



OPEN

# A dataset on corporate sustainability disclosure

DATA DESCRIPTOR

Jinfang Tian<sup>1,5</sup>, Qian Cheng<sup>1,5</sup>, Rui Xue<sup>2,5</sup>✉, Yilong Han<sup>3</sup>  & Yuli Shan<sup>4</sup> ✉

Enterprises, as key emitters, play a vital role in promoting sustainable development. Corporate sustainability disclosure provides a key channel for stakeholders to gain insights into a company's sustainability progress. However, few studies have been conducted to measure sustainability disclosure at the firm level. In this study, we apply the machine learning techniques to listed companies' management discussion and analysis (MD&A) documents and construct a dataset on corporate sustainability disclosure, including the *Corporate Sustainability Disclosure Index* (CSDI), *CSDI\_Economic Dimension* (CSDI\_ECO), *CSDI\_Environmental Dimension* (CSDI\_ENV), and *CSDI\_Social Dimension* (CSDI\_SOCI). The dataset will be updated annually. To the best of our knowledge, this is the first sustainability disclosure dataset constructed at the firm level. Our dataset reflects corporate managements' sustainability attitudes and promotes the implementation of corporate sustainability strategies and subsequent sustainable economic and social outcomes.

## Background & Summary

The fulfilment of sustainable development goals is a profound issue in today's economic and social development<sup>1</sup>. Corporate sustainable development helps promote the sustainable development of the economy and the society<sup>2,3</sup>. The disclosure of corporate sustainability related information thus can deliver key practices and performances that firms have contributed to the sustainable development. Accordingly, quantifying corporate sustainable development information can inform shareholders about how firms invest in sustainable transition and provide stakeholders with measurable quantitative benchmarks<sup>4</sup>, motivating firms to make feasible sustainability strategies and take active actions towards sustainable transition.

Existing research on sustainability measurement has largely focused on national and regional levels<sup>5–9</sup>, while research on sustainability at the firm level has remained at the qualitative level<sup>10</sup>. Studies on quantitative measures of corporate sustainability disclosure remain scarce and await empirical investigation.

As stakeholders' demand for corporate sustainability disclosure increases, more and more international organisations are providing guidelines for corporate sustainability disclosure<sup>11,12</sup>. The Sustainability Reporting Guidelines, developed by the Global Reporting Initiative (GRI), are the most popular among companies worldwide<sup>13,14</sup>. The GRI provides a framework for companies to improve the effectiveness of their sustainability practices<sup>15</sup>. It has also defined the 'triple bottom line' (economic, environmental, and social) for corporate sustainability<sup>16</sup>. Thus, constructing a dataset for corporate sustainability disclosure based on the principle of 'triple bottom line' can reflect the status of corporate sustainability initiatives in a relatively authoritative manner. The management discussion and analysis (MD&A) part of a company's annual report reflects the company's current status and strategic decisions, including sustainability information and strategies<sup>17,18</sup>. As such, quantifying the sustainability-related textual information covered in the MD&A documents can help provide insights into the importance placed by corporate management on sustainability strategies and identify a company's sustainability capabilities.

China is the world's largest emerging economy<sup>19</sup>. Meanwhile, it is also the world's largest carbon emitter<sup>20</sup> and faces severe hazards such as environmental degradation<sup>21</sup>. Therefore, China has attached great responsibilities to promote sustainable economic and social development<sup>22,23</sup>. Given that firms with superior sustainability performance are more inclined to disclose sustainability information, quantifying Chinese enterprises' sustainability

