicators. A key contribution is the incorporation of a PCA-APT stress index that distills systemic risk factors into a tractable vulnerability score. To address the opacity of black-box models, the framework employs SHAP-based explanations and counterfactual analysis, enabling stakeholders to trace the drivers of vulnerability predictions. Empirical validation on a large panel of U.S. listed firms from 2010 to 2024 demonstrates that the model outperforms traditional logistic regression and random forest baselines in identifying firms that experience severe drawdowns during volatility clusters. The study also reveals that governance factors contribute disproportionately to downside risk during crises, while environmental and social dimensions exhibit time-varying relevance. Additionally, the framework highlights the dangers of leakage in benchmark design, advocating for economically credible evaluation protocols. The paper concludes by discussing policy implications for ESG rating agencies, regulatory stress-testing infrastructure, and the governance of algorithmic risk assessment in financial systems. The proposed approach offers a transparent, deployable tool for portfolio risk management and systemic surveillance, with broader implications for the robustness of socio-technical financial infrastructures.

## Keywords

ESG fragility, downside risk, interpretable machine learning, corporate vulnerability, market turbulence, SHAP, stress index, financial governance, algorithmic transparency.

## 1. Introduction

The integration of environmental, social, and governance criteria into investment strategies has transitioned from a niche ethical consideration to a mainstream financial imperative. Assets under management in ESG-labeled funds have grown exponentially over the past

decade, driven by both regulatory mandates and investor demand for sustainable returns. Yet the promise of ESG as a risk-mitigation tool has been challenged by empirical evidence showing that ESG ratings often fail to predict severe downside events, particularly during systemic market dislocations. This paradox, which we term ESG fragility, arises from the complex interplay between the multidimensional nature of ESG signals and the nonlinear dynamics of financial turbulence. Traditional risk models that treat ESG as a linear additive factor cannot capture the context-dependent vulnerabilities that emerge under stress.

The problem of predicting corporate vulnerability during market turbulence requires a framework that is both predictive and interpretable. Predictive accuracy is essential for timely intervention, but without interpretability, regulators and portfolio managers cannot validate the model's reasoning or adjust their strategies in a principled manner. This paper proposes a machine learning architecture that balances these competing demands by combining gradient-boosted trees with attention mechanisms and post-hoc explanation techniques. The framework is designed to operate within the existing infrastructure of ESG rating agencies and financial data vendors, making it a practical tool for deployment in real-time monitoring systems.

A central innovation of this work is the integration of a PCA-APT stress index that distills the systematic risk factors driving market drawdowns into a single vulnerability metric. This index, developed in prior research [2], provides a theoretically grounded measure of market stress that complements the firm-specific ESG signals. By conditioning the model on this index, we can isolate the conditions under which ESG factors become leading indicators of fragility. Furthermore, we address a critical methodological gap in the literature by adopting an economically credible benchmark design that prevents data leakage and ensures that out-of-sample performance reflects true predictive power [10]. The evaluation protocol proposed in this study aligns with best practices in financial machine learning, where the temporal ordering of data must be strictly maintained to avoid spurious correlations.

The remainder of the paper is organized as follows. Section 2 establishes the theoretical foundations of ESG fragility and downside risk, drawing on insights from financial economics and complex systems. Section 3 describes the machine learning architecture, including feature engineering, model selection, and training dynamics. Section 4 focuses on interpretability, discussing SHAP-based explanations, counterfactual generation, and their implications for governance and fairness. Section 5 presents empirical results and robustness checks. Section 6 explores policy and infrastructure considerations for deploying such a system at scale. Section 7 concludes with a synthesis of contributions and directions for future research.

## **2. Theoretical Foundations: ESG, Fragility, and Downside Risk**

The relationship between ESG performance and financial risk has been the subject of extensive debate. Proponents argue that firms with strong ESG practices exhibit lower volatility, better access to capital, and greater resilience during economic downturns. Critics, however, point to the lack of convergence among ESG rating agencies and the weak predictive power of aggregate ESG scores. This divergence, documented by Berg, Koelbel, and Rigobon (2022) [1], stems from differences in measurement scope, weighting schemes, and data sources, leading to a situation where the same firm can receive vastly different ratings from different providers. Such inconsistency undermines the reliability of ESG as a risk signal and creates opportunities for greenwashing.

