



## Graduation Project Report

---

# Is Greenhushing a Smokescreen or a Panacea for Excess Executive Compensation?

---



### Elaborated by:

Yosri BEN HALIMA, Final Year EPT Student and VIS at Brock University.

### Supported the 28th of June 2024 before the Examination Board:

Prof. Zied SAADAOU, EPT President.

Prof. Amor MASSOUD, EPT Reviewer.

Prof. Samir TRABELSI, Brock University Supervisor.

Prof. Yassine HACHAICHI, EPT Supervisor.

**Academic Year: 2023/2024**

# Acknowledgement

*I would like to express my deepest gratitude to my supervisors, Professor Samir Trabelsi at Brock University and Professor Yassine Hachaichi at Ecole Polytechnique de Tunisie, for their invaluable guidance and support throughout this research. Their expertise and encouragement have been instrumental in the completion of this work.*

*A special thanks goes to Dr. Fares Belkhiria for his insightful advice and brotherly support, which have been a source of great motivation.*

*I am grateful to Professor Amna Chalwati as well for her valuable insights and the enriching conversations that have greatly contributed to my understanding and progress.*

*I am also profoundly thankful to the Goodman School of Business for providing me with the opportunity to conduct research and for offering a conducive academic environment.*

*And lastly, I would like to extend my appreciation to Ecole Polytechnique de Tunisie for the excellent education and for allowing me to embark on this internship experience at Brock University. The solid foundation provided by my home university has been pivotal to my academic and professional development.*

# Abstract

In this study, we develop solutions to proxy excess executive compensation and “greenhushing”—the deliberate underreporting of environmental performance by firms—and explore the relationship between these variables. Utilizing AI and Natural Language Processing techniques such as Latent Dirichlet Allocation and Continuous Bag of Words, we compare corporate disclosures against Sustainability Accounting Standards Board (SASB) guidelines and ESG Scores Gaps. Our findings reveal that both the presence and intensity of greenhushing significantly increase excess CEO compensation, suggesting that executives use greenhushing to mitigate reputational risks by securing higher pay. Internal governance factors, notably Board Size and CEO’s Power Status, further influence this relationship, while the economic determinants of CEO pay add complexity. These insights emphasize the need for boards and compensation committees to account for greenhushing when setting executive pay, promoting transparency and aligning incentives with long-term sustainability goals. This work also can be expanded using broader environmental disclosure sources to enhance the designed solutions, providing crucial insights for investors, regulators, and stakeholders in the ESG landscape.

**Keywords:** Greenhushing, Excess Compensation, Machine Learning, Natural Language Processing, ESG Performance, Executive Pay, Sustainability Goals, Environmental Disclosure.

## Table of Contents

General Introduction .....	7
Chapter 1. Literature Review .....	10
1.1. Executive Compensation & ESG Context.....	10
1.2. Excess Executive Compensation.....	11
1.3. Greenhushing & Silence in Disclosure.....	12
1.4. Machine Learning, Natural Language Processing & Disclosure Analysis.....	13
Chapter 2. Hypothesis Development .....	15
2.1. H1 & H2: Greenhushing, Disclosure, and Excess Compensation .....	15
2.2. H3: Internal Governance & Excess Compensation.....	15
2.3. H4: The Mediating Effect of the CEO's Power Status .....	15
Chapter 3. Data Collecting & Preprocessing.....	17
3.1. Sample Window .....	17
3.2. Data Gathering & Feature Engineering.....	17
Chapter 4. Solution Design for Key Variables Measurements .....	25
4.1. Estimating Excess Executive Compensation.....	25
4.2. Measuring Greenhushing.....	29
4.2.1. The First Pillar: The Environmental Score Gap.....	30
4.2.2. The Second Pillar: Topic Modeling Score .....	31
4.2.3. The Third Pillar: Unrelatedness to Corpus Score .....	39
4.2.4. Deriving and Exploring the Final Greenhushing Measure:.....	40
Chapter 5. Exploring the Relationship Between Greenhushing & Excess Compensation...43	
5.1. Econometric Models.....	43
5.1.1. The First Model: Studying the Impact of the Presence of Greenhushing on Excess Executive Compensation .....	44
5.1.2. The Second Model: Studying the Impact of the Greenhushing Intensity on Excess Executive Compensation .....	44
5.2. Results & Analysis .....	45
5.2.1. The Relationship Between Excess Executive Compensation & the Existence of Greenhushing.....	45
5.2.2. The Relationship Between Excess Executive Compensation & the Intensity of Greenhushing.....	49
5.3. Comparative Analysis for the Two Studies: Presence vs. Intensity.....	52

General Conclusion & Future Works.....	55
Bibliography.....	57

## List of Figures

Figure 1: S&P 500 Companies' Sustainability Reporting Percentage .....	17
Figure 2: S&P 500 Companies' Distribution per GICS Sector .....	18
Figure 3: Numbers of Companies without CDP Disclosure (Red) and Companies with Bad Environmental Performance Based on their Answers (Blue) .....	18
Figure 4: Environmental Pillar Score Distribution.....	20
Figure 5: Pair plot of the Expected Compensation Regression Data .....	21
Figure 6: Correlation Heatmap Expected Compensation Regression Data .....	22
Figure 7: Numbers of Companies where the CEO is a Past Employee (Blue) and Companies where the CEO is a Board Member (Red).....	22
Figure 8: Board Size and Independence Distributions .....	23
Figure 9: Correlation Heatmap for the Final Sample of the Study.....	24
Figure 10: Residuals vs. Dependent Variable Plot .....	26
Figure 11: Distributions & Q-Q Plots of the Residuals of the Executive Compensation Regression .....	27
Figure 12: Autocorrelation Function of Residuals .....	27
Figure 13: Excess CEO Compensation (in USD) Distribution.....	29
Figure 14: Example Documents' Dirichlet Distribution for Three Topics with its 2D Projection for $\alpha >$ , $=$ , and $< 1$ . .....	33
Figure 15: Architecture and Pipeline of the LDA Model .....	34
Figure 16: <i>GHScore</i> Distributions over the years.....	42
Figure 17: Number of Companies involved in Greenhushing Behavior Curve.....	42
Figure 18: Residuals vs Dependent Variable Plot for the First Model of Excess Compensation and Greenhushing .....	47
Figure 19: Distributions & Q-Q Plots of Residuals for the First Model of Excess Compensation and Greenhushing .....	47
Figure 20: Autocorrelation Function Plots of the First Excess Compensation and Greenhushing Model .....	48
Figure 21: Residuals vs Dependent Variable Plot for the Second Model of Excess Compensation and Greenhushing .....	50
Figure 22: Distributions & Q-Q Plots of Residuals for the Second Model of Excess Compensation and Greenhushing .....	51
Figure 23: Autocorrelation Function Plots of the Second Excess Compensation and Greenhushing Model.....	52

## List of Tables

Table 1: Executive Compensation Regression Table for all Models .....	26
Table 2: Goodness of Fit Metrics Table for all Models of Excess Compensation Estimation.....	28
Table 3: Regression Table for the First Excess Compensation & Greenhushing Model .....	45
Table 4: Goodness of Fit Metrics Table for all Specification of the First Model .....	48
Table 5: Regression Table for the Second Excess Compensation & Greenhushing Model.....	49
Table 6: Goodness of Fit Metrics Table for all Specification of the Second Model.....	52

## General Introduction

In concordance with the rising importance of Environmental, Social and Governance (ESG) in the modern operational aspect of a business, this study develops, in a first phase, innovative measures for Excess Executive Compensation and “Greenhushing” - a term coined by “TreeHugger”<sup>1</sup> back in 2008, to describe the purposeful concealment or underreporting of environmental performance or consequences by firms. It is typified by deliberate attempts of firms to omit environmental metrics, under-representing sustainable projects, or deflecting attention from environmental initiatives. In a second phase, we investigate the intricate and dynamic relationship between these two key variables.

This obscuring of the green initiatives stems from a firm’s fear of being accused of greenwashing<sup>2</sup>, hypocrisy or facing negative attention if it falls short in their environmental efforts and climate strategies. According to Eco-Business (2022), Increased scrutiny by the media, NGOs, the public, and consumer and market authorities has made companies more wary about communicating their targets, leading to a concerning trend of less public-facing communication and limited knowledge-sharing. As a result of that stigma, a report by South Pole (2022) found that nearly one in four large, private companies from various sectors have set net-zero targets but have chosen not to publicize their progress, proving the existence of greenhushing.

The fundamental question at the heart of this investigation is “Is greenhushing a panacea or a smokescreen for the issue of excessive executive compensation within corporate entities?”

This question encapsulates a multifaceted exploration into the mechanisms through which greenhushing interacts with and potentially incentivize or disincentivize the increase in CEO compensation packages, serving as a proxy for executive excess compensation. Using sophisticated regression models and leveraging cutting edge Natural Language Processing techniques such as Latent Dirichlet Allocation (LDA) for topic modeling, Continuous Bag of Words (CBoW) via Word2Vec model for topic embedding, Word Frequency Distribution (FreqDist), and Term Frequency-Inverse Document Frequency (TF-IDF) Vectorization for keyword extraction, along with a quantitative statistical analysis and methodical data extraction from public documents using web and document scraping methods, this work intends to quantify greenhushing, estimate excess CEO compensation, and shed light on the intricate

---

<sup>1</sup> A consultancy firm.

<sup>2</sup> A deceptive marketing practice where a company exaggerates or falsely claims the environmental benefits of its products, services, or overall practices to appear more environmentally responsible than it actually is.



interplay between executive remuneration packages and the narratives surrounding corporate sustainability initiatives and disclosure quality.

The significance of this research extends beyond academic inquiry to directly impact businesses in today's dynamic and interconnected global landscape. Understanding and addressing the pervasive phenomenon of greenhushing within corporate sustainability practices is paramount for businesses striving to uphold integrity, credibility, and long-term viability in the eyes of stakeholders. As emphasized by Compliance Week (2022), the alignment between corporate actions and stakeholder expectations, including investors, employees, and customers, is foundational to building and maintaining trust in business operations. Failure to transparently communicate ESG goals and performance can erode this trust, leading to reputational damage and diminished stakeholder confidence, with far-reaching implications for business sustainability and competitiveness.

Additionally, the detrimental impact of excessive executive compensation on both corporations and the broader economy is a topic of concern to this work. As highlighted by the Economic Policy Institute in 2019, CEO compensation packages exhibit a rapid upward trend relative to the wages of typical workers. The implications of this phenomenon, according to EPI, are far-reaching. Firstly, it exacerbates income inequality within corporations and the broader economy. In addition, excessive CEO compensation permeates corporate pay structures, influencing wage distributions throughout the organization. A reduction in CEO pay could potentially lead to more equitable compensation among other executives and, by extension, foster a fairer and more sustainable corporate environment and overall, a more room for shareholder returns.

Moreover, the interplay between greenhushing and executive compensation presents a critical juncture where corporate governance, accountability, and financial stewardship converge. Executives who leverage greenhushing to obscure sustainability shortcomings may inadvertently perpetuate a culture of opacity and misalignment with stakeholder interests out of fear of reputational risk. This potential deterioration of their reputation can have cascading effects on executive compensation practices, potentially incentivizing excessive pay packages that are not commensurate with genuine ESG performance but rather as a mechanism to mitigate such a career risk.

Ultimately, the aim of this work is to significantly broaden our knowledge of corporate governance, sustainability reporting, and executive compensation from a scholarly and practical standpoint. Firstly, this study will pioneer the development of a novel numerical measure for greenhushing, providing a quantitative framework for assessing the extent and impact of greenhushing behaviors within corporate sustainability initiatives and reporting. This work seeks to contribute to the

development of methodological approaches for detecting and mitigating instances of greenhushing and sustainability-related misconduct within corporate reporting. By leveraging cutting-edge techniques in data extraction, analysis, and machine learning, we aim to identify hidden patterns and signals that may indicate potential instances of deception or misrepresentation, thereby enhancing the transparency and reliability of sustainability disclosures.

Secondly, this study is expected to offer insights into the mechanisms through which executive compensation packages may influence corporate sustainability practices and vice versa, either by incentivizing genuine commitment to environmental stewardship or by inadvertently fostering underreporting behaviors. This understanding has profound implications for corporate governance practices, as well as for the design of executive compensation schemes that align with long-term sustainability objectives, promoting transparency whilst protecting executives from the potential reputational risk related to it.

From a practical standpoint, the findings are expected to offer valuable insights for investors, regulators, and corporate stakeholders seeking to navigate the increasingly complex landscape of ESG investing, and sustainable finance. By providing empirical evidence and actionable insights into the relationship between executive compensation, greenhushing, and corporate sustainability practices, we aim to equip decision-makers with the knowledge and tools necessary to make informed investment decisions, drive organizational change, and foster greater accountability and transparency in corporate reporting.

The remainder of this work is organized as follows. Chapter 1 gives a thorough review of the literature used for the solution design. Chapter 2 describes the hypothesis development. Chapter 3 discusses sample collection and variable design and visualizations. Chapter 4 documents the engineering of the key variables' measurements. Moreover, in chapter 5, we develop, identify, and estimate three models for the major study hypotheses: Pooled OLS, Fixed Effect, and Random Effect models. Section 5.2 analyzes how greenhushing affects CEO excess compensation and structure according to each model and assesses individual goodness of fit. Finally, the conclusion discusses findings and potential extensions or improvements to the solution at hand.

## Chapter 1. Literature Review

### Introduction

In this section, we explore the literature on excess compensation and greenhushing. By thoroughly examining their definitions, determinants, and implications in corporate governance, this review provides pivotal insights for the solution engineering to effectively measure excess compensation and greenhushing in corporate settings and lays a solid foundation for developing hypotheses on the relationships between these variables.

### 1.1. Executive Compensation & ESG Context

Executive compensation is a complex and multifaceted aspect of corporate governance, influenced by a wide range of determinants and evolving trends. A foundational study by Core, Guay, and Larcker (2008) has established that executive pay is shaped by various economic determinants. These include the size of the firm, stock price, accounting performance, investment opportunities, CEO tenure and experience, the specific skills and abilities required by the firm, the board's understanding of the CEO's capabilities, and the prevailing labor market conditions. These factors collectively create a nuanced framework for understanding how executive compensation is structured and adjusted over time.

Recent research has further expanded on these foundational insights, particularly regarding the integration of sustainability goals into executive pay. Roberto Barontini et al. (2023) highlight the growing trend of incorporating Environmental, Social, and Governance (ESG) targets into executive compensation packages. This shift is driven by increasing sustainability concerns and the pressure exerted by regulators and institutional investors who demand greater accountability and alignment with broader societal goals. The inclusion of ESG criteria in compensation structures aims to incentivize executives to prioritize long-term sustainability alongside traditional financial performance metrics.

Moreover, the study by Yujuan Wu et al. (2023) explores the direct effects and correlations within the realm of executive green incentives (EGI). The research finds a positive correlation between EGI and ESG performance, suggesting that executives are more likely to achieve higher ESG scores when their compensation is tied to such metrics.

These studies collectively underscore the evolving nature of executive compensation, which now increasingly incorporates ESG considerations as a critical component. This integration not only reflects the changing priorities of stakeholders but also represents a strategic approach to fostering sustainable business practices. As

such, understanding the determinants and trends in executive compensation, including the role of ESG targets, is crucial for comprehensively analyzing the factors that influence executive behavior and compensation. This context is particularly relevant when examining phenomena such as greenhushing and its impact on excess executive compensation.

## 1.2. Excess Executive Compensation

Excess executive compensation, defined as the disparity between total compensation and what is considered expected or reasonable based on firm and market benchmarks, is a critical issue in corporate governance. Syed Ali Rahat Jafri and Samir Trabelsi (2014) present a quantitative approach to defining excess compensation. They propose that excess compensation can be calculated as the difference between total compensation and expected compensation, with the latter being estimated using a logarithmic regression framework using economic determinants of CEO pay. This method provides a systematic way to identify and quantify excess pay, offering a basis for more detailed analysis and potential regulatory interventions.

This concept is also explored through various dimensions in recent and foundational academic literature. Adding to what was cited previously, Roberto Barontini et al. (2023) also discuss the potential risks associated with linking executive pay to ESG performance targets. They argue that such ESG-related pay could inadvertently increase managerial power and lead to higher executive payoffs, potentially allowing executives to extract private benefits without a genuine commitment to ESG values. This concern underscores the need to carefully design compensation packages that align with long-term sustainable goals without unduly inflating executive pay.

Moreover, Yujuan Wu et al. (2023) emphasize the importance of well-crafted incentive mechanisms in executive compensation. They highlight that these mechanisms play a significant role in mitigating conflicts of interest and motivating executives to undertake ventures that are beneficial for long-term corporate growth. By ensuring that incentives are properly aligned with corporate objectives, it is possible to mitigate the problem of excess compensation and promote sustained organizational success.

Additionally, Fried and Shilon (2011) outline the two primary costs associated with excess compensation. Firstly, there is the value diversion from shareholders to executives, meaning that any excess payment to the CEO directly reduces shareholder value. Secondly, excess compensation can lead to value destruction, where excessive pay undermines the positive effects of incentive pay and harms shareholders by reducing the overall efficiency and effectiveness of executive performance incentives.

Furthermore, Hirshleifer and Thakor (1992) provide insights into the risk aversion behavior of CEOs. They argue that due to career concerns, managers tend to be conservative in their project choices, often avoiding risky but potentially value-enhancing actions. In this context, the disclosure of green initiatives and the potential negative attention it might attract can threaten a CEO's reputation. To compensate for this perceived risk, CEOs may seek excess compensation, further exacerbating the issue.

These findings collectively highlight the multifaceted nature of excess executive compensation and the numerous factors that contribute to it. Therefore, understanding these dynamics is essential, particularly in the context of greenhushing, where the interplay between sustainability disclosures and executive pay practices can have significant implications for corporate governance and stakeholder relations. In conclusion, a nuanced approach to designing executive compensation packages is crucial to balance the goals of incentivizing sustainable practices and avoiding the pitfalls of excess pay.