<sup>1</sup>Research Center for Statistics and Interdisciplinary Sciences | School of Statistics and Mathematics, Shandong University of Finance and Economics, Jinan, 250014, China. <sup>2</sup>Centre for Corporate Sustainability and Environmental Finance, Department of Applied Finance, Macquarie University, Sydney, NSW, 2109, Australia. <sup>3</sup>School of Economics and Management, Tongji University, Shanghai, 200092, China. <sup>4</sup>School of Geography, Earth and Environmental Sciences, University of Birmingham, Birmingham, B15 2TT, UK. <sup>5</sup>These authors contributed equally: Jinfang Tian, Qian Cheng, Rui Xue. ✉e-mail: [rui.xue@mq.edu.au](mailto:rui.xue@mq.edu.au); [y.shan@bham.ac.uk](mailto:y.shan@bham.ac.uk)

information disclosure helps reflect and monitor the actual status of their sustainable development. Moreover, it contributes to the timely achievement of the ‘Double Carbon’ strategy<sup>24</sup>.

In this study, we follow the methods of Li *et al.*<sup>25</sup> about quantifying corporate culture based on latest machine learning techniques and Zhang *et al.*<sup>26</sup> about constructing a dictionary of environmental characteristics to measure corporate sustainability disclosure. Specifically, we first refer to Zhang *et al.*'s approach to identify sustainability related seed words and then follow Li *et al.*'s approach to process text documents and expand seed words to find their synonyms based on *word2vec* algorithm. Next, consistent with Li *et al.*, we apply the *tf.idf* weighting scheme to assign weights (importance) to each word and use the weighted sum of all words to calculate the final corporate sustainability information disclosure index of a specific company at a given year. Taken together, building on the sustainability dictionary constructed using *word2vec*, we build a sustainability disclosure dataset of listed firms in China using the *tf.idf* weighting scheme, including the aggregate *Corporate Sustainability Disclosure Index* (CSDI), *CSDI\_Economic Dimension* (CSDI\_ECO), *CSDI\_Environmental Dimension* (CSDI\_ENV), and *CSDI\_Social Dimension* (CSDI\_SOCI). We also conduct a series of validity tests. First, we test the accuracy of text mining techniques against manual extractions. Second, we examine whether CSDI, CSDI\_ECO, CSDI\_ENV, and CSDI\_SOCI are significantly correlated with the corresponding (real) outcome performance indicators. The results for the validation tests indicate that the constructed corporate sustainability disclosure dataset is valid and reliable, and can help stakeholders track and monitor the actual status of sustainable development of listed firms, which subsequently regulates corporate operations and ultimately promotes sustainable economic and social development.

Our study makes following contributions. First, building on the sustainability framework of ‘triple bottom line’ defined by the GRI guidelines, we employ text mining techniques to construct a corporate sustainability disclosure dataset. To the best of our knowledge, this is the first dataset on firm-level sustainability disclosure measurement. The dataset of corporate sustainable development information disclosure constructed in this study can be applied to investigate the progress and influences of corporate sustainable development, and provide data resources for promoting quantitative research of corporate sustainability<sup>27</sup>, which improves the efficiency of knowledge generation related to sustainable development and saves the social costs spent on related issues. Our dataset can also help entrepreneurs to better design sustainable development strategies at a lower cost and ultimately promote the achievement of global sustainable development<sup>28</sup>.

Second, we mine and quantify sustainability information derived from corporate MD&A documents. MD&A texts reflect the strategic directions and decisions of corporate management, which are closely associated with the company's sustainability strategies. Thus, our research methodology can be further extended to other texts that reflect executive decisions. For example, at the city level, because government work reports reflect the directions and strategies of government leaders, further studies can examine climate-related textual information in local government work reports to understand local attitudes and initiatives regarding climate governance.

Third, in constructing the ‘Corporate Sustainability Dictionary,’ we expand seed words—words that are closely related to ‘sustainability’—to a larger sustainability dictionary that includes their synonyms based on the *word2vec* technique. This helps avoid the omission of corporate sustainability information in the texts and reduces subjectivity. To more accurately calculate corporate sustainability disclosure indices, we apply the *tf.idf* weight counting scheme, which takes into account the importance of the corporate sustainability-related words in all text corpora<sup>29</sup>. Therefore, the method is able to distinguish between different levels of importance attached to different dimensions of corporate sustainability.

