# Predictive Models for Environmental Disclosure Quality and Market Reactions

*Lucas Morris*

*Lucy Bennett*

*Madeline Cooper*

A novel computational linguistics and bio-inspired optimization approach to ESG analysis

**Abstract**

This research introduces a novel, cross-disciplinary framework that applies computational linguistics and machine learning techniques, traditionally used in information retrieval and sentiment analysis, to the domain of corporate environmental disclosure. We depart from conventional econometric models by proposing a hybrid methodology that integrates Latent Dirichlet Allocation (LDA) for thematic decomposition of sustainability reports with a bio-inspired optimization algorithm—specifically, a modified Ant Colony Optimization (ACO) technique—to select predictive features for market reaction modeling. The core innovation lies in reformulating the problem of assessing disclosure quality not as a static classification task but as a dynamic, multi-dimensional signal extraction challenge, where the 'quality' is inferred from the semantic coherence, specificity, and forward-looking content of disclosures relative to industry-specific environmental materiality thresholds. Our model processes unstructured textual data from corporate environmental reports and regulatory filings to generate a continuous Disclosure Quality Index (DQI). We then examine the market's reaction to this index, hypothesizing that investors process nuanced qualitative information differently than quantitative metrics. Using a unique dataset compiled from 2000 to 2004, our results demonstrate that the DQI, derived from our hybrid LDA-ACO model, has superior predictive power for abnormal stock returns around disclosure events compared to traditional metrics based on word counts or binary compliance checks. Furthermore, we identify a non-linear, threshold-based market reaction, suggesting that investors discount disclosures until a certain level of specificity and coherence is achieved, a finding with significant implications for both corporate reporting strategies and regulatory policy. This work establishes a new paradigm for computational analysis in environmental, social, and governance (ESG) research, moving beyond keyword spotting towards a sophisticated understanding of informational substance.

**Keywords:** Environmental Disclosure, Computational Linguistics, Ant Colony Optimization, Market Reaction, Latent Dirichlet Allocation, ESG, Textual Analysis

# 1 Introduction

The landscape of corporate environmental disclosure has evolved from a peripheral public relations activity to a central component of investor communication and risk assessment. Traditional financial models, however, struggle to quantify the qualitative substance of these disclosures, often relying on binary indicators of presence or simple keyword frequencies. This research posits that the market's reaction to environmental information is not solely a function of its ex-

istence but is critically dependent on its qualitative attributes: specificity, thematic coherence, relevance to material issues, and forward-looking orientation. To capture these nuances, we propose a radical departure from established methods by constructing a predictive model grounded in computational linguistics and bio-inspired optimization—a synergy previously unexplored in the accounting and finance literature.

Our investigation is guided by two primary research questions that have not been extensively covered in a unified framework. First, can a hybrid model combining thematic topic modeling (LDA) with a feature selection mechanism inspired by swarm intelligence (ACO) reliably construct a continuous index of environmental disclosure quality from unstructured text? Second, does this computationally-derived quality index explain cross-sectional variations in market reactions (abnormal returns) better than conventional, reductionist measures, and is this relationship linear or governed by cognitive thresholds? By addressing these questions, we contribute to multiple fields: we advance methodological innovation in textual analysis for social science, provide new tools for investors and regulators to evaluate corporate transparency, and offer a novel theoretical perspective on how complex qualitative signals are processed in financial markets.

The remainder of this paper is structured as follows. The Methodology section details our innovative hybrid LDA-ACO framework and the construction of the Disclosure Quality Index (DQI). The Results section presents the empirical findings from applying our model to a unique dataset of corporate environmental reports from the early 2000s, comparing its predictive power against benchmarks. Finally, the Conclusion discusses the implications of our original contributions for theory, practice, and future research at the intersection of computer science, finance, and environmental policy.

## 2 Methodology

Our methodology is built upon a novel, two-stage hybrid framework designed to extract meaningful signals from the high-dimensional, noisy textual data of environmental disclosures and to model the market's reaction to these signals.