### 1.3. Greenhushing & Silence in Disclosure

Greenhushing, or the deliberate under-communication of sustainability efforts, is a nuanced phenomenon with significant implications for corporate transparency and stakeholder relations. The academic literature provides various perspectives on why companies might choose to remain silent about their green initiatives and the potential consequences of such behavior.

Firstly, Maia Gez et al. (2022) highlight the notable increase in environmental disclosures within ESG reporting. They found that all surveyed companies included environmental information in their 2022 filings, focusing on investments in sustainability, recycling, renewable energy use, and climate change evaluation. This indicates a growing awareness and reporting of environmental matters, yet it also sets a context for understanding why some firms may choose to under-communicate.

Simão and Lisboa (2023) identify several reasons why companies might under-communicate their sustainable development efforts. They suggest that a lack of perceived importance is a key factor: companies might not see their green achievements as significant enough to report. Additionally, to avoid scrutiny and hypocrisy accusations, firms may refrain from highlighting their green efforts, fearing backlash if their overall practices do not fully align with their green claims. Furthermore, there is a fear of negative consequences associated with active communication about sustainability, such as being targeted by activists or regulators.

Supporting these insights, Carlos and Lewis (2018) discuss "Strategic Silence" as a tactic to avoid perceptions of hypocrisy. They found that firms with a strong environmental reputation and higher certification legitimacy are more likely to

withhold publicizing their certification to avoid the heightened risk of being perceived as hypocritical.

Heli Wang et al. (2021) explore similar themes in their study on stakeholder management and quiet giving. They state that firms often remain silent about positive corporate behaviors to avoid stakeholder backlash, particularly in the context of corporate philanthropy. Firms that mistreat primary stakeholders tend to engage in quiet giving, choosing not to disclose charitable donations publicly to avert negative reactions.

Xavier Font et al. (2016) further elaborate on the under-communication of sustainability practices. Their research shows that businesses only communicate about 30% of their sustainability actions, often downplaying complex issues to reduce customer guilt and avoid perceptions of lower competence. This form of greenhushing protects businesses from cynical consumers who might interpret their statements as hypocritical.

Lastly, Andrea Ettinger et al. (2020) examine the desirability of CSR communication. Their study reveals that consumers favor hotels that engage in CSR communication and raise awareness about environmental issues. Compared to greenhushing, CSR communication leads to more favorable attitudes towards the hotels' CSR efforts and lower intentions to behave unethically.

To conclude, the literature on greenhushing and silence in disclosure highlights a range of strategic, reputational, and stakeholder management considerations that drive companies to under-communicate their sustainability efforts. While increased transparency in environmental matters is evident, firms often strategically choose silence to mitigate risks of scrutiny, backlash, and accusations of hypocrisy. Understanding these dynamics is crucial for addressing the broader implications of greenhushing on corporate governance and sustainability reporting.

#### 1.4. Machine Learning, Natural Language Processing & Disclosure Analysis

The integration of Machine Learning and Natural Language Processing in disclosure analysis has revolutionized how companies' public statements and reports are evaluated for sentiment and content. The literature shows several key benefits and advancements in this field, further highlighting the efficiency of our methodology that will be discussed later.

Firstly, Huang and Li (2022) provide a comprehensive performance comparison between machine learning methods and traditional dictionary methods in their study, "Disclosure Sentiment: Machine Learning vs. Dictionary Methods." They find that machine learning methods significantly outperform traditional methods in explaining

market returns at the times of 10-K filings and conference calls. Additionally, machine learning techniques are superior in handling language variations over time and across different industries, showcasing their robustness and adaptability.

Moreover, the use of AI and NLP in disclosure analysis is further exemplified by various studies that leverage these technologies for specialized applications. For instance, Bingler et al. (2022) employ the ClimateBERT model, a pre-trained language model specifically designed for analyzing climate-related text. This model demonstrates the power of AI in handling complex and domain-specific language, providing more accurate and contextually relevant insights from disclosed data.

Similarly, Bajic (2023) uses AI in "Climate Disclosure: A Machine Learning-Based Analysis of Company-Level GHG Emissions and ESG Data Disclosure." This study underscores the efficacy of machine learning in processing large volumes of data to extract meaningful patterns and trends related to greenhouse gas emissions and ESG disclosures. The machine learning-based analysis offers a detailed and nuanced understanding of companies' environmental impact and sustainability practices.

In summary, the growing body of research in this domain highlights the transformative potential of ML and NLP in enhancing the transparency and accountability of corporate disclosure practices, which we will leverage in our methodology.

## Conclusion

The reviewed literature reveals the complexities in the interplay between executive compensation, environmental disclosure and response to reputational risk as well as the use of AI as a toolkit for our measurements. These insights are essential for developing hypotheses and guiding empirical work on their impact on corporate behavior and sustainability practices.



## Chapter 2. Hypothesis Development

### Introduction

In this section, we embark on the crucial task of hypothesis development, laying the groundwork for our empirical analysis in Section 5. Drawing on insights from the literature reviewed earlier, we construct hypotheses to explore the relationships between key variables.

### 2.1. H1 & H2: Greenhushing, Disclosure, and Excess Compensation

Greenhushing highlights the risk of backlash associated with the disclosure of green initiatives, especially when companies fail to meet their sustainability goals, leading to negative attention and accusations of hypocrisy. This reputation risk, often attributed to CEOs who are risk-averse, prioritizes career considerations as noted by Hirshleifer and Thakor (1992). To mitigate this risk, CEOs may seek excess compensation as a buffer against potential career damage. Conversely, if a company is guilty of greenhushing, its reputation may suffer due to non-transparent communication, fueling skepticism about its credibility. Here again, the reputation risk falls on the executives, including the CEO, who might seek excess compensation as a risk management strategy. Thus, in both scenarios—whether greenhushing occurs or not—CEOs may pursue excess compensation to safeguard their position. Therefore, our primary hypotheses are:

**Hypothesis 1:** Greenhushing does not drive excess executive compensation.

**Hypothesis 2a:** The greenhushing behavior positively impacts excess compensation.

**Hypothesis 2b:** The greenhushing behavior negatively impacts excess compensation.

### 2.2. H3: Internal Governance & Excess Compensation

Strong internal governance mechanisms, exemplified by a smaller board size (Jensen, 1993) and higher board independence (Shams Pathan, 2009), are expected to exert greater oversight and constraints on executive compensation practices. This leads us to hypothesize:

**Hypothesis 3:** Stronger internal governance, represented by a smaller board size and higher board independence, has a negative impact on excess compensation.

### 2.3. H4: The Mediating Effect of the CEO's Power Status

The CEO's power status, which includes their level of influence within the organization, can significantly affect executive compensation decisions. Greater CEO power may exacerbate the link between greenhushing behavior and excessive executive pay. Therefore, we hypothesize:



**Hypothesis 4:** Greater CEO power status strengthens the relationship between greenhushing behavior and excess executive compensation.

## Conclusion

These hypotheses collectively address the dynamic interactions between greenhushing practices, internal governance mechanisms, and the movements of executive compensation, providing a broader framework for understanding how these elements interplay within corporate environments.

## Chapter 3. Data Collecting & Preprocessing

### Introduction

In this phase of the project, we explore the critical process of data collection, preprocessing, and feature engineering, where raw data is collected from various sources (WRDS databases, Refinitiv Eikon, Data Scrapping, etc.) and undergoes meticulous treatment to prepare it for analysis. By employing rigorous methods for data cleaning, scaling, and feature extraction, we aim to transform disparate datasets into a cohesive and informative resource.

### 3.1. Sample Window

We obtain data on S&P 500 firms from multiple sources, covering the period from 2017 to 2021. This timeframe is particularly significant as sustainability reporting among these firms surpassed 85% in 2017, marking a critical point in environmental transparency within the corporate sector. However, the availability of parsed Carbon Disclosure Project data sets the limit on the sample to 2021 but still enables for a comprehensive analysis of greenhushing practices.

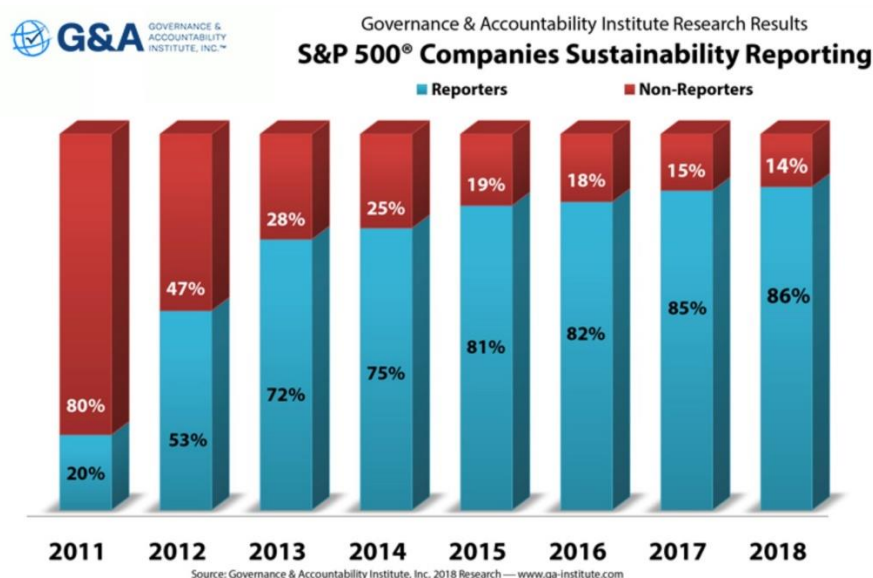


Figure 1: S&P 500 Companies' Sustainability Reporting Percentage

### 3.2. Data Gathering & Feature Engineering

We began the construction of our dataset by first gathering essential information on S&P 500 firms using web scraping techniques. This process involved obtaining tickers, company names, CIK codes, and GICS Sector and Sub-Industry, which provided a comprehensive overview of the companies within our sample.

To measure greenhushing, we utilized responses to the CDP Questionnaire, focusing on questions related to green efforts of S&P 500 companies from 2017 to

2021. The data collection process was meticulously optimized. The CDP data was stored in Excel files, with each year in a separate file. Opening a file entirely took about 7 to 8 minutes, but by accessing specific sheets containing the targeted questions and handling sheet name discrepancies through a recursive approach in Python triggered by exception handling, we reduced the data fetching duration to 6 minutes and 34 seconds for all files combined.

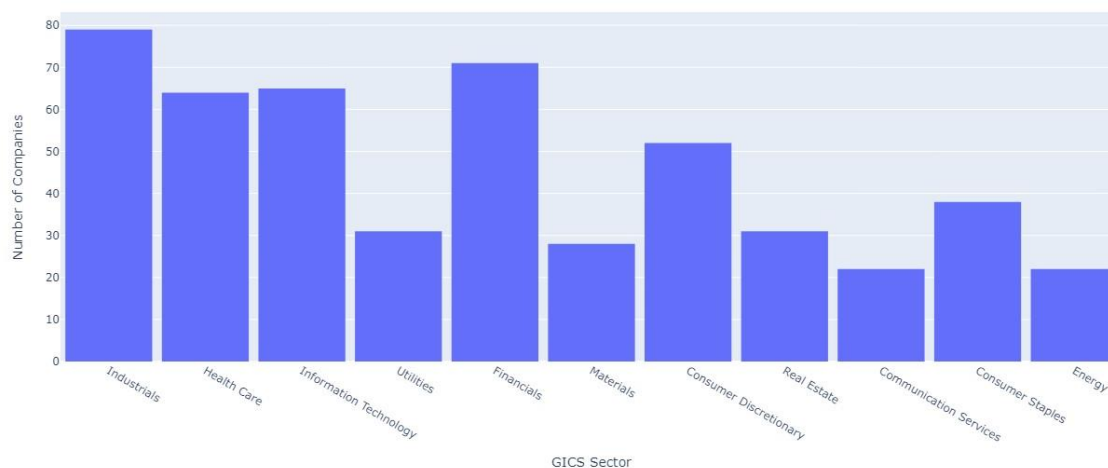


Figure 2: S&P 500 Companies' Distribution per GICS Sector

The targeted CDP questions<sup>3</sup> covered a broad spectrum of environmental topics including climate change governance, risk management, strategy, and emissions targets and sought to focus on companies' green initiatives in financial planning, strategy, and supplier engagement, detailing emissions reduction initiatives, investment methods as well as environmental risk management methods. The collected data facilitated an assessment of greenhushing by comparing disclosed topics against

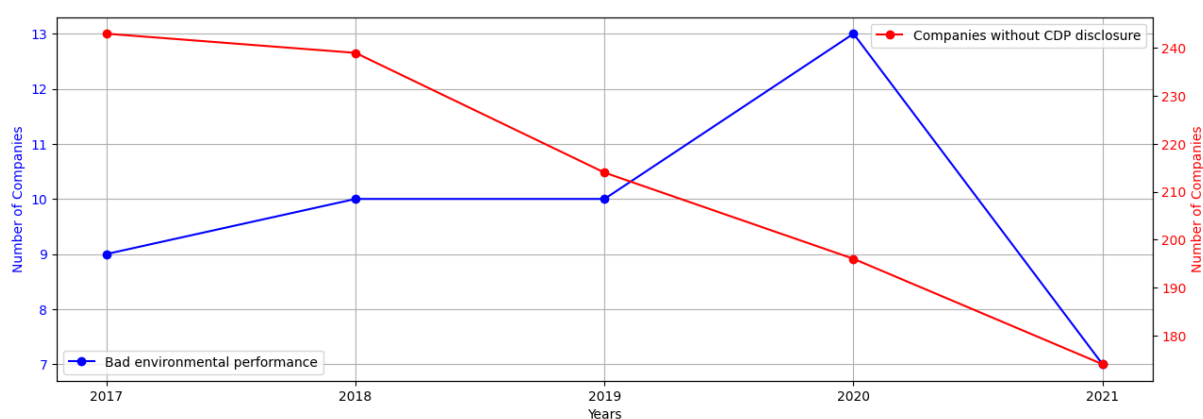


Figure 3: Numbers of Companies without CDP Disclosure (Red) and Companies with Bad Environmental Performance Based on their Answers (Blue)

industry-relevant topics and metrics from the Sustainability Accounting Standards

<sup>3</sup> The questions were manually collected focusing on the proactive aspects.

Board standards<sup>4</sup> and analyzing linguistic patterns and word choices in the disclosures<sup>5</sup>, drawing inspiration from Font et al. (2016). Upon exploring the gathered CDP data, we noticed a significant decrease in the company that do not respond to CDP questionnaire throughout the years. Also, the number of companies with bad environmental performance, proxied by the number of firms whose boards do not investigate environmental problems or do not have a process to assess and mitigate climate risk, is hovering between 7 and 13 companies across years, as demonstrated in Figure 3.

Industry-relevant SASB topics and metrics were meticulously collected by scraping the SASB Standards PDF (".pdf" extension) documents for each industry. These PDFs were manually collected to ensure the accuracy and relevance of the data. Each company was matched with its corresponding SASB Industrial Classification, obtained by scraping the SASB official website. This classification enabled us to identify the specific SASB topics and metrics pertinent to each company in the sample. These topics encompass a broad range of sustainability issues tailored to the unique characteristics and requirements of each industry.

Additionally, we collected Environmental pillar scores for S&P 500 firms from Refinitiv Eikon. Given the differences in fiscal years across companies, we collected data for the last ten fiscal years and organized it to determine the Environmental pillar score for each company within our sample window. This data was used to derive the E Gap Score, which measures the relative difference between a company's score and the mean excess above the 90th percentile, thereby quantifying greenhushing from a numerical perspective. This approach is based on the premise that enhanced disclosure of green innovation correlates with improved ESG scores, as demonstrated by Jianzhuang Zheng et al. (2022)

---

<sup>4</sup> in the Topic Modeling Pillar.

<sup>5</sup> in the Unrelatedness to Corpus Pillar.

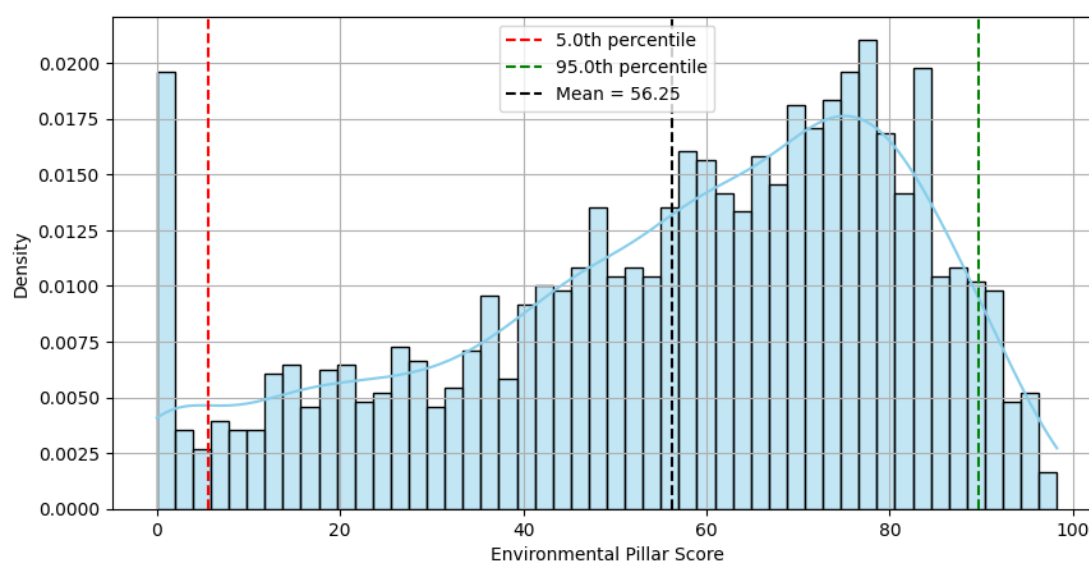


Figure 4: Environmental Pillar Score Distribution

To estimate excess executive compensation, we utilized a comprehensive dataset from Compustat Execucomp's Annual Compensation database, covering the years 2014 to 2023. This data included detailed information about compensation packages, the dates executives became CEOs, and their ages for each specific year. We applied feature engineering to derive the logarithm of CEO tenure, calculated as the difference between the year they became CEO and the current year of the data. Additionally, stock close prices and sales data were sourced from Compustat North America's Annual Fundamentals database. We further applied feature engineering to generate the lagged logarithm of sales and logarithmic stock returns (RET) and their first lags.