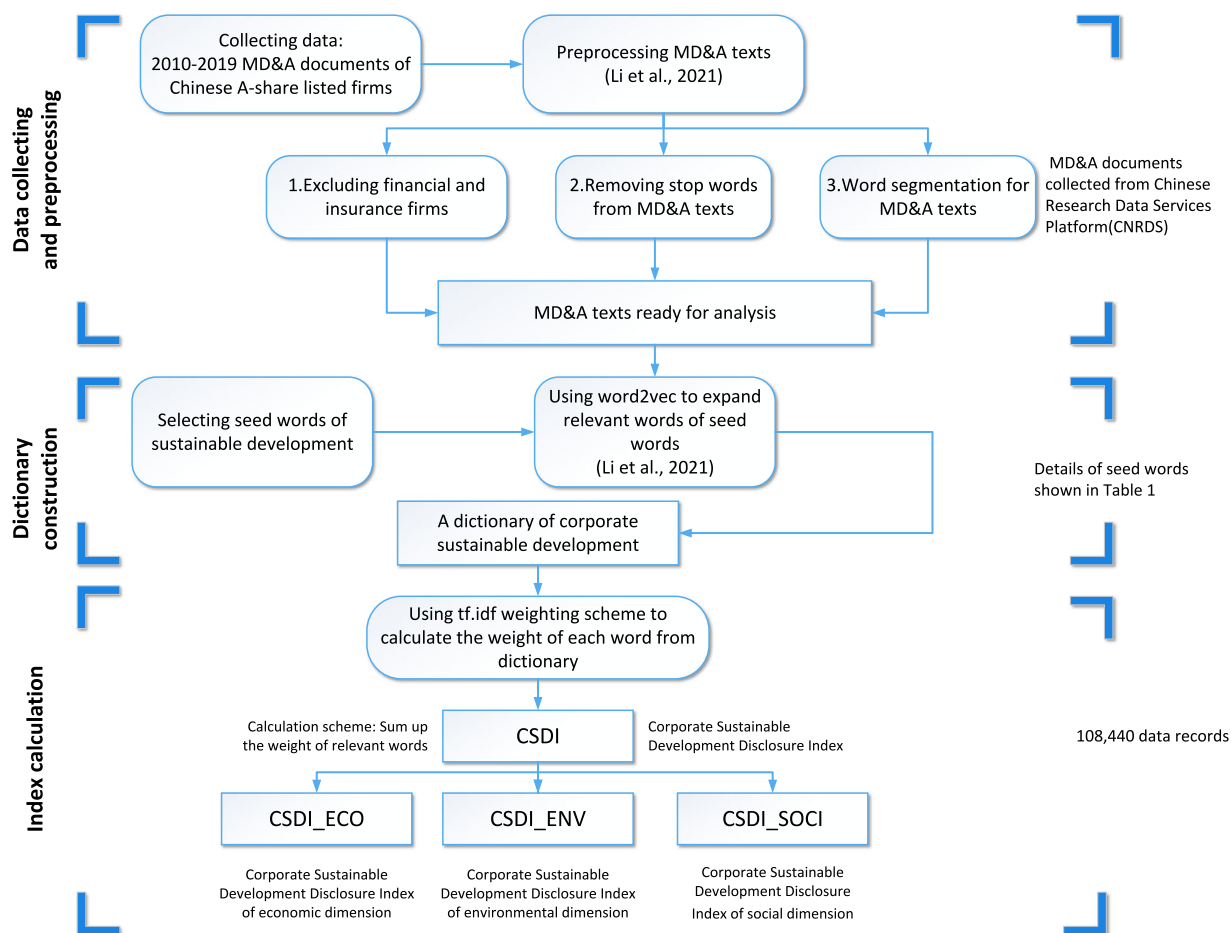
Fourth, the corporate sustainable development information disclosure dataset constructed in this study contributes to the knowledge management literature and expands the application of the serendipity-mindsponge-3D (SM3D) creativity management theory<sup>30</sup>. According to the SM3D framework, innovations are produced through 3-stage information processes: 1) information absorbing and filtering, 2) creativity processing, and 3) innovation outcome. The construction of the corporate sustainability information disclosure dataset follows the SM3D framework and applies the knowledge management theory to sustainable development areas. In the first stage, information on corporate sustainable development were collected from the MD&A documents and irrelevant information were screened out. The information generated from the first stage were then processed using machine learning techniques to create a corporate sustainable development dictionary and weights for each sustainability-related word. The final dataset of corporate sustainable development information disclosure was then produced as an outcome in the last stage. More generally, the SM3D knowledge management framework can be applied to future corporate sustainability management studies.