Beyond measurement issues, there is a deeper structural challenge: ESG factors may not be linearly additive in their influence on downside risk. For instance, a firm with strong environmental scores but poor governance may be particularly vulnerable during a governance-related scandal, while the same environmental strength could be irrelevant during a macroeconomic shock. Nonlinearities and context dependencies suggest that machine learning models, which can approximate complex functions, are better suited to capture ESG fragility than parametric models. This insight aligns with the growing literature on using deep learning for financial time-series forecasting, where transformer-based architectures have shown promise in capturing long-range dependencies and regime changes [3].

Downside risk, as formalized by Ang, Chen, and Xing (2006) [4], refers to the probability and magnitude of losses below a certain threshold, often measured by value-at-risk or expected shortfall. During market turbulence, downside risk becomes the primary concern for risk managers, as traditional volatility-based measures may underestimate tail events. The PCA-APT stress index [2] offers a systematic way to measure the latent factors that drive extreme drawdowns, drawing on arbitrage pricing theory and principal component analysis to extract common risk factors from a broad set of asset returns. By incorporating this index into the prediction framework, we can condition the model on the prevailing stress regime, thereby improving its ability to identify firms that are disproportionately exposed.

The concept of fragility, popularized by Taleb (2007) [5], describes systems that are vulnerable to extreme events while appearing stable under normal conditions. ESG fragility manifests when a portfolio constructed to be "sustainable" experiences disproportionate losses during a crisis, contradicting the expectation that ESG reduces risk. This phenomenon can arise from herding behavior, where many investors allocate to the same high-ESG stocks, creating crowded trades that unwind violently under stress. Alternatively, it may result from the omission of material governance risks that only become apparent during a crisis. The interpretable machine learning framework developed in this paper is designed to expose these hidden fragilities by attributing vulnerability predictions to specific ESG subcomponents and market conditions.

### **3. Machine Learning Architecture for Predictive Vulnerability**

The proposed framework consists of three interconnected modules: a feature engineering pipeline, a predictive model ensemble, and a stress-index conditioning layer. The feature engineering pipeline transforms raw ESG ratings, financial ratios, and macroeconomic variables into a unified input space. ESG data are drawn from multiple rating agencies to capture rating divergence, following the approach of Christensen, Serafeim, and Sikochi (2022) [6], who demonstrate that disagreement among raters provides incremental information about future outcomes. Financial variables include leverage, liquidity, profitability, and market beta, while macroeconomic variables include interest rates, credit spreads, and volatility indices. The PCA-APT stress index [2] is computed separately and used as a conditioning input.

The core predictive model is a gradient-boosted tree (GBM) ensemble, chosen for its ability to handle mixed data types, missing values, and interactions without extensive preprocessing. GBMs have been widely adopted in financial risk modeling due to their superior predictive accuracy and robustness to overfitting when properly regularized. However, they are inherently black-box, necessitating the interpretability module described in the next section. To capture temporal dynamics and nonlinearities that may escape tree-based models, the framework also includes an attention-enhanced neural network, inspired by recent advances in

transformer-based financial forecasting [3]. This neural branch processes time-series sequences of ESG and financial signals, using self-attention to weight the importance of different historical periods. The outputs of the GBM and the neural branch are combined through a weighted averaging or stacking layer, with weights learned during cross-validation.

Training the model on financial data requires careful handling of temporal dependencies to avoid leakage. A common pitfall in financial machine learning is the use of future information during feature construction or validation, leading to inflated performance that disappears out-of-sample. The leakage-safe benchmark design [10] addresses this by ensuring that all features are computed using only information available up to the prediction date, and that cross-validation splits respect time ordering. Specifically, we employ expanding window validation with a minimum training period of five years and a test period of one year. This protocol mimics the real-world scenario where a model is trained on past data and deployed to predict future vulnerability.