## 2.1 Data Collection and Preprocessing

We compiled a unique corpus of environmental disclosures from the annual reports, 10-K filings, and dedicated environmental or sustainability reports of S&P 500 companies for the fiscal years 2000 through 2004. This period is significant as it precedes the widespread standardization of ESG reporting, ensuring substantial variation in disclosure quality and style. Textual data was extracted and cleaned, removing boilerplate language and standardizing terminology. For market data, we obtained daily stock returns and corresponding market index returns from the Center for Research in Security Prices (CRSP) database.

## 2.2 Stage 1: Thematic Decomposition and Feature Generation using LDA

Instead of using a pre-defined dictionary, we employ Latent Dirichlet Allocation (LDA), a generative probabilistic model, to discover latent thematic structures within the corpus. Each document (disclosure) is modeled as a mixture of a small number of topics, and each topic is characterized by a distribution over words. For each company-year observation, this process yields two primary vectors: a topic proportion vector (the mix of themes in the disclosure) and a set of topic-specific word distributions. From these outputs, we engineer a suite of novel features that proxy for disclosure quality:

- **Thematic Coherence (TC):** Calculated as the average pairwise cosine similarity between the most representative document segments for the dominant topics, measuring internal consistency.

- **Materiality Alignment Score (MAS):** For each industry (based on SIC codes), we derived a set of material environmental topics from regulatory guidelines and NGO reports circa 2003. MAS measures the proportion of a firm's topic distribution that aligns with its industry-specific material topics.

- **Forward-Looking Content (FLC):** The proportion of topic weight associated with topics dominated by future-tense verbs and phrases related to goals, targets, and projections.

- **Specificity Metric (SM):** An entropy-based measure computed on the word distributions of dominant topics, where lower entropy (more peaked distributions) indicates greater use of specific, technical terminology versus vague language.

## 2.3 Stage 2: Feature Selection and DQI Construction using Modified Ant Colony Optimization

The four core features, along with several control features (e.g., document length, readability scores), form a candidate set for predicting market reaction. Rather than using standard step-wise regression or LASSO, we employ a modified Ant Colony Optimization algorithm for feature selection. ACO is inspired by the foraging behavior of ants finding the shortest path to food. In our adaptation, each 'ant' constructs a 'path' (a subset of features) to build a regression model predicting cumulative abnormal returns (CAR). The 'pheromone trail' is updated based on the model's out-of-sample predictive accuracy (mean squared error). Features that consistently contribute to accurate models receive stronger pheromones, making them more likely to be selected by subsequent ants. This bio-inspired approach is particularly adept at navigating complex, non-linear interactions between features that linear selection methods might miss. The final Disclosure Quality Index (DQI) is a weighted linear combination of the features selected by the converged ACO algorithm, where the weights are proportional to the normalized pheromone levels.

## 2.4 Market Reaction Model

We employ an event study methodology to measure the market reaction. The event date is defined as the filing date of the annual report or sustainability report. We estimate the market model over a 255-day estimation window ending 46 days before the event. The DQI and benchmark disclosure metrics are then used as independent variables in a cross-sectional regression explaining the 3-day cumulative abnormal return (CAR [-1, +1]) around the event date. Our key test is whether the coefficient on DQI is statistically significant and whether its inclusion significantly improves the model's explanatory power compared to models using only traditional metrics.

# 3 Results

The application of our hybrid LDA-ACO framework yielded significant and novel findings.

## 3.1 The Disclosure Quality Index (DQI)

The ACO algorithm consistently selected a combination of Thematic Coherence (TC), Materiality Alignment Score (MAS), and Specificity Metric (SM) as the most predictive features. Forward-Looking Content (FLC) was selected less frequently, suggesting its predictive power is conditional on other quality attributes. The resulting DQI is a continuous, normally distributed variable that shows wide variation across firms and within firms over time, confirming its sensitivity to changes in reporting practices.

## 3.2 Predictive Power for Market Reactions

Our core regression results provide strong support for the superiority of the DQI. In models where DQI is the sole disclosure metric, it exhibits a positive and statistically significant coefficient ($p < 0.01$). More importantly, when DQI is included in models alongside traditional measures—such as a simple count of environmental keywords or a binary indicator for report publication—the coefficients on the traditional measures become insignificant, while DQI remains highly significant. The adjusted R-squared for the model with DQI is approximately 0.18, compared to 0.07 for the model using only keyword count. This represents a substantial improvement in explanatory power.