Book-to-market ratios and return on assets (ROA) were also collected from Compustat North America's Financial Ratios and were reaggregated on an annual basis instead of daily, with the first lag of ROA being derived as well. These processed data points were used to aggregate total CEO compensation and determine the economic factors influencing CEO pay. As highlighted by Core et al. (2008), The logarithm of tenure served as a proxy for CEO entrenchment and risk aversion, under the premise that more risk-averse executives might seek excess compensation to offset their risk exposure. Logarithmic returns and their first lags, along with ROA and its first lag, were utilized to gauge the firm's financial performance. The lagged logarithm of sales acted as a proxy for firm size, while the book-to-market ratio was used to assess growth opportunities. This multifaceted dataset was paramount for analyzing the determinants of excess CEO compensation. After visual inspection of the pair-plot in Figure 6 we noticed that all variables present a bell-shaped PDF and there is a strong linear correlation between the ROA and its first lag across all companies within each single year which is confirmed in the correlation heatmap in Figure 6.

To analyze the impact of CEO characteristics and board structure on excess compensation, we collected detailed board data. This included the dates CEOs joined their companies and when they became CEOs, allowing us to create a binary variable indicating whether the CEO was a past employee of the company. This variable is important as prior research (Shams Pathan, 2009, Walker et al., 2002) suggests that being a past employee might grant CEOs more power within the company, potentially leading to higher compensation.

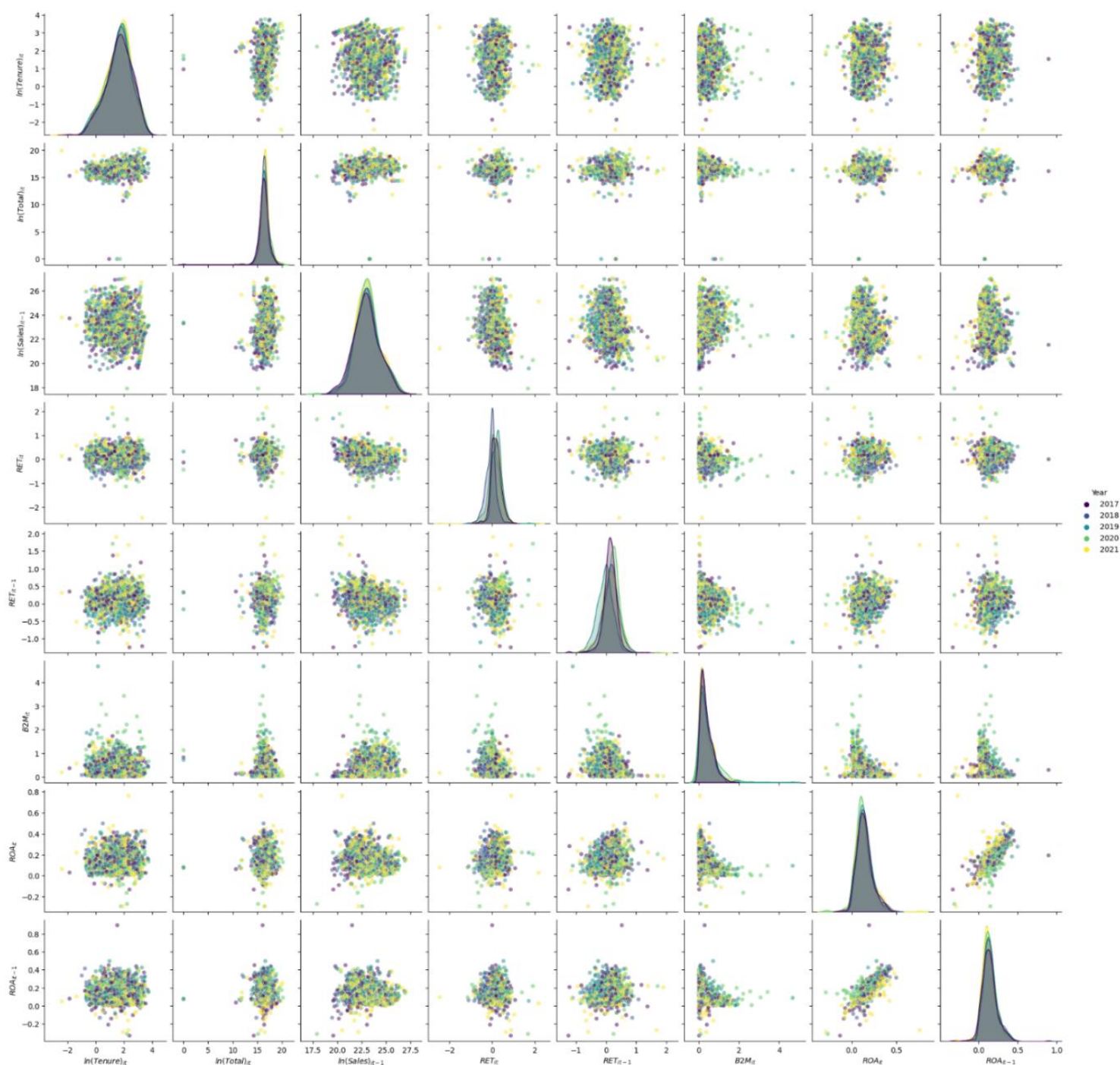


Figure 5: Pair plot of the Expected Compensation Regression Data

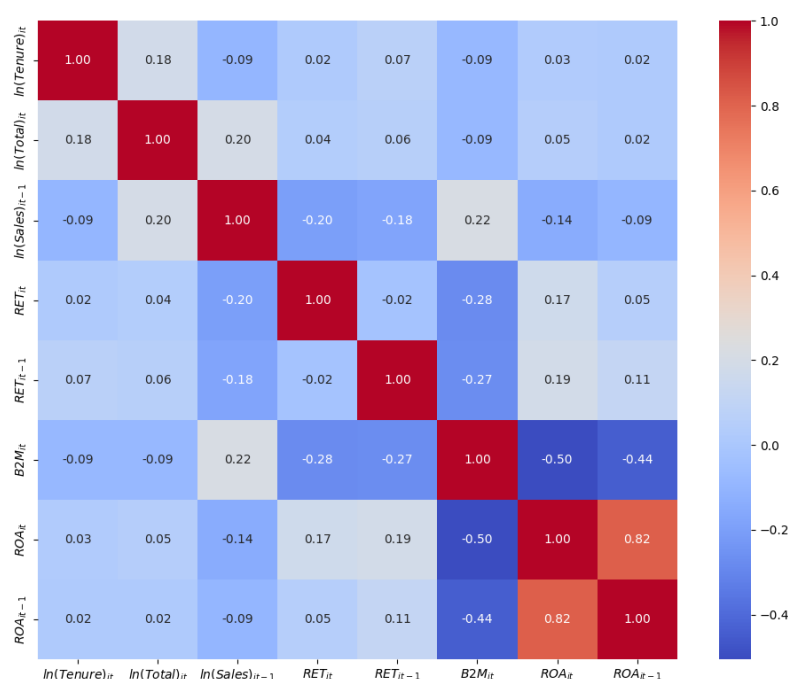


Figure 6: Correlation Heatmap Expected Compensation Regression Data

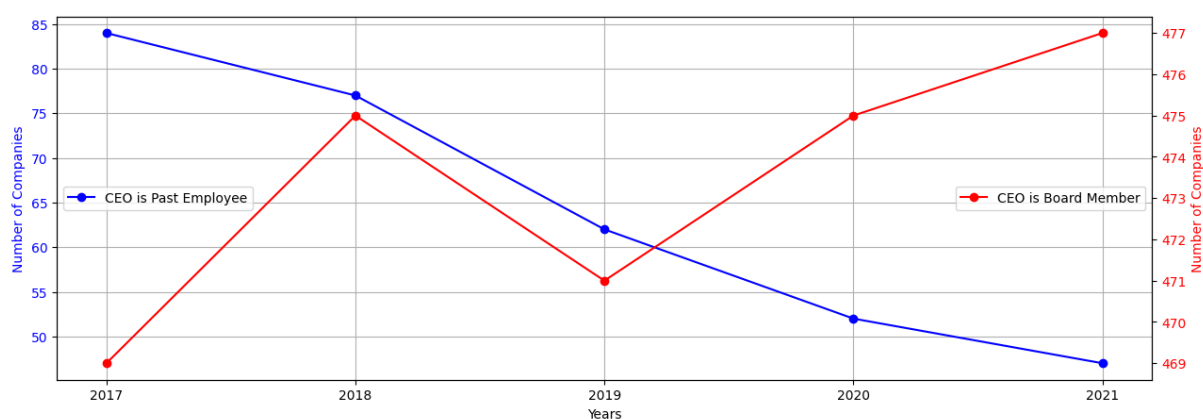


Figure 7: Numbers of Companies where the CEO is a Past Employee (Blue) and Companies where the CEO is a Board Member (Red)

Moreover, data on whether the CEO is a board member, the size of the board, and the board's independence were gathered from Refinitiv Eikon, like the Environmental pillar score query and treatment. Board size and independence are proxies for governance, following the methodology of Shams Pathan (2009). The inclusion of the binary variable indicating whether the CEO is also a board member is designed to further proxy CEO power within the company, complementing the variable of whether the CEO was a past employee as shown in prior literature.

Combining these data points, we measured greenhushing intensity on a scale from 0 to 1 and its existence as a binary variable indicating an intensity greater than 50%. We then estimated excess compensation to compile our final sample, which was winsorized using the IQR range to aggregate 1,841 firm-year observations. We also



addressed the challenge of disparate data scales by applying normalization before the regression. Normalizing the data eliminated distortions caused by vastly varying ranges and enhanced the interpretability of the regression coefficients making it easier to compare the relative impacts of different variables on excess executive compensation.

This unbalanced panel dataset includes the dependent variable of excess CEO compensation and independent variables such as economic determinants of CEO pay, governance variables, CEO's power proxies, and greenhushing measures. We will use this comprehensive dataset to analyze the relationship between greenhushing and excess executive compensation, contributing valuable insights to the discourse on corporate governance and environmental transparency. The correlation matrix for this final sample is shown in Figure 9.

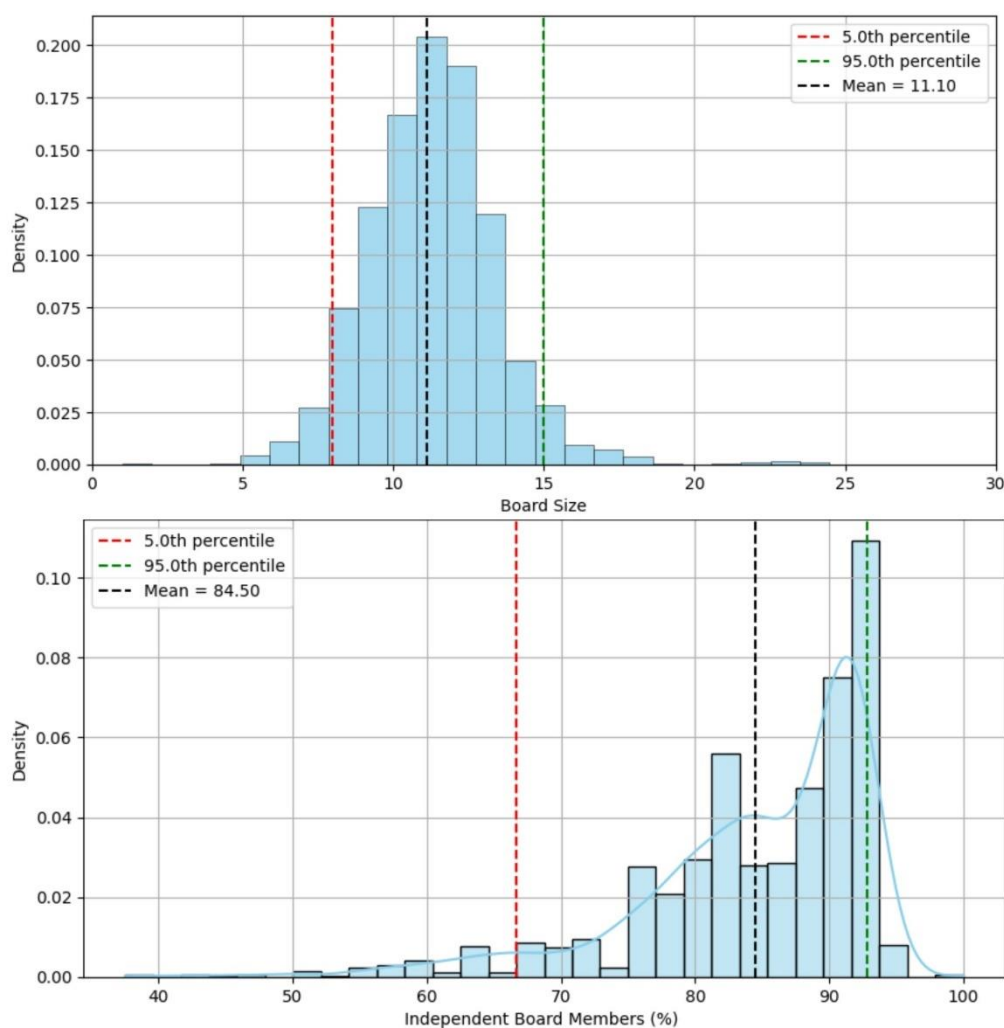


Figure 8: Board Size and Independence Distributions



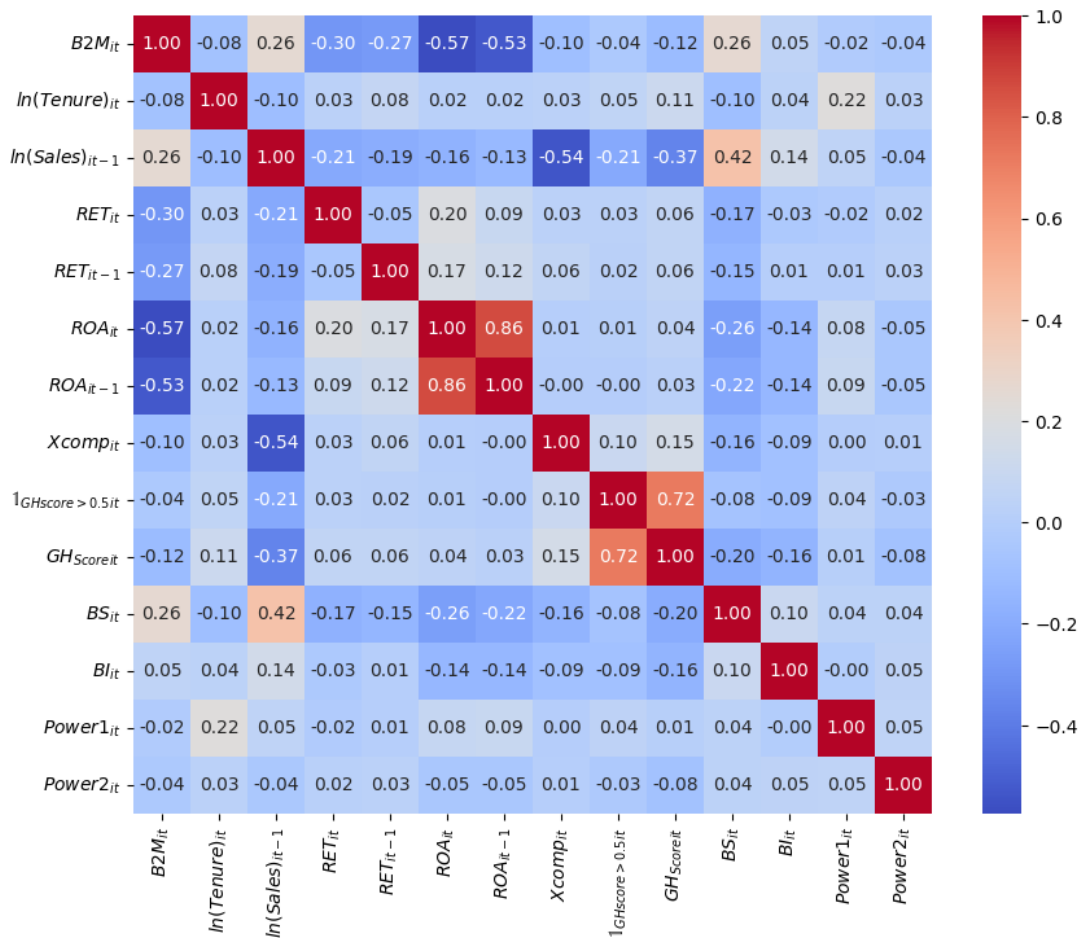


Figure 9: Correlation Heatmap for the Final Sample of the Study

## Conclusion

Through meticulous data collection and feature extraction, we've successfully transformed raw data into a structured dataset ready for analysis. This refined dataset will serve as the foundation for measuring key variables and exploring their interplay econometrically in the next sections.

## Chapter 4. Solution Design for Key Variables Measurements

### Introduction

In this section, we delve into the design of our solution for measuring key variables, namely excess compensation and greenhushing. Our approach begins with an econometric estimation to quantify excess compensation, leveraging rigorous statistical techniques. Additionally, we employ a combination of Natural Language Processing (NLP) and quantitative tools to quantify the extent of greenhushing within corporate disclosures.