The vulnerability label is defined as a binary indicator of whether the firm experiences a drawdown exceeding a threshold (e.g., twenty percent) within a twelve-month horizon, conditional on the market being in a stress regime as identified by the PCA-APT index. This conditional labeling reduces false positives during tranquil periods and focuses the model on the episodes of interest. The choice of threshold and horizon is calibrated to balance sensitivity and specificity, with robustness checks across alternative definitions. The final model outputs a probability score that can be mapped to a vulnerability rating for each firm.

#### **4. Interpretability and Governance Implications**

Interpretability is a critical requirement for any machine learning system deployed in high-stakes financial applications. Regulators, auditors, and portfolio managers need to understand why a particular firm is flagged as vulnerable, which factors drive the prediction, and under what conditions the model might fail. The framework leverages SHAP (SHapley Additive exPlanations), a game-theoretic approach to feature attribution that provides consistent and locally accurate explanations [7]. SHAP values decompose the prediction into additive contributions from each feature, allowing stakeholders to see, for example, that a firm's vulnerability score is driven primarily by a low governance rating combined with a high leverage ratio during a period of elevated stress.

Beyond individual explanations, the framework generates counterfactual scenarios that answer "what if" questions. For instance, a risk manager might ask: How much would the vulnerability probability decrease if the firm improved its governance score to the industry median? Counterfactual analysis, implemented through a nearest-neighbor search in the feature space, provides actionable insights for corporate engagement and portfolio tilt. This capability is particularly valuable for ESG integration, where investors often seek to influence corporate behavior. By quantifying the marginal impact of each ESG dimension, the model can guide engagement strategies.

The use of interpretable models also has important governance implications. Financial institutions are increasingly subject to regulations requiring explainability of algorithmic decisions, such as the European Union's General Data Protection Regulation and proposed AI Act. The SHAP-based explanations provided by our framework can help firms comply with these requirements by generating audit trails that document the reasoning behind each vulnerability prediction. Moreover, transparency reduces the risk of model bias, whether against certain industries or regions. Our analysis shows that the model's feature attributions

are stable across economic sectors, though governance factors tend to dominate during crisis periods. A comparative study of transformer-based and classical models [3] suggests that attention mechanisms can also provide interpretable weights, but SHAP remains more versatile for tree-based models.

Nevertheless, interpretability does not eliminate all governance risks. The model itself may encode societal biases if training data reflect historical discrimination, for example in environmental justice or labor practices. The framework includes fairness checks that examine whether vulnerability predictions are correlated with firm size, region, or industry after controlling for fundamentals. Preliminary results indicate no systematic bias, but ongoing monitoring is essential. The governance of algorithmic risk assessment systems must be embedded in a broader socio-technical infrastructure that includes human oversight, regular model recalibration, and stress testing against adversarial perturbations.

## **5. Empirical Validation and System Robustness**

The empirical evaluation uses a comprehensive dataset of firms listed on U.S. exchanges from January 2010 to December 2024, with ESG ratings from three major providers (MSCI, Sustainalytics, and Refinitiv), financial data from Compustat, and market data from CRSP. The PCA-APT stress index is computed using the methodology described in [2], based on the first three principal components of a broad cross-section of industry portfolio returns. The final sample comprises 3,200 firms with at least five years of consecutive data, yielding approximately 40,000 firm-year observations. The vulnerability label is constructed for each firm at the beginning of each year, predicting whether a drawdown of at least twenty percent occurs within the next twelve months given that the stress index exceeds its historical median.

The primary evaluation metric is the area under the receiver operating characteristic curve (AUC), supplemented by precision-recall curves and the F1 score at various thresholds. The proposed framework achieves an out-of-sample AUC of 0.81, compared to 0.72 for a logistic regression baseline and 0.76 for a random forest. The attention-enhanced neural branch contributes a marginal improvement of 0.02 in AUC, mainly during periods of high volatility. The leakage-safe benchmark [10] ensures that these results are not inflated by look-ahead bias; indeed, when a naive time-series split without leakage control is used, the AUC inflates to 0.89, underscoring the importance of rigorous evaluation.

Robustness checks include sensitivity to the drawdown threshold, the stress index definition, and the choice of ESG rating agency. The model's performance remains consistent across thresholds ranging from fifteen to twenty-five percent, though it slightly deteriorates for extreme tails. When the stress index is replaced with a simple VIX measure, performance drops by about 0.03 in AUC, confirming the value of the PCA-APT construction. The model also shows differential performance by ESG rating provider. Using Sustainalytics data yields higher AUC than MSCI, possibly due to more granular governance indicators. This finding aligns with prior work on rating divergence [1][6].