Last but important, the dataset of corporate sustainable development information disclosure can motivate the private sector to make greater contributions to global sustainable development. As major emission contributors, businesses lack incentives to make green investments willingly as it is hard to quantify their efforts towards sustainable development. The dataset constructed in this study can not only demonstrate the attitudes and strategic direction of corporate management towards sustainable transition but also provide stakeholders such as investors and regulators with quantitative benchmarks. To satisfy stakeholder expectations and so to maintain a positive public image and secure sustained capital flows, corporate executives have a stronger willingness to make voluntary low-carbon transition. Accordingly, the construction of this dataset can help promote green transition of the private sector and provide solutions to mitigating global sustainability problems such as climate change<sup>31</sup> and help developing countries to better achieve sustainable development<sup>32</sup>.

## Methods

The corporate sustainability disclosure dataset is constructed based on companies' MD&A documents. We collect the data on 29,134 MD&A texts over the period 2010–2019 from the China Research Data Service (CNRDS) platform<sup>33</sup>. The CNRDS database covers sub-databases related to economics, finance and business research, such as economic development section, corporate characteristics section, news and media section, and many more.

## Construction of the dataset on corporate sustainability disclosure



**Fig. 1** Workflow of the construction of the corporate sustainability disclosure dataset.

The MD&A text data used in this study is collected from the listed company text information module of the corporate characteristics section in CNRDS. It is worth mentioning that account registration is needed to log into the CNRDS platform to retrieve and use the data. We exclude financial and insurance companies because the financial sector has adopted different accounting and disclosure rules<sup>34</sup>. Finally, the data on the remaining 27,110 MD&A documents are obtained. Our workflow diagram for constructing the dataset is shown in Fig. 1.

The flowchart presents a brief overview of the process of building the corporate sustainability disclosure dataset. Next, we will explain the process of constructing the corporate sustainability disclosure dataset in detail.

**Pre-processing.** To facilitate *word2vec* to ‘read’ the neural network of MD&A texts, ‘learn’ the meaning of the corporate sustainability seed word set, and predict its similarity, we need to first clean the MD&A data and separate the words<sup>35</sup>.

Given that removing stop words improves the accuracy of text mining<sup>36</sup>, we select the following stop word lists that are currently extensively used: the Chinese stop word list, Baidu stop word list, Harbin Institute of Technology stop word list, and Sichuan University stop word list. The removal of stop words in MD&A texts is implemented using Python. Referring to the method of Li *et al.*<sup>25</sup>, we apply the *jieba* word segmentation library in Python to perform word segmentation on the text documents because *jieba* word segmentation technique is widely used in Chinese word segmentation<sup>37,38</sup>. Eventually, the clean MD&A texts of Chinese listed companies from 2010–2019 are obtained following word segmentation.

**The corporate sustainability dictionary.** The development of a corporate sustainability dictionary paves the basis for the construction of a corporate sustainability disclosure dataset. This development is achieved through the following two steps: first, we select the corporate sustainability closely-related seed words, and second, we expand the seed words to include their synonyms for building a corporate sustainability dictionary. Next, we will explain the process of building a corporate sustainability dictionary in detail.