#### 4.1. Estimating Excess Executive Compensation

Mathematically, and within the framework outlined by Syed Ali Rahat Jafri and Samir Trabelsi (2014), excess compensation is computed as the difference between total compensation and expected compensation. To derive the expected compensation, a regression model is employed, where the natural logarithm of total CEO compensation  $\ln(\text{Total Compensation})_{it}$  is regressed on various proxies representing economic determinants of CEO pay  $ED_{it}$ . These proxies encapsulate critical factors such as firm performance, firm size, growth opportunities, and industry controls. Specifically,

$$\ln(\text{Total Compensation})_{it} = \beta_0 + \beta_1 \cdot ED_{it} + \varepsilon_{it}$$

where  $ED_{it}$  encompasses variables such as book to market  $B2M_{it}$  to proxy for growth opportunities, the logarithm of CEO's tenure  $\ln(\text{Tenure})_{it}$  as CEOs with longer tenure are more likely to be entrenched and will seek to avoid risk (Berger et al. 1997), sales growth  $\ln(\text{Sales})_{it-1}$  to proxy for firm size, stock return  $RET_{it}$  and  $RET_{it-1}$ , and return on asset  $ROA_{it}$  and  $ROA_{it-1}$  as proxies for firm performance (Core et al. 2008).

$$ED_{it} = \begin{pmatrix} B2M_{it} \\ \ln(\text{Tenure})_{it} \\ \ln(\text{Sales})_{it-1} \\ RET_{it} \\ RET_{it-1} \\ ROA_{it} \\ ROA_{it-1} \end{pmatrix}$$

Table 1: Executive Compensation Regression Table for all Models<sup>6</sup>

The table displays regression results for three econometric models: Pooled OLS, Fixed Effects, and Random Effects. The dependent variable is the logarithm of total CEO compensation, and the independent variables form a vector of economic determinants of CEO pay. For each model, the table provides coefficient estimates, t-values, and significance levels for the economic determinants.

Dependent Variable $\ln(\text{Total Compensation})_{it}$			
Independent Variables	Pooled OLS Model	Fixed Effects Model	Random Effects Model
$B2M_{it}$	-0.1914** (-2.1539)	-0.1622 (-1.4191)	-0.2749*** (-2.9407)
$\ln(\text{Tenure})_{it}$	0.3238*** (11.286)	0.2305*** (8.2780)	0.2608*** (10.081)
$\ln(\text{Sales})_{it-1}$	0.6768*** (168.05)	0.3868*** (3.7986)	0.6905*** (159.27)
$RET_{it}$	0.6465*** (6.5564)	0.2050*** (2.7100)	0.3196*** (4.3338)
$RET_{it-1}$	0.6077*** (6.1196)	0.1647** (2.1563)	0.2446*** (3.2982)
$ROA_{it}$	1.4387*** (2.6400)	0.7484 (1.6458)	1.0552** (2.5144)
$ROA_{it-1}$	-0.2805 (-0.5387)	-0.7938* (-1.6535)	-0.6070 (-1.4794)

The regression and residual analysis results for the total compensation regression on economic determinants were thoroughly assessed using three different models: pooled ordinary least squares, fixed effects, and random effects. Upon visual investigation, it was observed that residuals for all models exhibited some degree of linear correlation with the dependent variable. Notably, pooled OLS model provided the sparsest residuals around the dependent variable, followed by random effects model, indicating the poorest correlations among the three models, as seen in Residuals vs. Dependent Variable ( $\ln(\text{Total Compensation})_{it}$ ) plot in Figure 10.

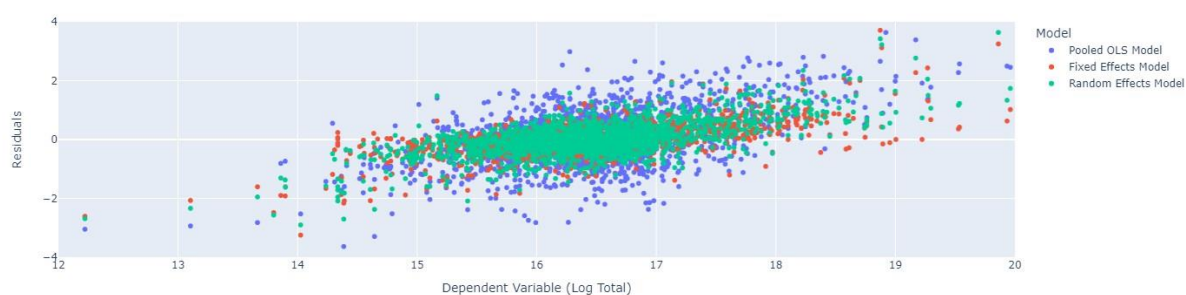


Figure 10: Residuals vs. Dependent Variable Plot

<sup>6</sup> The table displays the coefficient and Student's T (between parenthesis) for each variable in the three models. \*, \*\*, and \*\*\* indicate two-tailed statistical significance at 0.1, 0.05 and 0.01 levels respectively.

Normality tests, such as the Shapiro-Wilk, Anderson-Darling, and Kolmogorov-Smirnov, denied the normality of residuals in all models. However, the Q-Q Plot in Figure 11 shows that the residual distributions of all models have quantiles that are closest to a gaussian distribution around the bulk with a degree of deviation around the tails. Although the residuals deviated from normality, they remained within acceptable bounds.

Furthermore, the presence of fast decaying autocorrelation in the residuals was confirmed for all models with random effects having only one significant lag by inspecting the ACF plots in Figure 12 and applying a variant of the Durbin Watson test suitable for large samples, as described in Lee's (2016) work.<sup>7</sup>

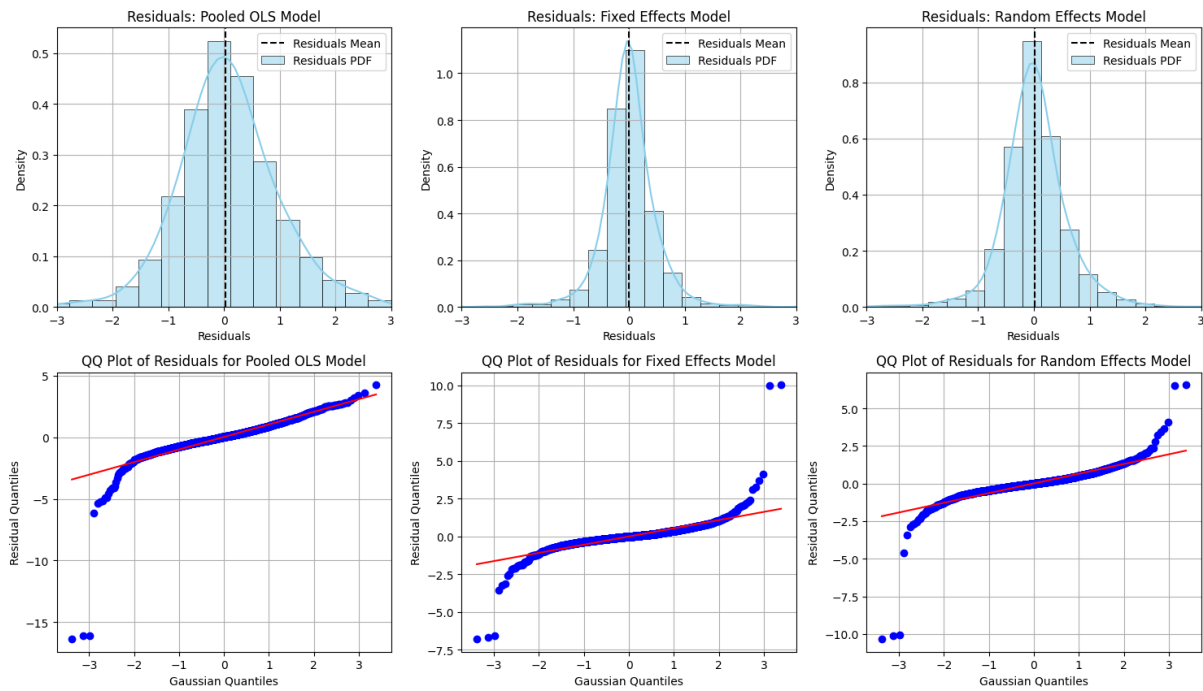


Figure 11: Distributions & Q-Q Plots of the Residuals of the Executive Compensation Regression

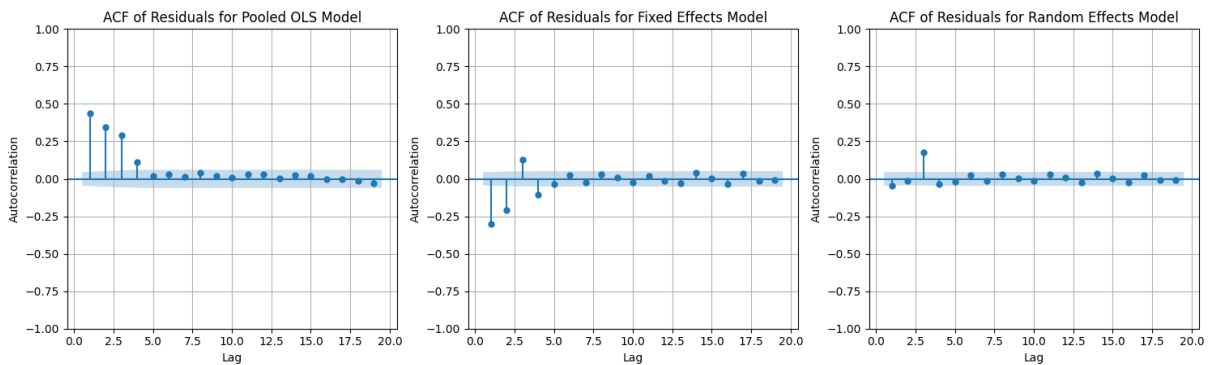


Figure 12: Autocorrelation Function of Residuals

<sup>7</sup> According to Lee (2016),  $z = \frac{(2-d)\sqrt{df}}{2} \sim \mathcal{N}(0,1)$ , where  $d$  is the Durbin-Watson statistic and  $df$  is the number of degrees of freedom of the model. The test rejects the null hypothesis if  $|z| > \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$  where  $\Phi$  is the CDF of a Standard Gaussian and  $\alpha = 5\%$ .

Based on these studies, and on the fact that the random effects model yields the lowest symmetrical mean absolute percentage error (sMAPE), and lower AIC and BIC scores than the Pooled OLS, fellow model in sMAPE, and higher log likelihood value, that model was found to be the best fit for this study.

Table 2: Goodness of Fit Metrics Table for all Models of Excess Compensation Estimation<sup>8</sup>

The table summarizes the performance metrics of three econometric models: Pooled OLS, Fixed Effects, and Random Effects. The metrics included are Log-Likelihood, Akaike Information Criterion (AIC), Bayesian Information Criterion (BIC), and symmetric Mean Absolute Percentage Error (sMAPE).

Metric	Pooled OLS Model	Fixed Effects Model	Random Effects Model
Log Likelihood	-3014.25	-1932.17**	-2184.70
BIC	6081.40	3917.25*	4422.30
AIC	6042.49	3878.34*	4383.39
sMAPE	0.06	0.56	0.06*

The regression analysis of the logarithm of total CEO compensation using random effects model reveals several significant relationships between the dependent variable and a set of independent variables. The coefficient for the book to market ratio is significant at the 1% level and negative, indicating that a higher book to market ratio negatively impacts total CEO compensation. Conversely, the coefficient for the logarithm of CEO's tenure is significant at the 1% level and positive, suggesting that longer CEO tenure positively impacts total compensation. Similarly, the lagged logarithm of sales is significant at the 1% level and positive, implying that higher past sales positively affect total CEO compensation. Both the stock log returns and its first lag are significant at the 1% level and positive, indicating that both current and past stock performance positively impact total CEO compensation. The coefficient for ROA is also significant at the 5% level and positive, demonstrating that higher ROA, reflecting better profitability, positively impacts total CEO compensation. In contrast, the coefficient for the first lag of ROA is negative but not significant, suggesting that the lagged ROA does not provide additional explanatory power for CEO compensation beyond the current ROA due to high collinearity with the current ROA.

Through this regression, the expected compensation can be derived as:

$$\ln(\text{Expected Compensation})_{it} = \beta_0 + \beta_1 \cdot ED_{it}$$

which can be transformed into expected compensation as  $e^{\ln(\text{Expected Compensation})}$ .

This estimation approach enables an understanding of the factors driving CEO compensation and allows for the identification of any excess compensation beyond

<sup>8</sup> \*\* highlights the maximum value for Log Likelihood. \* points out the lowest value for AIC, BIC and sMAPE.

what is expected based on economic determinants leading to the following distribution of excess CEO pay (Figure 13).

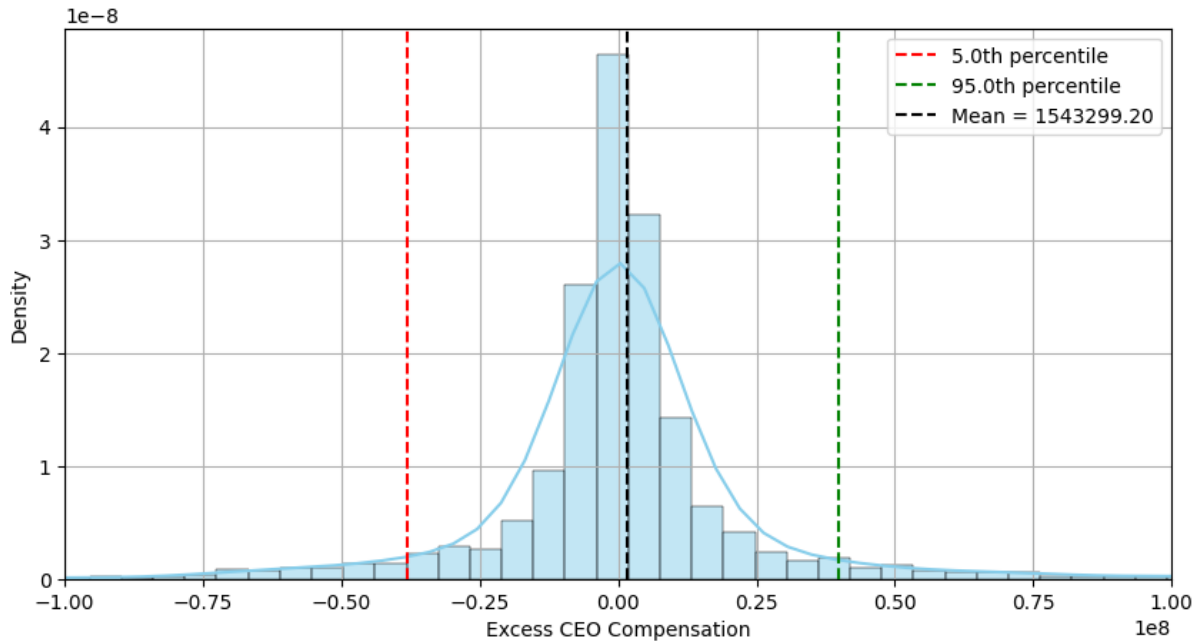


Figure 13: Excess CEO Compensation (in USD) Distribution

## 4.2. Measuring Greenhushing

Our approach to measuring greenhushing is inspired by the work of X. Font et al. (2016), which utilizes text analysis to explore corporate communication strategies. While both methodologies are grounded in text analysis, there are notable differences in our approach that enhance its precision and applicability.

Firstly, unlike X. Font et al.'s work, which focuses on comparing company disclosures in audits versus advertisements, our methodology compares what companies say they do with what they should say they do according to the Sustainability Accounting Standards Board (SASB). This alignment with SASB standards ensures that our assessment is based on established benchmarks of sustainability efforts, providing a more objective and relevant measure of greenhushing. Secondly, our approach leverages advanced AI techniques to provide a numerical figure for measuring greenhushing. This integration of AI allows for more robust and scalable analysis, capable of processing large volumes of text data to detect discrepancies between reported and expected sustainability practices.

The Greenhushing Score  $GH_{Score}$ , the numerical guideline developed in our methodology, is computed using the following formula:

$$GH_{Score} = \alpha \cdot EG_{Score} + \beta \cdot TM_{Score} + \gamma \cdot UC_{Score}$$

$$\alpha = \beta = \gamma = \frac{1}{3}$$

Where:

- $EG_{Score}$  is the normalized E (pillar of the ESG) Gap Score between the company at hand and the score benchmark used as ideal score of a company.
- $TM_{Score}$  is the Topic Modeling Score quantifying the difference between the relevant topics for the industry set by regulators (SASB) and those disclosed in the "CDP answers data" to the public.
- $UC_{Score}$  is the Unrelatedness-to-Corpus Score measuring the unrelatedness of the text in the CDP answers data to a corpus of words related to the environmental efforts relative to the SASB Industrial Classification sector of the company.

It is worth noting that all these scores will be normalized (ranging between 0 and 1) after computation and one can notice that the Greenhushing Score is a weighted sum that will range from 0 to 1 with 1 indicating total greenhushing behavior and 0 indicating that no deliberate under-communication is detected.

In the following, we will dissect each of the 3 pillars of the Greenhushing Score and explain its contribution in capturing the under-communication behavior in a company.

#### 4.2.1. The First Pillar: The Environmental Score Gap

The  $EG_{Score}$ , representing the Gap in the E Score, offers one of the pillars of the measure for capturing greenhushing within a company's environmental communication framework.

It is given by the following formula:

$$EG_{Score_{it}} = \max\left(0, \frac{\bar{E}_t - E_{it}}{\bar{E}_t}\right)$$

Where:

- $\bar{E}_t$  is the desired standard of the E score for all companies for year  $t$ .
- $E_{it}$  is the company's actual E score.

It identifies discrepancies between actual environmental performance and the desired standard. By focusing on the right tail of the E score distribution at 10%, the  $EG_{Score}$  highlights instances where a company falls short of expected levels of transparency and disclosure regarding its environmental practices indicating potential deliberate under-communication behavior.

This is explained by the findings of Jianzhuang Zheng et al. (2022) suggesting that green innovation, which can only be identified if the company documents and

discloses those initiatives, can significantly enhance ESG scores. Consequently, the  $EG_{score}$  can capture a chunk of the greenhushing patterns.

#### 4.2.2. The Second Pillar: Topic Modeling Score

The Topic Modelling Score,  $TM_{score}$ , is the second pillar of our methodology for quantifying and detecting instances of greenhushing within a company's communication channels. Drawing inspiration from X. Font et al. (2016), who pioneered identifying greenhushing by juxtaposing actual green initiatives against communicated messages and actions via a linguistic analysis, this approach focuses on comparing the topics disclosed in the CDP questionnaire with those deemed relevant by SASB for the company's industrial sector.