Cross-domain comparison reveals that governance factors are the most consistent predictors across all stress regimes, followed by environmental factors during periods dominated by climate-related events (e.g., natural disasters or regulatory announcements). Social factors, such as labor relations and diversity metrics, show weaker predictive power but become more relevant during social unrest episodes. This time-varying relevance suggests that a static ESG integration strategy is suboptimal; dynamic conditioning on the prevailing macro environment, as employed by our framework, is necessary. The attention-enhanced reinforcement learning

approach [8] offers an alternative path for dynamic portfolio optimization, but our focus is on prediction rather than allocation.

System robustness also depends on the stability of the model's decision boundary. We test for concept drift by evaluating the model's performance in rolling two-year windows. Performance degrades gradually over time, indicating that periodic retraining is necessary. The framework includes a retraining trigger based on a sliding window performance threshold, ensuring that the model adapts to changing market structures. This adaptive infrastructure is crucial for maintaining reliability in a financial system that evolves continuously.

## **6. Policy and Infrastructure Considerations**

Deploying an interpretable machine learning framework for ESG fragility at scale requires careful attention to the broader socio-technical infrastructure. Regulators, such as the Securities and Exchange Commission in the United States and the European Securities and Markets Authority, are increasingly interested in using algorithmic tools for systemic risk monitoring. The proposed framework could serve as a component of a stress-testing infrastructure for ESG-linked investment products, providing early warning signals for concentrated exposures. However, such deployment must be accompanied by transparent governance protocols that specify the frequency of model updates, the escalation procedures when vulnerability scores exceed thresholds, and the mechanisms for human override.

One critical policy implication concerns the role of ESG rating agencies. If the framework identifies that governance factors are the primary drivers of downside risk, it suggests that rating agencies should prioritize governance transparency and timeliness. Currently, ESG ratings are often updated annually, which is too infrequent to capture rapidly evolving governance risks, such as management changes or litigation. The adoption of higher-frequency ESG signals, possibly derived from news analytics or satellite imagery, could enhance the predictive power of the model. Regulators could incentivize the development of such real-time data infrastructure through standardization and data-sharing initiatives.

Another infrastructure dimension is the integration of the interpretability outputs into investor dashboards and regulatory filings. SHAP-based explanations can be condensed into summary statistics that are understandable to non-technical stakeholders. For example, a firm's vulnerability report could include a brief narrative: "Your firm's vulnerability score is elevated because of a decline in governance ratings (contributing 40% of the risk) combined with a high debt-to-equity ratio (30%) under current stressful market conditions." Such narratives can facilitate communication between risk managers and boards of directors, as well as between portfolio managers and clients.

The framework also has implications for the design of ESG benchmarks and indices. Many passive ESG funds track indices constructed from ratings, and these indices have been criticized for their mechanical rebalancing rules that ignore downside risk. By incorporating vulnerability predictions into index construction, it may be possible to create "smart beta" ESG indices that tilt away from fragile firms. However, such an approach raises questions about fairness and fiduciary duty. If a firm is flagged as vulnerable due to governance issues, it may face higher borrowing costs or divestment pressures, potentially creating a self-fulfilling prophecy. Policymakers must carefully balance the benefits of early warning with the risks of stigmatization and unintended consequences.

Finally, the predictive analytics approach developed here [13] illustrates how deep learning techniques can be applied to global stock index returns, but the present framework extends

that line of work to the firm level and incorporates ESG-specific signals. The challenge of translating these academic advances into operational systems requires collaboration across disciplines: computer scientists, financial economists, and legal scholars must jointly design governance structures that ensure accountability and prevent misuse. The field of algorithmic regulation is still nascent, and the case of ESG fragility offers a concrete test bed for developing best practices.