*Step 1: Selection of corporate sustainability seed words.* We select corporate sustainability seed words based on the three dimensions of economic, environmental, and social, in accordance with the ‘triple bottom line’ defined

Corporate sustainability dimension	Subdimension	Seed words	Sources	
Economic dimension	Innovation management	Innovation	Guiso <i>et al.</i> <sup>52</sup>	
		Efficiency	Nini <i>et al.</i> <sup>79</sup>	
	Risk management	Repayment		
		Risks	Power (2009) <sup>45</sup>	
	Profitability	Profit	Grunig (1979) <sup>80</sup>	
	Corporate governance	Corporate governance	Performance	Nini <i>et al.</i> <sup>79</sup>
			Growth	
Development				
Expenses				
Management				
Environmental dimension	Ecology and environmental protection	Ecology	Sharma & Henriques (2005) <sup>49</sup>	
		Climate		
		Green	Prasad & Elmes (2005) <sup>81</sup>	
		Environmental protection		
	Pollution control	Emission reduction	Sharma & Henriques (2005) <sup>49</sup> ; Chan(2005) <sup>82</sup>	
		Pollution		
	Recycling	Low carbon	Sharma & Henriques (2005) <sup>49</sup> ; Chan(2005) <sup>82</sup>	
		Waste		
		Energy saving	Sharma & Henriques (2005) <sup>49</sup> ; Chan(2005) <sup>82</sup>	
		Renewable	Bansal (2005) <sup>1</sup> ; Chan (2005) <sup>82</sup>	
Social dimension	Community relations	Community	Guiso <i>et al.</i> <sup>52</sup>	
		Caring		
	Philanthropy	Ethics	Cochran & Wood (1984) <sup>51</sup>	
		Donation		
	Product quality	Safety		
		Responsibility		
	Information disclosure	Transparency	Bansal (2005) <sup>1</sup> ; Chan (2005) <sup>82</sup>	
		Fairness		
	Employment relations	Welfare	Sonenshein (2016) <sup>83</sup>	
		Team	Guiso <i>et al.</i> <sup>52</sup>	

**Table 1.** Selection of corporate sustainability seed words for each dimension.

in the GRI guidelines. To ensure that the set of corporate sustainability seed words is relatively authoritative and convincing, we take the following procedure. First, we extensively check-through the corporate sustainability literature and apply a triangulation process with several experts in the field of corporate sustainability<sup>39</sup>. As a result, the economic dimension of corporate sustainability<sup>40</sup> includes corporate innovation management<sup>41–44</sup>, risk management<sup>45</sup>, profitability<sup>46,47</sup>, and corporate governance<sup>48</sup>; the environmental dimension covers ecological protection, pollution control, and recycling<sup>49,50</sup>; and the social dimension includes community relations, philanthropy, product quality, information disclosure, and employment relations<sup>51,52</sup>. Second, we finalise corporate sustainability seed words based on the following principles: (1) the selected seed word must appear in the MD&A texts; and (2) after training, the synonyms must complement the meaning of the seed words, following the extension of the *word2vec* model. Finally, a total of 30 corporate sustainability seed words are obtained. Among these, 10 seed words each are dedicated to the economic, environmental, and social dimension, respectively. Table 1 presents the literature sources for the selection of the corporate sustainability seed word set.

**Step 2: Generation of corporate sustainability dictionary.** As a latest machine learning technique in natural language processing (NLP) and an open-source tool to produce word vectors<sup>53</sup>, *word2vec* uses neural networks to more accurately learn low-dimensional vectors that represent word meanings, converts words into vector representations, and predicts similarity between words based on the *cosine* similarity method<sup>54</sup>. In this study, we extend seed words and find their synonyms through using *word2vec* to obtain the semantic similarity between words in the MD&A texts. Accordingly, we use the trained *word2vec* model to extend the seed words to include their synonyms based on the full MD&A documents to construct the corporate sustainability dictionary.