By using Latent Dirichlet Allocation (LDA) for topic modelling, a Natural Language Processing technique, this method identifies and extracts key themes and subjects within documents. It then quantifies the disparity between topics disclosed in CDP answers and the SASB topics for that industry using numerical embedding techniques via Word2Vec's Continuous Bag of Words (CBoW) model to get their vectorial forms. After that we compute cosine similarity of topics to derive  $TM_{score}$ . The formula for computing the score captures the divergence in topic representation, with a higher score indicating a greater dissonance between official documents and those targeting customers.

$$TM_{score_{it}} = 1 - TS_{score_{it}}$$

Where  $TS_{score_{it}}$  is the Topics Similarity Score for company  $i$  and year  $t$  calculated by embedding topics in vectors all having the same dimensions and then performing cross vectors cosine similarity to quantify similar ones between the two types of documents and average them.

This methodology's strength lies in its ability to uncover instances where green initiatives guided by the compliance with SASB<sup>9</sup> are omitted or downplayed in communications, and it describes the following process:

1. Extracting topics from CDP answers for each company.
2. Numerically embedding CDP topics and SASB topics.
3. Computing Topic Similarity Score and Topic Modeling Score.

##### 4.2.2.1. Extracting topics from CDP answers using Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA) is a generative probabilistic model and one of the most popular approaches for topic modeling in a corpus, introduced by Blei, Ng, and Jordan in 2003. It represents documents as random mixtures over latent topics,

<sup>9</sup> According to the BDO Insights article, 452 out of the S&P 500 companies aligned their ESG reporting with the SASB standards in 2021. This represents over 90% adoption among the S&P 500



with each topic characterized by a distribution over words. The high-probability words in a topic help identify the topic's nature. There are two core concepts that govern LDA, according to H. Jelodar et al. 2018:

1. **Documents as Mixtures of Topics:** Each document in the corpus is viewed as a mixture of several topics. The model assumes that documents are generated by selecting topics and then words from these topics.
2. **Topics as Distributions Over Words:** Each topic is characterized by a distribution over words. This means that some words have a higher probability of appearing in a particular topic than others.

LDA can be understood geometrically by envisioning the creation of simplices. A  $K$ -Simplex, where  $K$  is the number of topics, is formed where each vertex represents a distinct topic. Documents are then placed within this simplex closer to the vertices (topics) that they are most representative of, indicating their topic distribution (Dirichlet Distribution) as shown in the 3-dimensional case in Figure 14. Similarly, a  $V$ -Simplex, where  $V$  is the size of the vocabulary, is created where each vertex represents a unique word in the vocabulary. Topics are placed within this simplex closer to the vertices (words) that they frequently contain. This geometric interpretation helps visualize how LDA assigns documents to topics and topics to words, with proximity to the vertices indicating a stronger association.

In this model:

- ❖  $D$  is the original set of CDP answer texts for a given company  $i$  and a given year  $t$ .
- ❖  $M$  is the number of documents in  $D$ .
- ❖  $N_d$  is the number of words in the  $d$ -th document  $d \in \{1, \dots, M\}$ .
- ❖  $N = \sum_{d=1}^M N_d$  is the total number of words in the corpus.
- ❖  $K$  is the number of topics.
- ❖  $V$  is the size of the vocabulary. It is the total number of unique words across all documents in the corpus.
- ❖  $\theta_d$  is the topic distribution for document  $d$  and  $\phi_k$  is the word distribution for topic  $k$  (latent variables). A Dirichlet distribution of parameter  $\alpha$  is given by the following PDF:

$$f(x|\alpha) = \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K x_i^{\alpha_i-1} = \frac{1}{B(\alpha)} \prod_{i=1}^K x_i^{\alpha_i-1}$$

where:

- $B(\alpha) = \frac{\prod_{i=1}^K \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^K \alpha_i)}$  is the multivariate beta function.
- $\Gamma(\cdot)$  is the gamma function.

- ❖  $\alpha = (\alpha_1, \dots, \alpha_K)$  and  $\beta = (\beta_1, \dots, \beta_V)$  are hyperparameters that control the sparsity of the distributions where  $\alpha_1 = \dots = \alpha_K$  and  $\beta_1 = \dots = \beta_V$ . For mathematical simplicity we will treat  $\alpha$  and  $\beta$  as scalars.

Using the text data fed to (original corpus of documents) LDA, this model follows a generative process for creating new observed corpuses of documents, as shown in Figure 15:

1. **Word Distribution for Each Topic:** For each topic  $k$  (where  $k \in \{1, \dots, K\}$ ), select a multinomial word distribution  $\phi_k$  from a Dirichlet distribution with parameter  $\beta$ .
2. **Topic Distribution for Each Document:** For each document  $d$  (where  $d \in \{1, \dots, M\}$ ), choose a multinomial topic distribution  $\theta_d$  from a Dirichlet distribution with parameter  $\alpha$ .
3. **Word Assignment in Documents:** For each word position  $n$  (where  $n \in \{1, \dots, N_d\}$ ) in document  $d$ :
  1. Select a topic  $z_{dn}$  from the distribution  $\theta_d$ .
  2. Select a word  $w_{dn}$  from the distribution  $\phi_{z_{dn}}$ .

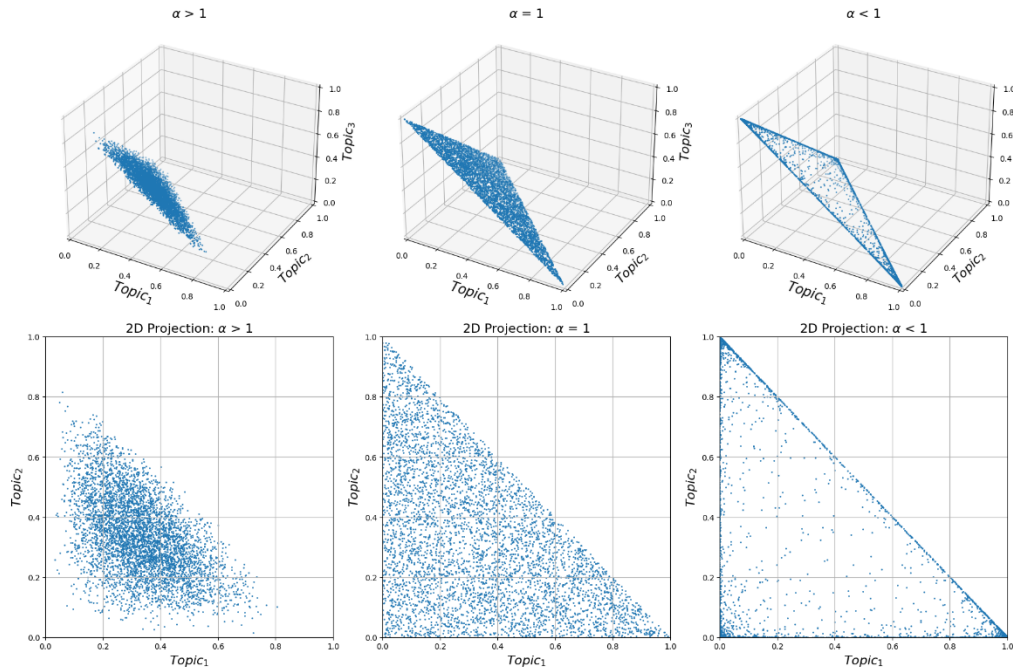


Figure 14: Example Documents' Dirichlet Distribution for Three Topics with its 2D Projection for  $\alpha >$ ,  $=$ , and  $< 1$ .

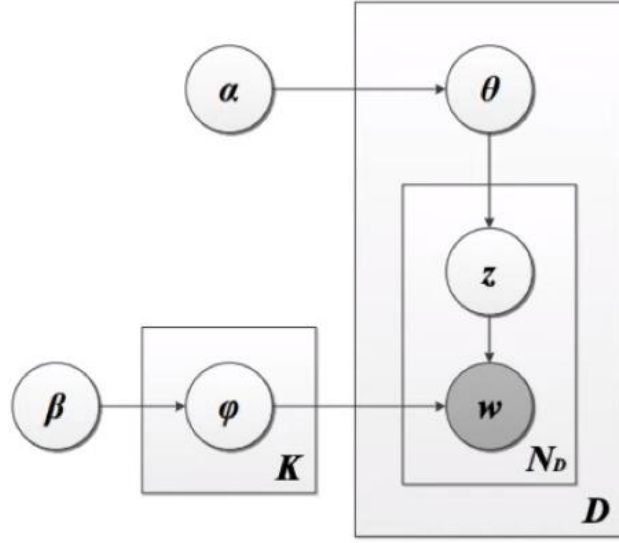


Figure 15: Architecture and Pipeline of the LDA Model

The goal of inference in LDA is to estimate the set of latent topic distributions  $\theta$  and word distributions  $\phi$  that are most likely to have generated the original corpus. This involves computing the posterior distribution of the latent variables given the observed data. The probability of the observing  $D$  is given by:

$$p(D|\alpha, \beta) = \prod_{d=1}^M \int p(\theta_d|\alpha) \left( \prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn}|\theta_d) p(w_{dn}|z_{dn}, \beta) \right) d\theta_d$$

Where:

- $p(D|\alpha, \beta)$ : This is the probability of observing the entire original document set  $D$  given the hyperparameters  $\alpha$  and  $\beta$ .
- $p(w_{dn}|z_{dn}, \beta)$ : Given a topic assignment  $z_{dn}$  for each word  $w_{dn}$  in document  $d$ , this term calculates the probability of observing that word under the topic-word distribution  $\beta$ . It tells us how likely it is to see a particular word given its assigned topic.
- $p(z_{dn}|\theta_d)$ : This term represents the probability of assigning word  $w_{dn}$  in document  $d$  to a particular topic  $z_{dn}$ , given the document's topic mixture  $\theta_d$ . It reflects how probable it is for a word in a document to belong to a specific topic based on the document's overall topic distribution.

The formula encapsulates the joint probability of observing all the documents in the corpus under the LDA model. It combines the probabilities of generating each word in each document from its corresponding topic, weighted by the document-topic distribution and the topic-word distribution. The integral over  $\theta_d$  represents

uncertainty in the topic distribution of each document and the multiplying these probabilities over all documents in the corpus gives us the overall likelihood of observing the entire document set.

The goal in LDA is to estimate the parameters  $\alpha$  and  $\beta$  that maximize this probability, all while keeping each document (resp. topic) as monochromatic as possible (e.g., associated to a minimal number of topics, resp. words, ideally one unique topic resp. word). However, since this integral is often intractable to compute directly, approximate inference methods like Gibbs sampling are typically used.

Gibbs sampling in LDA is a Markov Chain Monte Carlo (MCMC) algorithm that generates samples from the posterior distribution to help us infer the hidden topic structure of a corpus by iteratively updating the topic assignments of words in documents based on their conditional probabilities which is influenced by the proportion of words assigned to each topic in the document (document-topic distribution) and the proportion of times a word is assigned to a topic across the entire corpus (topic-word distribution). This process allows us to approximate the joint distribution of topic assignments, enabling us to estimate the underlying distributions that define the topics and their relationships to the words and documents. The sampling procedure in Gibbs Sampling method describes the following process:

1. **Initialize  $z_{dn}$** : Start with random topic assignments for each word in each document:  $\alpha^{(0)}$  and  $\beta^{(0)}$

2. **Update  $z_{dn}$** : For each word  $w_{dn}$  in each document  $d \in \{1, \dots, M\}$ :

- Remove  $w_{dn}$  from its current topic count
- Calculate the conditional probability of assigning  $w_{dn}$  to each topic  $k \in \{1, \dots, T\}$  given by:

$$p(z_{dn} = k | z_{-dn}, w_{dn}, \alpha, \beta) \propto \frac{n_{d,k}^{-dn} + \alpha_k}{N_d + \sum_k \alpha_k} \cdot \frac{n_{k,w_{dn}}^{-dn} + \beta_{w_{dn}}}{n_k + \sum_w \beta_w}$$

where:

- $n_{d,k}^{-dn}$  is the count of words in document  $d$  assigned to topic  $k$  excluding the current word.
- $n_{k,w_{dn}}^{-dn}$  is the count of word  $w_{dn}$  assigned to topic  $k$  excluding the current word.
- $N_d$  is the total number of words in document  $d$ .
- $n_k$  is the total count of words assigned to topic  $k$ .
- Update the parameters  $\alpha$  and  $\beta$ : The update equation for  $\alpha$  involves maximizing the log-likelihood of the document-topic distributions. One common approach is to use Newton's method:

$$\alpha_k^{(t+1)} = \alpha_k^{(t)} + \frac{N}{M} \left( \psi \left( \sum_{i=1}^M \alpha_k^{(t)} \right) - \psi(\alpha_k^{(t)}) + \frac{1}{M} \sum_{d=1}^M \left( \psi(n_{d,k}^{(t)} + \alpha_k^{(t)}) - \psi \left( \sum_{j=1}^K n_{d,j}^{(t)} + \sum_{j=1}^K \alpha_j^{(t)} \right) \right) \right)$$

$$\beta_w^{(t+1)} = \beta_w^{(t)} + \frac{V}{K} \left( \psi \left( \sum_{i=1}^K \beta_w^{(t)} \right) - \psi(\beta_w^{(t)}) + \frac{1}{K} \sum_{k=1}^K \left( \psi(n_{k,w}^{(t)} + \beta_w^{(t)}) - \psi \left( n_k^{(t)} + \sum_{v=1}^V \beta_v^{(t)} \right) \right) \right)$$

Where:

- $\alpha_k^{(t)}$  is the value of  $\alpha$  for topic  $k$  at iteration  $t$ .
  - $\beta_w^{(t)}$  is the value of  $\beta$  for word  $w$  at iteration  $t$ .
  - $N = \sum_{d=1}^M N_d$  is the total number of words in the corpus. It is the sum of the word counts across all documents.
  - $V$  is the size of the vocabulary. It is the total number of unique words across all documents in the corpus.
  - $K$  is the number of topics.
  - $M$  is the total number of documents in the corpus.
  - $n_{d,k}^{(t)}$  is the count of words assigned to topic  $k$  in document  $d$  at iteration  $t$ .
  - $n_{k,w}^{(t)}$  is the count of word  $w$  assigned to topic  $k$  at iteration  $t$ .
  - $\psi(\cdot)$  is the digamma function, the derivative of the log gamma function.
- Sample a new topic assignment for  $w_{dn}$  based on these probabilities.
  - Update the topic counts with the new assignment.

3. **Convergence:** After many iterations, the topic assignments  $z_{dn}$  will approximate samples from the true posterior distribution leveraging the Ergodic Theorem for Markov Chains. These samples can then be used to estimate the parameters of the document-topic distributions  $\alpha$  and the topic-word distributions  $\beta$ .

Gibbs sampling can also be intuitively understood as a type of "coloring problem", especially in the context of the LDA model, hence the previous use of "monochromatic" to qualify the best topic distribution. Think of each word in the corpus as a node in a graph and the topic assignments  $z_{dn}$  as colors that can be assigned to these nodes. In Gibbs sampling, the topic (color) of each word (node) is updated based on the current topic assignments (colors) of all other words (nodes). Start with an initial assignment of topics (colors) to each word (node). This can be done randomly. For each word, compute the conditional probability distribution of topics given the current assignments of all other words then reassign the word to a new topic (color) based on this distribution. After many iterations of reassigning topics (re-coloring nodes), the algorithm converges to a stable state where the topic assignments reflect the underlying structure of the corpus. In a coloring problem, this would

correspond to a state where the colors are consistently assigned according to the rules of the problem.

In the context of our methodology, to systematically extract topics from CDP disclosures for each company in the S&P 500, a structured approach is employed. The first step is data preparation, where CDP responses are collected for each company and subjected to preprocessing to eliminate noise and irrelevant information, ensuring the dataset is clean and focused. Following this, the text undergoes tokenization, breaking it down into individual words or tokens. This step includes the removal of common stop words—words that do not carry significant meaning—and non-alphanumeric characters. Additionally, lemmatization is performed to reduce words to their base or root forms, thereby standardizing the vocabulary.

Subsequently, the cleaned and tokenized text is converted into a document-term matrix using the Bag of Words (BoW) model. In this matrix, each row corresponds to a document, and each column corresponds to a unique term from the vocabulary. The entries in the matrix represent the frequency of each term in each document, facilitating a structured representation of the text data. The next phase involves training a LDA model on the document-term matrix. This is achieved through Gibbs Sampling.

The LDA model learns the topic distribution for each document and the word distribution for each topic, effectively identifying clusters of terms that frequently co-occur across the corpus. Finally, the topics are extracted and are represented by the top words associated with each inferred topic in the company's CDP disclosures.

#### 4.2.2.2. *SASB and CDP Topics' Numerical Embedding using Word2Vec's CBoW*

To get the vectorial form of each extracted topic we need to get the numerical embeddings of each word in it which will be performed using Word2Vec's CBoW. CBoW model of Word2Vec learns word embeddings by training a neural network to predict target words given their surrounding context words. By maximizing the conditional probability of target words through iterative optimization, the model effectively captures the semantic relationships within the corpus, resulting in distributed representations of words.

Suppose we have a vocabulary consisting of  $V$  unique words, where each word is represented by a one-hot encoded vector. For any given word  $w$  in the vocabulary, its one-hot encoded vector  $x_w$  has a dimensionality of  $V$ , with all elements being zero except for the position corresponding to the index of  $w$ , which is one. The context window size, denoted as  $C$ , determines the number of context words considered around each target word.

The architecture of the CBoW model comprises three layers: an input layer, a hidden layer, and an output layer. The input layer includes the one-hot encoded vectors

of the context words surrounding the target word. The hidden layer consists of  $D$  neurons, where  $D$  is the dimensionality of the word embeddings. The output layer contains  $V$  neurons, each corresponding to a word in the vocabulary.