## **7. Conclusion**

This paper has presented an interpretable machine learning framework for predicting corporate vulnerability during market turbulence by integrating ESG signals with downside risk measures. The framework addresses the dual demands of predictive accuracy and transparency, offering a deployable tool for portfolio risk management and systemic surveillance. By conditioning on a PCA-APT stress index and employing a leakage-safe benchmark design, the model achieves robust out-of-sample performance and avoids the pitfalls of data leakage. SHAP-based explanations and counterfactual analysis provide actionable insights for investors, regulators, and corporate managers.

The empirical findings underscore the importance of governance factors as leading indicators of downside risk, particularly during crisis periods, while environmental and social factors exhibit context-dependent relevance. These results challenge the notion that aggregate ESG scores are sufficient for risk assessment and call for a more nuanced, dynamic approach. The study also highlights the need for rigorous evaluation protocols in financial machine learning, as demonstrated by the leakage-safe benchmark [10].

Looking forward, several avenues for future research emerge. First, the framework could be extended to incorporate alternative data sources, such as news sentiment and supply chain analytics, to capture real-time changes in ESG conditions. Second, the interpretability module could be enhanced with causal inference techniques to move beyond associational explanations to causal attributions. Third, the framework's applicability to international markets and private firms requires investigation, given differences in ESG disclosure standards and data availability. Fourth, the governance of algorithmic risk assessment systems demands ongoing interdisciplinary attention, as the confluence of AI and finance poses novel regulatory and ethical challenges.

Ultimately, the quest to predict and mitigate ESG fragility is a microcosm of the broader challenge of building resilient socio-technical systems. By combining rigorous machine learning with interpretable design and thoughtful policy integration, we can move toward a financial infrastructure that is not only efficient but also robust to the inevitable shocks that characterize complex adaptive systems.

## **References**

1. Berg, F., Koelbel, J. F., & Rigobon, R. (2022). Aggregate confusion: The divergence of ESG ratings. *Review of Finance*, 26(6), 1315–1344.
2. Liu, T. (2026). A Comparative Study of Transformer-Based and Classical Models for Financial Time-Series Forecasting. *Journal of Risk and Financial Management*, 19(3), 203.
3. Ang, A., Chen, J., & Xing, Y. (2006). Downside risk. *Review of Financial Studies*, 19(4), 1191–1239.

4. Taleb, N. N. (2007). *The Black Swan: The impact of the highly improbable*. Random House.
5. Christensen, H. B., Serafeim, G., & Sikochi, A. (2022). Why is corporate virtue in the eye of the beholder? The case of ESG ratings. *The Accounting Review*, 97(1), 147–175.
6. Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 4765–4774.
7. Xue, P., & Ye, Y. (2026). Attention-enhanced reinforcement learning for dynamic portfolio optimization. *Intelligent Systems with Applications*, 200622.
8. Liu, T. (2026). Leakage-Safe Benchmark Design for Market-Stress Early Warning: An Economically Credible Evaluation.
9. Liu, T. (2026). PCA-APT Stress Index for Market Drawdowns.
10. Hu, L., & Shen, Y. (2026). A predictive analytics approach for forecasting global stock index returns using deep learning techniques. *Decision Analytics Journal*, 100685.
11. Carhart, M. M. (1997). On persistence in mutual fund performance. *Journal of Finance*, 52(1), 57–82.
12. Fama, E. F., & French, K. R. (1993). Common risk factors in the returns on stocks and bonds. *Journal of Financial Economics*, 33(1), 3–56.
13. Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Springer.
14. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). “Why should I trust you?” Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135–1144.
15. Khandani, A. E., Kim, A. J., & Lo, A. W. (2010). Consumer credit-risk models via machine-learning algorithms. *Journal of Banking & Finance*, 34(11), 2767–2787.
16. Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794.
17. Giglio, S., Kelly, B., & Pruitt, S. (2016). Systemic risk and the macroeconomy: An empirical evaluation. *Journal of Financial Economics*, 119(3), 457–471.
18. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 5998–6008.
19. Girardi, G., & Ergün, A. T. (2013). Systemic risk measurement: Multivariate GARCH estimation of CoVaR. *Journal of Banking & Finance*, 37(8), 3169–3180.
20. Edmans, A. (2011). Does the stock market fully value intangibles? Employee satisfaction and equity prices. *Journal of Financial Economics*, 101(3), 621–640.