To facilitate understanding, we take the extension of the synonyms in the economic dimension as an example. First, we assume that the MD&A corpus contains  $V$  words, and the initial word vector dimension of each seed word is  $V$ . The *word2vec* model reduces the dimension of each word vector to ensure that the word vector dimension is set in such a way that it can summarise the meaning of the seed words in a more comprehensive way without being excessively redundant. Following Li *et al.*'s<sup>25</sup> method of reducing the dimension of the word vector of each corporate cultural seed word, we set the word vector dimension to 300. Then, we have the word vector of 'innovation' is

$V^{(1)} = [x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots, x_{300}^{(1)}]$ , word vector of ‘efficiency’ is  $V^{(2)} = [x_1^{(2)}, x_2^{(2)}, x_3^{(2)}, \dots, x_{300}^{(2)}]$ , and the word vector of ‘management’ is  $V^{(10)} = [x_1^{(10)}, x_2^{(10)}, x_3^{(10)}, \dots, x_{300}^{(10)}]$ . Here,  $x_i^{(j)}$  denotes the dimension  $i$  of the  $j$ th seed word. Next, we calculate the similarity between the word vectors of each word in the MD&A corpus and the ten seed words. We follow the approach of Li *et al.*<sup>25</sup> in calculating the similarity between words using the *cosine* similarity between word vectors. The *cosine* similarity is expressed as the *cosine* of the angle between two word vectors  $A = (a_1, a_2, \dots, a_n)$  and  $B = (b_1, b_2, \dots, b_n)$ , following the formula shown below in Eq. (1). The higher the *cosine* similarity, the closer the *cosine* value of the angle between the two word vectors is to 1. This indicates that the more the angle between the two word vectors converges to 0, the more similar the two words are.

$$\text{sim}(A, B) = \frac{A \cdot B}{\|A\| \times \|B\|} \quad (1)$$

To ensure a sufficient number of synonyms, we first select 30 synonym words with the highest *cosine* similarity to each seed word. After obtaining the 30 synonyms, we manually check-through all of them to exclude words with inappropriate or irrelevant meaning to sustainability. As a last screening step, for each seed word, we retain top 10 synonyms with the highest *cosine* similarity to the vector of the seed word. Following the same approach, we obtain 100 words each for the environmental and social dimensions of corporate sustainability, respectively. Ultimately, we construct a corporate sustainability dictionary comprising a total of 330 words (30 seed words + 30\*10 synonyms of the seed words).

**Construction of the corporate sustainability disclosure dataset.** In this sub-section, we explain the process of constructing the Corporate Sustainability Disclosure Dataset (CSDI, CSDI\_ECO, CSDI\_ENV, CSDI\_SOC) in two steps.

*Step 1: Weighting scheme.* Referring to Li *et al.*'s<sup>25</sup> method of quantifying corporate culture, we utilise the *tf.idf* weighting scheme to calculate the weight of each word. This weighting scheme takes into account the importance of corporate sustainability-related words in a particular company's MD&A document and the entire MD&A corpus. Specifically, in this study, *tf* (Term Frequency) denotes the frequency of a word in a company's MD&A for a given year. Considering the word ‘innovation’ as an example, its frequency *tf* in text *j* is:

$$tf_{innovation,j} = \frac{\text{the frequency of the word "innovation" in MD \& A texts}}{\text{Total word counts of this MD \& A text}} \quad (2)$$

*idf* (inverse document frequency) denotes the inverse document frequency of a word in the texts of the MD&A corpus. In simple terms, *idf* represents the general importance of a word in the MD&A corpus. Considering the example of ‘innovation’ again, the inverse document frequency of ‘innovation’ in the MD&A corpus is:

$$idf_{innovation} = \log \left( \frac{\text{Total number of MD \& A texts in the corpus}}{\text{Number of texts containing "innovation" + 1}} \right) \quad (3)$$

Finally, we obtain the weight of ‘innovation’ in text *j* as follows:

$$tfidf_{innovation,j} = tf_{innovation,j} \times idf_{innovation} \quad (4)$$

We apply the *tf.idf* weighting scheme to calculate the weight of each word in the corporate sustainability dictionary.