The training objective of the CBoW model is to maximize the probability of predicting the target word given its context words. Mathematically, the conditional probability  $P(w_t|\{w_c\})$  of predicting the target word  $w_t$  given its context words  $\{w_c\}$  is defined using the softmax function:

$$P(w_t|\{w_c\}) = \frac{e^{v_{w_t}^T h}}{\sum_{v=1}^V e^{v_v^T h}}$$

Where:

- $v_{w_t}$  is the output vector corresponding to the target word  $w_t$ .
- $h$  is the average of the input vectors (embeddings) corresponding to the context words  $\{w_c\}$ .

The objective function aims to maximize the average log probability of predicting target words given their context words across the entire corpus:

$$\frac{1}{T} \sum_{t=1}^T \log P(w_t|\{w_c\})$$

where  $T$  represents the total number of target words in the corpus. Optimization of this objective function is typically performed using stochastic gradient descent, which iteratively updates the model parameters (word embeddings).

Mathematically, let's represent a topic as a sequence of words:  $t = (w_1, w_2, \dots, w_n)$ , where  $n$  is the number of words in the topic and  $v_{w_i}$  is the CBoW embedding vector for word  $w_i$  in the topic  $t$ . The numerical embedding of the topic is given by the following function:

$$f(t) = f(w_1, w_2, \dots, w_n) = \frac{1}{n} \sum_{i=1}^n v_{w_i}$$

#### 4.2.2.3. Computing Topic Similarity Score and Topic Modeling Score

After extracting the topics disclosed by a certain company within the sample for a certain year using LDA, and gathering the SASB industry-relevant topics, we quantify topics similarity by computing cosine similarity between the numerical representations of each  $(CDP_{Topic}, SASB_{Topic})$  pair. We then identify SASB topics in the CDP disclosed topics by picking for each SASB topic the CDP topic that maximizes the cosine similarity of the pair. Finally, we average those similarities across SASB topics to derive the  $TS_{Score}$  and consequently the  $TM_{Score}$ . The formula for  $TS_{Score}$  is given by:



$$TS_{Score} = \frac{1}{m} \sum_{j=1}^m \max_{i \in \{1, \dots, n\}} \left( C \left( f(t_i), f(t'_j) \right) \right)$$

Where:

- $\{t_1, t_2, \dots, t_n\}$  is the set of topics discussed by the company in the CDP for the given year.
- $\{t'_1, t'_2, \dots, t'_m\}$  is the set of the SASB target topics for the industry of that company.
- $f(.)$  is the function that computes the numerical embedding of a topic.
- $C(.,.)$  is the cosine similarity function.

Hence:

$$TM_{Score} = 1 - \frac{1}{m} \sum_{j=1}^m \max_{i \in \{1, \dots, n\}} \left( C \left( f(t_i), f(t'_j) \right) \right)$$

#### 4.2.3. The Third Pillar: Unrelatedness to Corpus Score

The  $UC_{Score}$ , or Unrelatedness-to-Corpus Score, plays a pivotal role in capturing aspects of greenhushing behavior within the company's disclosure, again inspired by X. Font et al.'s linguistic analysis. By measuring the alignment of the CDP text to a corpus of environmental sustainability-related words tailored to the company's industry sector, this metric offers insights into the degree to which green initiatives are communicated transparently and comprehensively via word choice.

Specifically, the  $UC_{Score}$  evaluates the absence of sector-specific environmental concerns within the text, providing a quantitative measure of the extent to which the company addresses the disclosure of its initiatives for relevant environmental sustainability issues. It quantifies the proportion of words of the corpus that are absent in the set of keywords (extracted using TF-IDF) of the AI-preprocessed CDP text (Tokenization, Stopwords Elimination, and Non-Alphanumeric Words Filtering):

$$UC_{Score_{it}} = 1 - \frac{\text{Card}(S_{it})}{\text{Card}(C_i)} = \frac{\text{Card}(C_i \setminus S_{it})}{\text{Card}(C_i)}$$

$$S_{it} = \{w \in C_i | w \in K_{it}\}$$

Where:

- $C_i$  is the corpus of vocabulary containing words related to SASB's target metrics for the industry of company  $i$ .
- $K_{it}$  is the set of the treated keywords in the CDP disclosure (constructed using TF-IDF Vectorization) of the company  $i$  for year  $t$ .



- $S_{it}$  is the set of corpus words semantically existing in the CDP text of the company  $i$  for year  $t$  (identified using a Word Frequency Distribution model FreqDist).

Let  $D$  be the set of the text documents for each answer to the CDP questions for a given company-year pair. To build the  $K_{it}$  set, we use TF-IDF Vectorization method which consists of building a matrix where the rows are the documents contained in  $D$ , the columns are all the words contained in all documents, and the entries for the matrix are the TF-IDF scores for  $i$ -th word in the  $j$ -th document in  $D$ , and then pick for each document the top words that yield the highest scores.

TF-IDF is a numerical statistic in NLP that is intended to reflect the importance of a word in a document relative to a corpus of documents. It is a product of two statistics, term frequency (TF) and inverse document frequency (IDF). The TF-IDF score for a term  $t$  in document  $d \in D$  is given by:

$$\text{TF-IDF}(t, d, D) = \text{TF}(t, d) \times \text{IDF}(t, D) = \frac{\text{Card}(\{w \in d \mid w = t\})}{N_d} \times \log\left(\frac{N}{\text{Card}(\{d \in D \mid t \in d\})}\right)$$

Where:

- $N_d$  is the total number of terms in document  $d$ .
- $N$  is the total number of documents in the corpus  $D$ .
- $\{d \in D \mid t \in d\}$  is the set of documents in  $D$  where the term  $t$  appears.
- $\{w \in d \mid w = t\}$  is the set of occurrences of term  $t$  in document  $d$ .

After aggregating the set  $K_{it}$ , we estimate Word Frequency Distribution to identify words from the corpus  $C_i$  occurring in it and store them in the set  $S_{it}$  to derive the final  $UC_{Score}$ .

A low  $UC_{Score}$  (closer to 0) signifies a strong correlation between the text and the corpus, indicating a thorough and accurate representation of environmental efforts relative to industry standards. Conversely, a high  $UC_{Score}$  (closer to 1) suggests potential greenhushing behavior, where crucial environmental information may be under-communicated or omitted.

#### 4.2.4. Deriving and Exploring the Final Greenhushing Measure:

After averaging the three pillars of greenhushing behavior, we derive the Greenhushing Score  $GH_{Score}$ , which serves as a quantitative measure of the intensity of greenhushing behavior among companies. This metric allows us to assess and compare the extent to which companies are downplaying or underreporting their environmental initiatives and performance.

When analyzing the  $GH_{Score}$  distributions over the years, we observe that they predominantly exhibit a bell curve shape. This bell curve distribution suggests that

most companies have moderate  $GH_{score}$  values, with fewer companies exhibiting extremely high or low scores.

A significant trend observed is the decline in the average  $GH_{score}$  from 0.59 in 2017 to 0.51 in 2021, which represents an almost 14% drop. This downward trend indicates that, on average, companies have been increasingly transparent about their environmental practices over the years.

Moreover, we notice a thinning of the right tail of the  $GH_{score}$  distribution over the years. This suggests a decrease in the number of companies with remarkably high  $GH_{score}$ , implying that fewer companies are engaging in extreme greenhushing behaviors. This trend may reflect growing pressure from stakeholders for greater environmental transparency and accountability, as well as the increasing importance of sustainability in corporate governance.

Another important variable derived from the  $GH_{score}$  is a binary indicator that signifies whether the  $GH_{score}$  exceeds 0.5. This threshold is crucial as it marks the 50% discrepancy level, allowing us to identify instances of significant greenhushing behavior.

Encouragingly, the trend of the number of companies exhibiting significant greenhushing behavior shows a clear downward trajectory, declining from just under 360 companies in 2017 to a little over 270 companies in 2021 as shown in Figure 17. This decline aligns with the observed decrease in the mean  $GH_{score}$  and the thinning of the distribution's right tail over the years. These trends collectively suggest a reduction in the intensity and prevalence of greenhushing behaviors among companies, reflecting a positive shift towards greater environmental transparency.

However, despite this encouraging trend, the fact that over 270 corporations continue to engage in major greenhushing behavior as of 2021 is troubling. This figure emphasizes the persistent problem of establishing full openness in corporate environmental practices, as well as the importance of continuing efforts to encourage and enforce honest reporting by all corporations. The presence of many corporations with considerable greenhushing behavior suggests that, while progress has been made, there is still much space for improvement in corporate sustainability reporting and responsibility.

## Conclusion

In this section, we have outlined our solution design for measuring excess CEO compensation and greenhushing, through a blend of econometric estimation and machine learning techniques. By quantifying these dimensions of corporate behavior, we have laid a solid foundation for understanding their impact on corporate governance and sustainability practices. Moving forward, our analysis will extend

beyond mere measurement; we will explore the intricate relationships between these variables using econometric models in the following section.

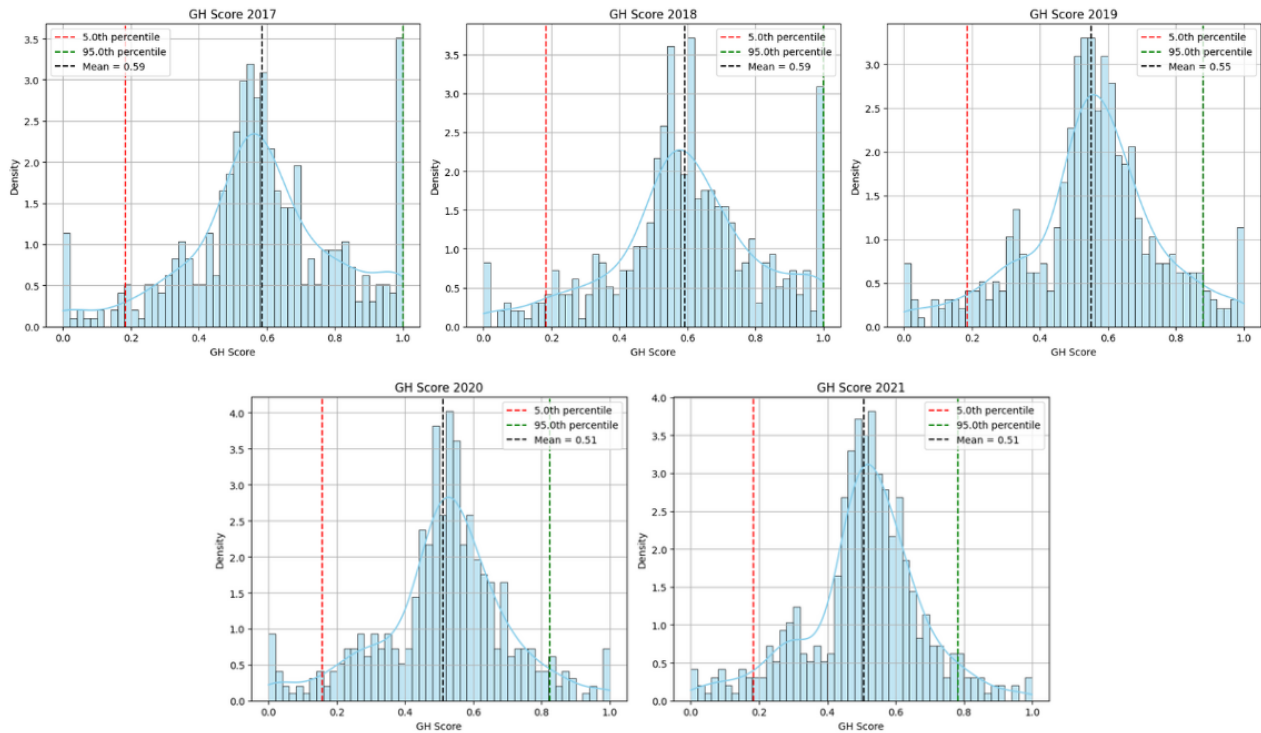


Figure 16:  $GH_{Score}$  Distributions over the years.

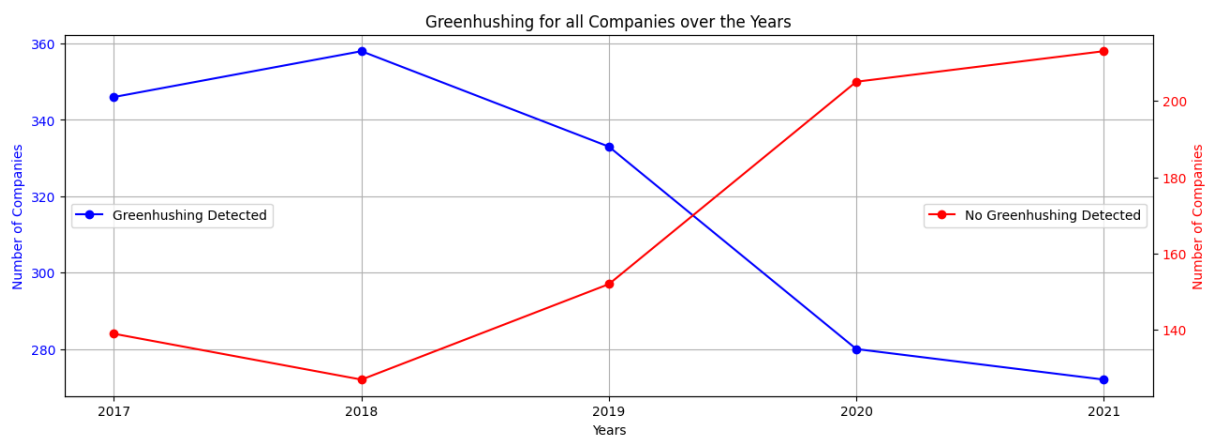


Figure 17: Number of Companies involved in Greenhushing Behavior Curve

## Chapter 5. Exploring the Relationship Between Greenhushing & Excess Compensation

### Introduction

In this section, we delve into the relationship between greenhushing and excess compensation. We introduce two models—one focusing on greenhushing's existence, the other on its intensity within corporate disclosures. Through regression analysis, we explore the influence of greenhushing on excess compensation along with other factors to test the hypotheses developed earlier.

### 5.1. Econometric Models

In our econometric model, we aim to comprehensively explore the intricate relationship between greenhushing behavior and excess compensation within corporations. By focusing into this relationship, we seek to unravel the potential implications of greenhushing on executive pay structures and corporate governance practices.

At the heart of our analysis lies the dependent variable, excess compensation  $Xcomp$ , which serves as a crucial indicator of the remuneration received by executives beyond what is deemed reasonable or justified. This excess compensation is influenced by a myriad of factors, including the extent of greenhushing behavior exhibited within the organization.

To capture the nuanced nature of greenhushing, we introduce the Greenhushing Score  $GH_{score}$  as an independent variable. This score quantifies the degree of greenhushing behavior present within the corporate environment, with a value of 0 indicating no greenhushing and a value of 1 representing total greenhushing. Through this measure, we aim to shed light on how the concealment or downplaying of environmental initiatives may impact executive compensation practices.

Furthermore, we include the internal governance  $GOV_{it}$  of firms using proxies represented by Board Size  $BS$  and Board Independence  $BI$  which is the proportion of the independent directors. These governance variables provide insights into the composition and independence of the board of directors, offering valuable context for understanding the dynamics at play within the organization.

$$GOV_{it} = \begin{pmatrix} BS_{it} \\ BI_{it} \end{pmatrix}$$

The same Economic determinants  $ED$  mentioned earlier in the estimation of the dependent variable are also incorporated into our model. These variables serve as

proxies for firm performance and growth opportunities, which may influence executive compensation decisions.

$$ED_{it} = \begin{pmatrix} B2M_{it} \\ \ln(Tenure)_{it} \\ \ln(Sales)_{it-1} \\ RET_{it} \\ RET_{it-1} \\ ROA_{it} \\ ROA_{it-1} \end{pmatrix}$$

Additionally, we consider the CEOs' Power Status  $CEOchar_{it}$  which is an indicator of the CEO's power position within the company made of two binary variables indicating whether the CEO was a past employee  $Power1_{it}$  at the company and if he is a board member or not  $Power2_{it}$ . This aspect of the model allows us to explore how the power dynamics within corporations may interact with greenhushing behavior to shape executive pay structures.

$$CEOchar_{it} = \begin{pmatrix} Power1_{it} \\ Power2_{it} \end{pmatrix}$$

Utilizing a linear regression framework, we formulate two distinct models to probe the relationship between greenhushing and excess compensation.

### 5.1.1. The First Model: Studying the Impact of the Presence of Greenhushing on Excess Executive Compensation

The first model is expressed as:

$$Xcomp_{it} = \beta_0 + \beta_1 \cdot \mathbb{1}_{GHscore > 0.5_{it}} + \beta_2 \cdot GOV_{it} + \beta_3 \cdot ED_{it} + \beta_4 \cdot CEOchar_{it} + \epsilon_{it}$$

It was designed to investigate the impact of the presence of greenhushing behavior on excess compensation. Here,  $Xcomp_{it}$  represents excess compensation,  $\mathbb{1}_{GHscore > 0.5_{it}}$  denotes the presence of greenhushing behavior,  $GOV_{it}$  represents internal governance proxies,  $ED_{it}$  encompasses economic determinants,  $CEOchar_{it}$  indicates CEOs' power status and  $\epsilon_{it}$  is the error term.

### 5.1.2. The Second Model: Studying the Impact of the Greenhushing Intensity on Excess Executive Compensation

The second model, expressed as:

$$Xcomp_{it} = \beta_0 + \beta_1 \cdot GHscore_{it} + \beta_2 \cdot GOV_{it} + \beta_3 \cdot ED_{it} + \beta_4 \cdot CEOchar_{it} + \epsilon_{it}$$

studies the extent of greenhushing behavior and its influence on executive pay. In this model,  $GHscore_{it}$  quantifies the intensity of greenhushing behavior, while the remaining variables retain their previous definitions. Through these models, we aim to

elucidate the intricate relationship between greenhushing behavior and excess compensation within corporate settings.