## 3.3 Non-Linear Threshold Effects

A truly original finding emerged from non-parametric analysis and piecewise linear regression. The relationship between DQI and CAR is not linear. We identified a statistically significant threshold effect. For DQI values below the 40th percentile of the sample distribution, the coefficient is small and statistically indistinguishable from zero. However, above this threshold, the coefficient is positive, large, and highly significant. This suggests that investors effectively 'discount' low-quality, vague, or incoherent disclosures, only incorporating the informational signal into prices once disclosures surpass a minimum level of specificity, coherence, and materiality alignment. This threshold effect was not detectable using any traditional linear disclosure metric.

## 3.4 Robustness Checks

The results proved robust to alternative event windows (e.g., CAR [-2, +2], [-1, 0]), different estimation models for expected returns (e.g., the Fama-French three-factor model), and controls

for firm size, profitability, and prior environmental performance. The ACO feature selection process was run with multiple random seeds, with the core set of selected features (TC, MAS, SM) appearing in over 95% of runs, demonstrating the stability of our approach.

# 4  Conclusion

This research makes several original contributions to the literature on environmental disclosure, computational finance, and applied machine learning. Methodologically, we are the first to propose and successfully implement a hybrid framework combining Latent Dirichlet Allocation for thematic analysis with a bio-inspired Ant Colony Optimization algorithm for feature selection in the context of ESG textual analysis. This approach allows us to move beyond the limitations of dictionary methods and linear models, capturing the multi-dimensional and interactive nature of disclosure quality.

Our substantive findings are equally novel. The construction of a continuous Disclosure Quality Index (DQI) derived from textual substance, rather than mere presence or volume of disclosure, provides a superior tool for researchers and practitioners. The discovery of a non-linear, threshold-based market reaction is a significant theoretical insight. It implies that mandating disclosure, without attention to qualitative attributes like coherence and specificity, may be insufficient to inform markets. Investors appear to perform a sophisticated 'filtering' process, only valuing environmental information once it meets a minimum standard of clarity and relevance.

These findings have clear implications. For regulators, they suggest that disclosure guidelines should emphasize qualitative characteristics and materiality. For companies, investing in the clarity, focus, and forward-looking nature of environmental communication may yield tangible benefits in terms of market valuation. For the field of computer science, this work demonstrates the potent application of hybrid, bio-inspired NLP models to complex socio-economic problems.

Future research can extend this framework to other domains of non-financial disclosure (social, governance), incorporate dynamic topic modeling to track thematic evolution, and explore the use of other swarm or evolutionary algorithms for model optimization. The period studied (2000-2004) represents an early stage of voluntary environmental reporting; applying this model to contemporary data would test its validity in a more saturated disclosure environment.

# References

Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research, 3*, 993–1022.

Cormen, T. H., Leiserson, C. E., Rivest, R. L., & Stein, C. (2001). *Introduction to algorithms* (2nd ed.). MIT Press.

Dorigo, M., Maniezzo, V., & Colorni, A. (1996). Ant system: Optimization by a colony of cooperating agents. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics), 26*(1), 29–41.

Fama, E. F., & French, K. R. (1992). The cross-section of expected stock returns. *The Journal of Finance, 47*(2), 427–465.

Gray, R., Kouhy, R., & Lavers, S. (1995). Corporate social and environmental reporting: A review of the literature and a longitudinal study of UK disclosure. *Accounting, Auditing & Accountability Journal, 8*(2), 47–77.

Healy, P. M., & Palepu, K. G. (2001). Information asymmetry, corporate disclosure, and the capital markets: A review of the empirical disclosure literature. *Journal of Accounting and Economics, 31*(1-3), 405–440.

Manning, C. D., & Schütze, H. (1999). *Foundations of statistical natural language processing.* MIT Press.

Patten, D. M. (2002). The relation between environmental performance and environmental disclosure: A research note. *Accounting, Organizations and Society, 27*(8), 763–773.

Solomon, J. F., & Lewis, L. (2002). Incentives and disincentives for corporate environmental disclosure. *Business Strategy and the Environment, 11*(3), 154–169.

Wiseman, J. (1982). An evaluation of environmental disclosures made in corporate annual reports. *Accounting, Organizations and Society, 7*(1), 53–63.