*Step 2: Calculation of the corporate sustainability disclosure index (CSDI).* We then use the weighted sum of all words' frequency in the corresponding sustainability dictionary to derive the CSDI at the firm-year level. We also apply the same methodology to each of the three dimensions to derive the *CSDI\_Economic Dimension* (CSDI\_ECO), *CSDI\_Environmental Dimension* (CSDI\_ENV), and *CSDI\_Social Dimension* (CSDI\_SOC). Accordingly, the corporate sustainability disclosure dataset (CSDI, CSDI\_ECO, CSDI\_ENV, CSDI\_SOC) is generated.

**Statistical analysis of the corporate sustainability disclosure dataset.** In this section, we report summary statistics, summarise the trend of each dimension index, and analyse the correlation between each dimension index to provide a preliminary understanding of the corporate sustainability disclosure dataset. The results are presented in Table 2.

First, Table 2 Panel A presents the descriptive statistics for the corporate sustainability disclosure dataset. The mean value of economic dimension (CSDI\_ECO) is the largest, suggesting that economic sustainability is the basis for corporate sustainable development<sup>55</sup>. Moreover, the mean value of social dimension (CSDI\_SOC) is the smallest, which provides evidence to support the statement of Wartick and Cochran<sup>47</sup> that corporate social responsibility is less important than the maximisation of corporate profits. Furthermore, the variance of environmental dimension (CSDI\_ENV) is the largest, indicating that a significant difference exists in the degree of importance placed on environmental protection by different corporate executives. This finding also implies that promoting corporate environmental management is indispensable in the process of advancing sustainable economic and social development.

Second, we analyse the patterns in the corporate sustainability disclosure dataset. We investigate the consistency of CSDI, CSDI\_ECO, CSDI\_ENV, and CSDI\_SOC over the five-year window and display the correlation

Panel A – Summary statistics				
Variable	Mean	Std. Dev.	Min	Max
CSDI	0.233	0.122	0.000	1.447
CSDI_ECO	0.130	0.052	0.000	1.124
CSDI_ENV	0.063	0.095	0.000	1.365
CSDI_SOCI	0.041	0.043	0.000	0.528
Panel B – Autocorrelations of CSDI (CSDI_ECO, CSDI_ENV, CSDI_SOCI)				
Variables in year t	Year t-1	Year t-2	Year t-3	Year t-4
CSDI	0.107***	0.110***	0.080***	0.075***
CSDI_ECO	0.085***	0.100***	0.084***	0.072***
CSDI_ENV	0.058***	0.055***	0.037***	0.019***
CSDI_SOCI	0.181***	0.172***	0.165***	0.175***
Panel C – Correlations of CSDI, CSDI_ECO, CSDI_ENV, CSDI_SOCI				
Variables	CSDI	CSDI_ECO	CSDI_ENV	CSDI_SOCI
CSDI	1.000			
CSDI_ECO	0.479***	1.000		
CSDI_ENV	0.810***	0.017**	1.000	
CSDI_SOCI	0.490***	0.119***	0.114***	1.000

**Table 2.** Descriptive statistics. Notes: Panel A lists the descriptive statistics of CSDI, CSDI\_ECO, CSDI\_ENV, CSDI\_SOCI; Panel B presents the changing trends of CSDI, CSDI\_ECO, CSDI\_ENV, CSDI\_SOCI; Panel C shows the correlations between CSDI and CSDI\_ECO, CSDI\_ENV, CSDI\_SOCI; and \*\*\* $p < 0.01$ , \*\* $p < 0.05$ , \* $p < 0.1$ .

coefficients of each dimension of the corporate sustainability disclosure index for the current period with its lagged periods in Panel B of Table 2. We find that the correlation coefficients between current year and years  $t-1$  to  $t-4$  are all significantly positive, suggesting that a consistent attitude of corporate management towards sustainability.