## 5.2. Results & Analysis

### 5.2.1. The Relationship Between Excess Executive Compensation & the Existence of Greenhushing

The econometric results of linking excess CEO compensation and the existence of greenhushing behavior reveal nuanced insights as displayed in Table 3, highlighting the dynamic interplay of several numerical figures in that relationship. This analysis incorporates three primary models: the Pooled OLS, Fixed Effects, and Random Effects models.

Table 3: Regression Table for the First Excess Compensation & Greenhushing Model<sup>10</sup>

The table presents regression results for three econometric specifications for the first model: Pooled OLS, Fixed Effects, and Random Effects. The dependent variable is excess compensation. The first independent variable is greenhushing presence (binary), followed by two internal governance variables. The next seven variables are economic determinants, and the final two variables serve as proxies for CEO's power. The table includes coefficient estimates, t-values, and significance levels for each variable across the models.

Dependent Variable: $Xcomp_{it}$			
Independent Variables	Pooled OLS Model	Fixed Effects Model	Random Effects Model
$\mathbb{1}_{GH_{Score} > 0.5}_{it}$	0.0298*** (2.6541)	-0.0028 (-0.2240)	0.0202* (1.7666)
$BS_{it}$	0.1208*** (4.9896)	0.0797** (2.2540)	0.1178*** (4.1207)
$BI_{it}$	0.0338* (1.6892)	-0.0060 (-0.1932)	0.0420* (1.7313)
$B2M_{it}$	0.1577*** (6.5575)	0.0106 (0.2074)	0.1341*** (4.1940)
$\ln(Tenure)_{it}$	0.1112*** (4.5152)	-0.0075 (-0.2661)	0.0693*** (2.8188)
$\ln(Sales)_{it-1}$	-0.6717*** (-23.557)	-0.5631*** (-3.8349)	-0.5952*** (-13.532)
$RET_{it}$	0.1191*** (4.4592)	-0.0167 (-0.7056)	0.0608*** (2.7098)
$RET_{it-1}$	0.1341*** (5.2253)	0.0023 (0.1050)	0.0762*** (3.5868)
$ROA_{it}$	0.1065** (1.9818)	0.0347 (0.7263)	0.1313*** (2.9733)
$ROA_{it-1}$	0.0856* (1.7212)	-0.0079 (-0.1675)	0.0800** (1.9713)
$Power1_{it}$	0.0071 (0.4625)	0.0359 (1.2823)	0.0095 (0.4677)

<sup>10</sup> The table displays the coefficient and Student's T (between parenthesis) for each variable in the three models. \*, \*\*, and \*\*\* indicate two-tailed statistical significance at 0.1, 0.05 and 0.01 levels respectively.

<i>Power2<sub>it</sub></i>	0.4005*** (13.116)	0.0321 (0.5655)	0.4369*** (12.763)
----------------------------	-----------------------	--------------------	-----------------------

In the Pooled OLS model, several significant insights emerge. Firstly, the existence of greenhushing positively impacts excess CEO compensation, with the coefficient significant at the 1% level. This finding supports Hypothesis 2a, suggesting that greenhushing behavior drives higher excess compensation, thereby rejecting Hypothesis 1, which posits no relationship. Furthermore, internal governance, specifically Board Size, significantly and positively impacts excess compensation at the 1% level, opposing Hypothesis 3, which anticipated a negative impact. Board Independence is also significant at 10% level, further opposing the latter hypothesis. Additionally, the CEO power status variable, indicating the CEO's role as a board member *Power2* (measured by whether the CEO is also a board member), significantly and positively impacts excess compensation at 1% level, confirming Hypothesis 4. Moreover, the economic determinants show a significant negative impact on excess compensation through the logarithm of lagged sales at the 1% level and positively impact it through the book to market ratio at 1%, the logarithm of tenure at 1%, with return and its first lag significant at 1%, and ROA and its first lag are at levels 5% and 10% respectively. This adds another layer of complexity to the analysis.

On the other hand, the Fixed Effects model presents a distinct perspective. In this model, greenhushing is not a significant predictor of excess CEO compensation. This result rejects Hypotheses 2a and 2b and confirms Hypothesis 1, suggesting that greenhushing behavior does not influence excess compensation. Additionally, internal governance, specifically Board Size, significantly and positively impacts excess compensation at the 5% level, opposing Hypothesis 3, which anticipated a negative impact. However, Board Independence is insignificant. Moreover, CEO power status variables are also not significant, rejecting Hypothesis 4. However, it is noteworthy that economic determinants significantly impact excess compensation only through the logarithm of lagged sales negatively at the 1% level.

Meanwhile, the Random Effects model offers insights consistent with the Pooled OLS model regarding the positive impact of greenhushing existence on excess CEO compensation. In this model, greenhushing's binary presence is significant at the 10% level, again confirming Hypothesis 2a and rejecting Hypothesis 1. For internal governance variables, Board Size once again shows a significant positive impact on excess compensation at the 1% level, contradicting Hypothesis 3, while Board Independence is also significant at 10% level. Additionally, CEO power status, particularly the *Power2* variable, significantly and positively impacts excess compensation at the 1% level, thereby confirming Hypothesis 4. Furthermore, the economic determinants show a significant negative impact on excess compensation through the book to market ratio at 1%, the logarithm of lagged sales at the 1% level and positively impact it through the logarithm of tenure at 1%, with return and its first lag significant at 1%, and ROA and its



first lag are at levels 1% and 5% respectively. Thus, the Random Effects model reinforces the complexity observed in the other models.

In evaluating the relationship between excess executive compensation and the binary presence of greenhushing, the residual analysis and goodness of fit provide critical insights into the model performances. According to the Fischer test at the 5% significance level, the Pooled OLS, Random Effects, and Fixed Effects models are all globally significant. However, Fixed Effects model has a P-value that is very close to 5%. This indicates that the Pooled OLS and Random Effects models provide a better overall fit compared to the Fixed Effects model.

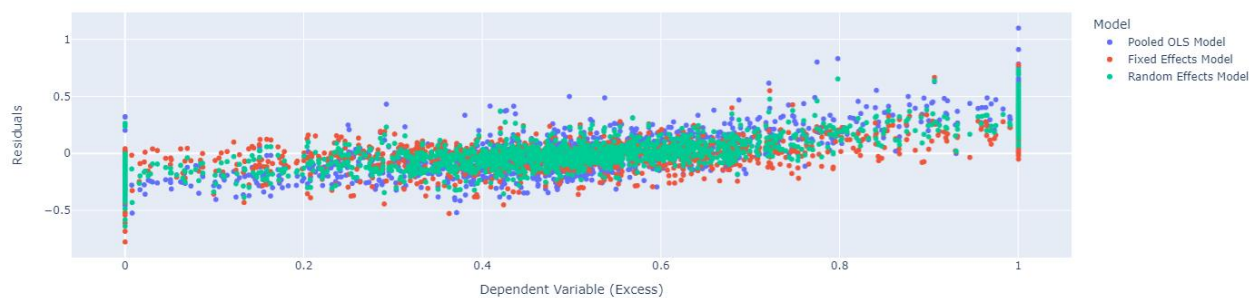


Figure 18: Residuals vs Dependent Variable Plot for the First Model of Excess Compensation and Greenhushing

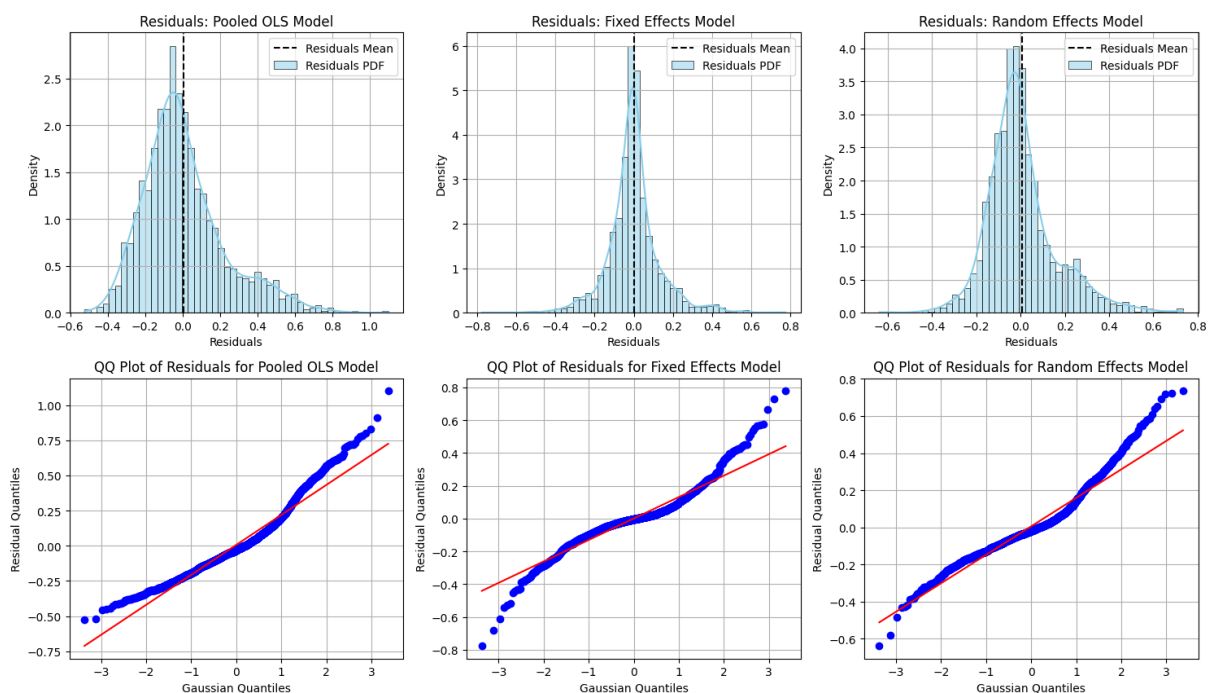


Figure 19: Distributions & Q-Q Plots of Residuals for the First Model of Excess Compensation and Greenhushing

All three models—Pooled OLS, Fixed Effects, and Random Effects—exhibit some linear correlation between the residuals and the dependent variable, excess compensation (Figure 18). The distribution of residuals for all three models is not normal, as evidenced by the Shapiro-Wilk, Anderson-Darling, and Kolmogorov-Smirnov tests. These tests indicate that the residuals revolve around means that are



relatively close to zero. The Q-Q plots reveal that the Pooled OLS model yields residuals with quantiles closest to normal distribution, followed by the Random Effects model, and then the Fixed Effects model (Figure 19). This normality proximity in residuals further supports the robustness of the Pooled OLS model.

Analyzing the autocorrelation function of residuals, both the Fixed Effects and Pooled OLS models exhibit autocorrelation, with positive autocorrelation in the Pooled OLS model and negative autocorrelation in the Fixed Effects model. These findings are confirmed by the Durbin-Watson test variant for large samples. For the Random Effects model, the test is indecisive; however, the ACF plot indicates no observable autocorrelation (Figure 20).

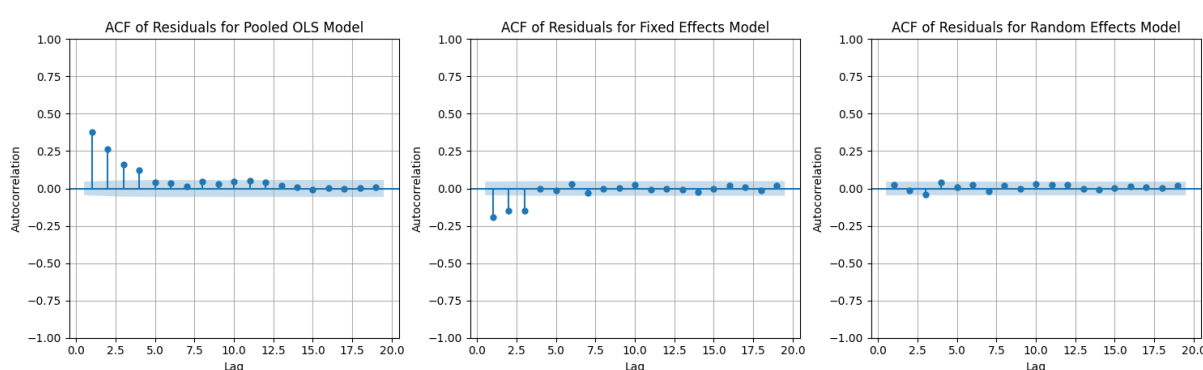


Figure 20: Autocorrelation Function Plots of the First Excess Compensation and Greenhushing Model

When comparing model fit metrics, the Random Effects model stands out by having higher Log-Likelihood, and lower AIC and BIC values than Pooled OLS model, by not showing any residual autocorrelation, and by minimizing the symmetric Mean Absolute Percentage Error (Table 4), suggesting superior predictive accuracy.

Table 4: Goodness of Fit Metrics Table for all Specification of the First Model<sup>11</sup>

The table summarizes the performance metrics of three econometric models: Pooled OLS, Fixed Effects, and Random Effects. The metrics included are Log-Likelihood, Akaike Information Criterion (AIC), Bayesian Information Criterion (BIC), and symmetric Mean Absolute Percentage Error (sMAPE).

Metric	Pooled OLS Model	Fixed Effects Model	Random Effects Model
<b>Log Likelihood</b>	186.15	1059.50**	782.01
<b>BIC</b>	-282.09	-2028.78*	-1473.81
<b>AIC</b>	-348.30	-2095.00*	-1540.02
<b>sMAPE</b>	0.55	2.00	0.54*

<sup>11</sup> \*\* highlights the maximum value for Log Likelihood. \* points out the lowest value for AIC, BIC and sMAPE.

### 5.2.2. The Relationship Between Excess Executive Compensation & the Intensity of Greenhushing

Similar to the previous sub-subsection, we analyze the econometric findings of linking excess CEO compensation to the intensity of greenhushing. This analysis also involves the same three primary models: the Pooled OLS, Fixed Effects, and Random Effects models, each offering unique perspectives.

Table 5: Regression Table for the Second Excess Compensation & Greenhushing Model<sup>12</sup>

The table presents regression results for three econometric specifications for the second model: Pooled OLS, Fixed Effects, and Random Effects. The dependent variable is excess compensation. The first independent variable is greenhushing intensity, followed by two internal governance variables. The next seven variables are economic determinants, and the final two variables serve as proxies for CEO's power. The table includes coefficient estimates, t-values, and significance levels for each variable across the models.

Dependent Variable: $Xcomp_{it}$			
Independent Variables	Pooled OLS Model	Fixed Effects Model	Random Effects Model
$GH_{Score_{it}}$	0.0847*** (3.8373)	-0.0575* (-1.7730)	0.0665*** (2.6007)
$BS_{it}$	0.1223*** (5.0687)	0.0761** (2.1532)	0.1198*** (4.2013)
$BI_{it}$	0.0360* (1.8010)	-0.0054 (-0.1754)	0.0411* (1.6971)
$B2M_{it}$	0.1525*** (6.3366)	0.0069 (0.1354)	0.1312*** (4.1063)
$ln(Tenure)_{it}$	0.1011*** (4.0762)	-0.0081 (-0.2869)	0.0656*** (2.6604)
$ln(Sales)_{it-1}$	-0.6588*** (-22.851)	-0.5599*** (-3.8200)	-0.5874*** (-13.320)
$RET_{it}$	0.1106*** (4.1229)	-0.0149 (-0.6336)	0.0568** (2.5197)
$RET_{it-1}$	0.1277*** (4.9657)	0.0008 (0.0341)	0.0741*** (3.4863)
$ROA_{it}$	0.1079** (2.0138)	0.0331 (0.6917)	0.1302*** (2.9511)
$ROA_{it-1}$	0.0780 (1.5699)	-0.0087 (-0.1868)	0.0745* (1.8352)
$Power1_{it}$	0.0095 (0.6183)	0.0381 (1.3672)	0.0113 (0.5575)
$Power2_{it}$	0.3876*** (12.573)	0.0301 (0.5305)	0.4229*** (12.105)

In the Pooled OLS model, several significant relationships emerge. Firstly, the intensity of greenhushing has a positive and significant impact on excess CEO compensation at the 1% level. This result confirms Hypothesis 2a, indicating that greater

<sup>12</sup> The table displays the coefficient and Student's T (between parenthesis) for each variable in the three models. \*, \*\*, and \*\*\* indicate two-tailed statistical significance at 0.1, 0.05 and 0.01 levels respectively.

greenhushing intensity drives higher excess compensation, thereby rejecting Hypothesis 1, which posits no relationship. Furthermore, internal governance factors play a crucial role, with Board Size significantly and positively impacting excess compensation at the 1% level. This finding contradicts Hypothesis 3, which anticipated a negative impact. Board Independence is also significant at 10% level, further opposing the latter hypothesis. Additionally, the CEO power status variable, particularly the indicator of the CEO being a board member *Power2*, significantly and positively impacts excess compensation at the 1% level, confirming Hypothesis 4. Moreover, the economic determinants show a significant negative impact on excess compensation through the logarithm of lagged sales at the 1% level and positively impact it through the book to market ratio at 1%, the logarithm of tenure at 1%, with return and its first lag significant at 1%, and ROA at level 5%. This adds another layer of complexity to the analysis.

Conversely, the Fixed Effects model presents a different scenario. Here, the intensity of greenhushing negatively impacts excess CEO compensation, with the coefficient significant at the 10% level. This finding supports Hypothesis 2b, suggesting that higher greenhushing intensity leads to lower excess compensation, thereby rejecting Hypothesis 1. Additionally, governance variables, specifically Board Size, significantly and positively impacts excess compensation at the 5% level, opposing Hypothesis 3, which anticipated a negative impact. However, Board Independence is insignificant. Moreover, CEO power status variables are also not significant, rejecting Hypothesis 4. Nevertheless, economic determinants significantly impact excess compensation only through the logarithm of lagged sales negatively at the 1% level.

Meanwhile, the Random Effects model offers insights consistent with the Pooled OLS model regarding the positive impact of greenhushing intensity on excess CEO compensation. In this model, greenhushing intensity is significant at the 10% level, again confirming Hypothesis 2a and rejecting Hypothesis 1. internal governance factors play a crucial role, with Board Size significantly and positively impacting excess compensation at the 1% level. This finding contradicts Hypothesis 3, which anticipated a negative impact. Board Independence is also significant at 10% level, further opposing Hypothesis 3. Additionally, the CEO power status, particularly the *Power2* variable, significantly and positively impacts excess compensation at the 1% level, confirming Hypothesis 4. Furthermore, the economic determinants show a significant negative impact on excess compensation through the book to market ratio at 1%, the logarithm of lagged sales at the 1% level and positively impact it through the logarithm of tenure at 1%, with return and its

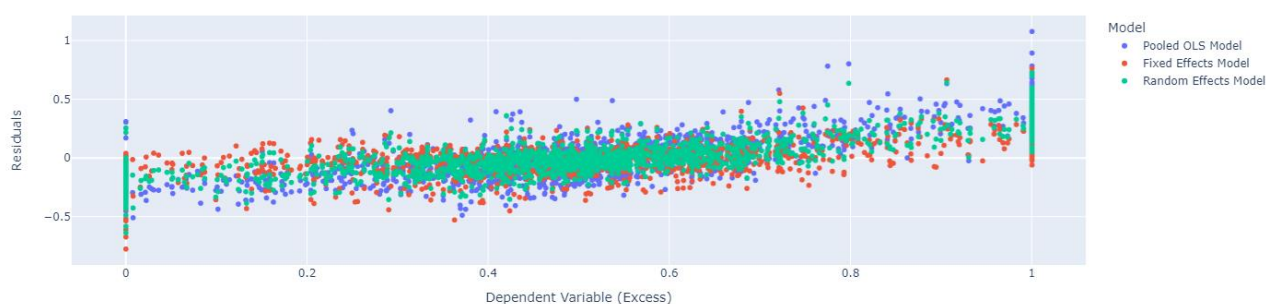


Figure 21: Residuals vs Dependent Variable Plot for the Second Model of Excess Compensation and Greenhushing

first lag significant at 1%, and ROA and its first lag are at levels 1% and 10% respectively. Thus, the Random Effects model reinforces the complexity observed in the other models.

When investigating the association between excess executive salary and the degree of greenhushing, residual analysis and goodness of fit demonstrate the efficacy of several models. At the 5% significance level, the Fischer test indicates that the Pooled OLS, Random Effects, and Fixed Effects models are all globally significant. This suggests that each model gives a good match for the data in this case.

All three models—Pooled OLS, Fixed Effects, and Random Effects—show a linear relationship between the residuals and the dependent variable, excess compensation (Figure 21).

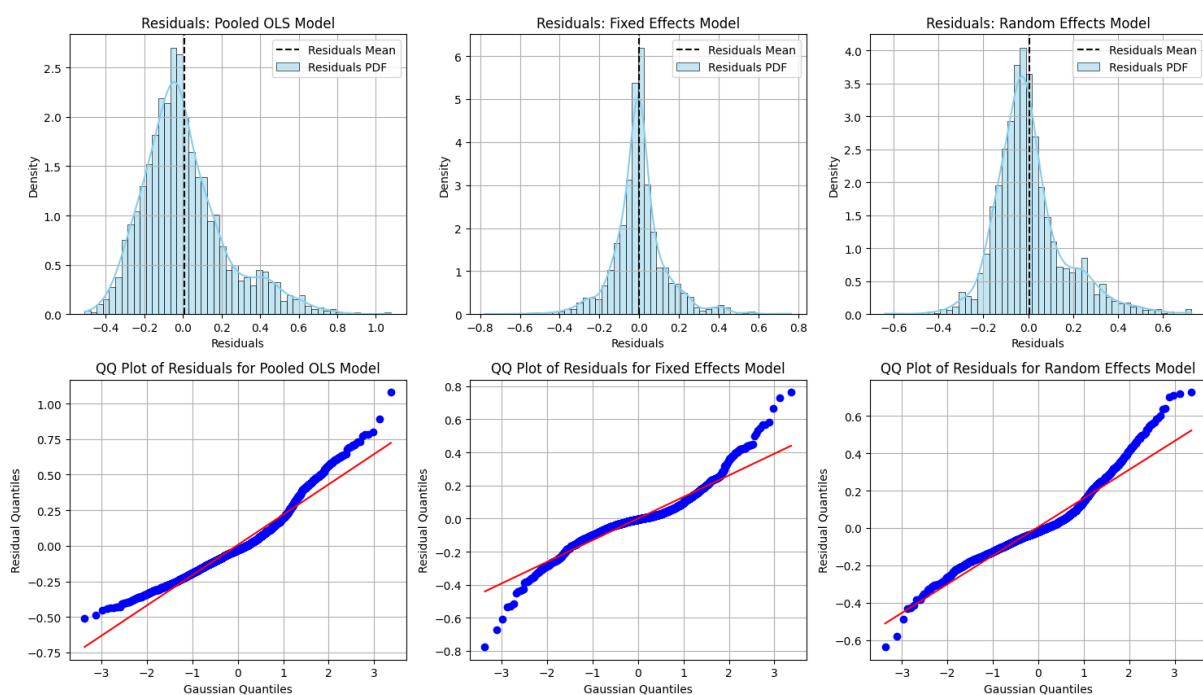


Figure 22: Distributions & Q-Q Plots of Residuals for the Second Model of Excess Compensation and Greenhushing

The Shapiro-Wilk, Anderson-Darling, and Kolmogorov-Smirnov tests show that the residual distributions in all three models are not normal. These tests demonstrate that the residuals are clustered around values relatively close to zero. The Q-Q plots show that the Pooled OLS model produces residuals with quantiles closest to a normal distribution, followed by the Random Effects and Fixed Effects models (Figure 22). This implies that the Pooled OLS model provides the best approximation of normal residuals of the three.

When examining the autocorrelation function of residuals, both the Fixed Effects and Pooled OLS models show autocorrelation—positive in the Pooled OLS model and negative in the Fixed Effects model. The Durbin-Watson test variation for large samples

provides support for these findings. For the Random Effects model, the test results are indecisive, but the ACF plot indicates no observable autocorrelation (Figure 23).

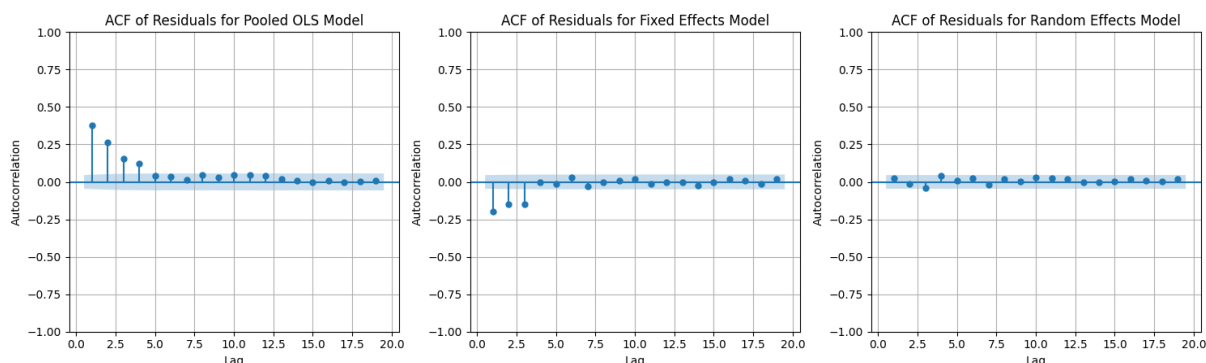


Figure 23: Autocorrelation Function Plots of the Second Excess Compensation and Greenhushing Model

When comparing model fit metrics, the Random Effects model stands out by having higher Log-Likelihood, and lower AIC and BIC values than Pooled OLS model, by not showing any residual autocorrelation, and by minimizing the symmetric Mean Absolute Percentage Error (Table 6), suggesting superior predictive accuracy.

Table 6: Goodness of Fit Metrics Table for all Specification of the Second Model<sup>13</sup>

The table summarizes the performance metrics of three econometric models: Pooled OLS, Fixed Effects, and Random Effects. The metrics included are Log-Likelihood, Akaike Information Criterion (AIC), Bayesian Information Criterion (BIC), and symmetric Mean Absolute Percentage Error (sMAPE).

Metric	Pooled OLS Model	Fixed Effects Model	Random Effects Model
Log Likelihood	189.99	1061.48**	783.56
BIC	-289.77	-2032.75*	-1476.91
AIC	-355.99	-2098.96*	-1543.12
sMAPE	0.55	2.00	0.54*

### 5.3. Comparative Analysis for the Two Studies: Presence vs. Intensity

The econometric analyses of greenhushing's presence and intensity in relation to excess CEO compensation reveal nuanced insights into how various aspects of greenhushing behavior impact executive pay. Here, we compare the results from these analyses across three econometric models: Pooled OLS, Fixed Effects, and Random Effects.

First, examining the effect of greenhushing, both analyses indicate that greenhushing, whether measured as its presence or intensity, significantly influences excess CEO compensation. However, the direction and significance of this influence vary across models. For instance, the presence of greenhushing positively impacts

<sup>13</sup> \*\* highlights the maximum value for Log Likelihood. \* points out the lowest value for AIC, BIC and sMAPE.

excess CEO compensation in the Pooled OLS model, supporting Hypothesis 2a and rejecting Hypothesis 1. Conversely, in the Fixed Effects model, greenhushing presence is not a significant predictor, confirming Hypothesis 1 and rejecting Hypotheses 2a and 2b. In the Random Effects model, greenhushing presence is significant, aligning with the Pooled OLS model. On the other hand, the intensity of greenhushing shows a positive impact on excess CEO compensation in both the Pooled OLS and Random Effects models, supporting Hypothesis 2a and rejecting Hypothesis 1, while it shows a negative impact in the Fixed Effects model, supporting Hypothesis 2b. This variation underscores the complexity of greenhushing's role in executive compensation dynamics.

Second, regarding governance variables, Board Size consistently shows a significant and positive impact on excess compensation when it is significant, regardless of whether the focus is on the presence or intensity of greenhushing. This consistent result opposes Hypothesis 3, suggesting that the presence of the CEO on the board leads to higher excess compensation. However, Board Independence remains, at best, significant at 10% in both analyses, indicating that the proportion of independent directors in the board does not play a substantial role in moderating excess CEO compensation, and in some cases, opposes Hypothesis 3 even further.

Third, CEO Power Status also exhibits consistent patterns. The variable *Power2* (indicating the CEO is a board member) consistently shows a positive and significant impact on excess compensation across the Pooled OLS and Random Effects models in both analyses, thus confirming Hypothesis 4. Conversely, *Power1* (indicating that the CEO is a past employee at the company) is not significant in either analysis, suggesting that other aspects of CEO power may be less influential in determining excess compensation.

Lastly, considering economic determinants, several variables consistently demonstrate significant impacts on excess compensation across models. In the presence model, the logarithm of lagged sales negatively impacts excess compensation at the 1% level in the Pooled OLS, Fixed and Random Effects models. The intensity model similarly shows significant negative impacts of lagged sales on excess compensation at the 1% level in these models. Additionally, stock returns and their first lag, return on assets and its first lag, book to market ratio and the logarithm of tenure show consistent significant positive impacts in Pooled OLS and Random Effects models only. These results highlight the importance of these financial metrics in moderating executive excess pay, irrespective of greenhushing measures.

## Conclusion

In conclusion, while both the presence and intensity of greenhushing are significant determinants of excess CEO compensation, the specifics of their impacts

differ depending on the econometric model employed. Internal governance and CEO power status also play crucial roles, with consistent findings across different measures of greenhushing. The results underscore the complexity and multifaceted nature of the factors influencing executive compensation, illustrating that both the structural attributes of governance and the behavioral tendencies related to greenhushing must be considered to fully understand the dynamics of excess CEO pay.



## General Conclusion & Future Works

The purpose of this work is to examine the relationship between greenhushing and excess executive compensation by leveraging Natural Language Processing techniques to analyze corporate disclosures and employing quantitative approaches to quantify greenhushing and link it econometrically to excess pay along with other variables. We provide strong empirical evidence that both the presence and intensity of greenhushing significantly influence excess CEO compensation.

Firstly, examining the existence of greenhushing, we find that it positively impacts excess CEO compensation, with the coefficient being significant. This finding supports the notion that greenhushing behavior drives higher excess compensation, thereby rejecting the hypothesis that greenhushing does not drive excess executive compensation. Furthermore, internal governance factors, significantly and positively impact excess compensation via Board Size, opposing the hypothesis that stronger internal governance has a negative impact on excess compensation. Additionally, the CEO's power status variable, indicating the CEO's role as a board member, significantly and positively impacts excess compensation, confirming that greater CEO power strengthens the relationship between greenhushing behavior and excess executive compensation. Moreover, the economic determinants show a significant negative impact on excess compensation through the book to market ratio, the logarithm of lagged sales and positively impact it through the logarithm of tenure, with stock return and its first lag, and the return on assets and its first lag. This adds another layer of complexity to the analysis.

Similarly, when we examine the severity of greenhushing, we see that it has a positive and significant effect on excess CEO salary. This finding demonstrates that more greenhushing intensity leads to higher excess compensation, refuting the theory that greenhushing does not drive excess executive remuneration. Internal governance characteristics continue to play an important effect, with Board Size having a considerable and positive impact on excess remuneration, contradicting the expectation that stronger internal governance reduces excess compensation. Furthermore, the CEO power status variable, namely the CEO's status as a board member, has a significant and positive impact on excess remuneration, confirming that increased CEO power improves the association between greenhushing conduct and excess executive compensation. Economic determinants reveal that logarithmic lagged sales negatively affect excess compensation and positively affect it the rest of determinants of CEO pay, adding further complexity to the model.

These findings suggest a number of implications for practitioners. First, it is crucial for boards of directors, particularly compensation committees, to consider the impact of greenhushing when designing CEO compensation plans. The presence and



intensity of greenhushing should be factored into decisions about executive pay to prevent CEOs from gaining excess compensation through under-communication of their companies' sustainability efforts. Second, enhancing transparency and reducing greenhushing could mitigate the risks associated with excessive executive pay and align CEO incentives with long-term corporate sustainability goals.

Finally, we recommend further research to expand the sources of environmental disclosure beyond the CDP for a broader analysis using text mining on various disclosures, enlarge the sample size in terms of both timeframe and the number of individuals, and improve econometric model specifications to explore additional variables. Additionally, enhancing our solution for Greenhushing measuring with better embedding algorithms such as ClimateBERT, a BERT model fine-tuned to environmental issues, could potentially offer better insights and improved accuracy in understanding and analyzing environmental disclosures.

## Bibliography

Bajic, A. (2023). "Climate Disclosure: A Machine Learning-Based Analysis of Company-Level GHG Emissions and ESG Data Disclosure." Available at SSRN: <https://ssrn.com/abstract=4534697>

Barontini, R., & Hill, J. G. (2023). "Sustainability and Executive Compensation." ECGI Law Working Paper N° 747/2023.

Carlos, K., & Lewis, B. R. (2018). "Strategic Silence: Withholding Certification Status as a Hypocrisy Avoidance Tactic." *Administrative Science Quarterly*, 63(4), 683-722.

Core, J. E., Guay, W., & Larcker, D. F. (2008). "The power of pen and executive compensation." *Journal of Financial Economics*, 88, 1-25.

Ettinger, A., Grabner-Kräuter, S., Okazaki, S., & Terlutter, R. (2021). "The Desirability of CSR Communication versus Greenhushing in the Hospitality Industry: The Customers' Perspective." *Journal of Travel Research*, 60(3), 619-638.

Font, X., Elgammal, I., & Lamond, I. (2017). "Greenhushing: the deliberate under communicating of sustainability practices by tourism businesses." *Journal of Sustainable Tourism*, 25(7), 1007-1023.

Fried, J., & Shilon, N. (2011). "Excess-pay clawbacks." *The Journal of Corporation Law*, 36(4), 722-741.

Gez, M., Anagnosti, E., & Pullins, T. (2022). "ESG Disclosure Trends in SEC Filings." *Harvard Law School Forum on Corporate Governance*.

Hirshleifer, D., & Thakor, A. (1992). "Managerial conservatism, project choice and debt." *The Review of Financial Studies*, 5, 437-470.

Huang, A., & Li, Y. (2022). "Disclosure Sentiment: Machine Learning vs. Dictionary Methods." *Management Science*, 68(7), 5313-5337. <https://doi.org/10.1287/mnsc.2021.4181>

Jafri, S. R. A., & Trabelsi, S. (2014). "Managerial Risk-Taking and CEO Excess Compensation." Unpublished manuscript, Goodman School of Business, Brock University, St. Catharines, Ontario.

Jelodar, H., Wang, Y., Yuan, C., Feng, X., Jiang, X., Li, Y. and Zhao, L. (2018). "Latent Dirichlet Allocation (LDA) and Topic modeling: models, applications, a survey." *arXiv preprint arXiv:1711.04305v2*.

Lee, M-Y., (2016). "Durbin Watson statistic based on a Z-test in large samples." *Int. J. Computational Economics and Econometrics*, Vol. 6, No. 1, pp.114-121.

Pathan, S. (2009). "Strong boards, CEO power and bank risk-taking." *Journal of Banking and Finance*, 33, 1340-1350.

Simão, L., & Lisboa, I. (2023). "Why companies might under-communicate their efforts for sustainable development and what can be done."

Wang, H., Jia, M., & Zhang, Z. (2021). "Good deeds done in silence: Stakeholder management and quiet giving by Chinese firms." *Organization Science* 32(3), 649-674.

Wu, Y., & Tham, J. (2023). "The Impact of Executive Green Incentives and Top Management Team Characteristics on Corporate Value in China: The Mediating Role of Environment, Social and Government Performance." *Sustainability*, 15(15), 12518. <https://doi.org/10.3390/su151612518>

Zheng, J., Khurram, M.U., & Chen, L. (2022). "Can Green Innovation Affect ESG Ratings and Financial Performance? Evidence from Chinese GEM Listed Companies." *Sustainability*, 14(14), 8677. <https://doi.org/10.3390/su14148677>