Third, we examine the correlation between CSDI and each dimension index. Table 2 Panel C presents the correlations between CSDI and CSDI\_ECO, CSDI\_ENV, and CSDI\_SOCI. The results indicate that CSDI is significantly and positively correlated with CSDI\_ECO, CSDI\_ENV, and CSDI\_SOCI, and there also exists a significant positive correlation between CSDI\_ECO, CSDI\_ENV, and CSDI\_SOCI. This corroborates the fact that in the long run, enterprises undertaking environmental and social responsibilities will help contribute to the maximisation of the company value; that is, corporate environmental and social performance will translate into financial performance<sup>56</sup>.

### Data Records

We upload a total of 108,440 data points from the Corporate Sustainability Disclosure dataset. Besides, the corporate sustainability dictionary is also uploaded. The data records consist of the following datasets:

- The Corporate Sustainability Disclosure Dataset (CSDI, CSDI\_ECO, CSDI\_ENV, CSDI\_SOCI) contains a total of 108,440 (27,110\*4) data points<sup>57</sup>. Notably, this dataset has an individual file for each year (2010–2019) with column headings of “Year” describing the year of the dataset, “Corporate Code” and “Corporate Name” representing the name and stock code of the company, “CSDI” displaying the overall disclosure index of corporate sustainable development, “CSDI\_ECO”, “CSDI\_ENV”, and “CSDI\_SOCI” listing the respective dimension disclosure index, and “Industry Code” indicating the industry classification code of the company.
- The Corporate Sustainability Dictionary includes a total of 330 words (30 seed words + 30\*10 synonyms of the seed words)<sup>58</sup>.
- The data used in the Technical Validation section includes corporate indicators of economic performance, environmental performance, and social performance, as well as a battery of control variables<sup>59</sup>.

### Technical Validation

In this section, we examine the accuracy of the text mining techniques used to construct the dataset and validate the relationship between the corporate sustainability disclosure dataset and corresponding corporate actual performance.

**Validation of the text mining techniques.** To verify the accuracy of the text mining technique employed in this study, we randomly select ten MD&A documents from different industries across different years, and manually retrieve ten words included in the constructed corporate sustainability dictionary. The differences between the manual check and text mining results are compared and reported in Table 3. Panel A of Table 3 presents the number of occurrences of sustainability-related words in the corresponding MD&A documents using the manual retrieval method, while Panel B presents the corresponding results based on the text mining techniques. It is clear that the results of word counts based on manual retrieval and the text mining techniques are identical.

Panel A – Manual Search Word Count										
	Chemical fibre manufacturing	Air transport industry	Capital markets services	General equipment manufacturing	Catering	Road transport industry	Wholesale trade	Rubber and plastic products industry	Real estate	Computer, communications, & other electronic equipment manufacturing
Energy saving	0	0	4	16	1	0	0	0	0	0
Profit	6	6	5	0	0	3	1	3	4	9
Responsibility	0	7	3	2	3	0	2	0	0	0
Efficiency	0	1	1	1	0	0	0	0	3	4
Environmental protection	1	1	7	10	1	7	7	0	0	0
Safety	0	37	3	3	5	10	5	2	1	0
Development	9	28	36	23	32	27	8	36	27	14
Ecology	0	0	0	0	0	0	0	0	0	0
Fairness	0	0	0	0	0	0	0	0	1	0
Panel B – Text Mining Word Count										
Energy saving	0	0	4	16	1	0	0	0	0	0
Profit	6	6	5	0	0	3	1	3	4	9
Responsibility	0	7	3	2	3	0	2	0	0	0
Efficiency	0	1	1	1	0	0	0	0	3	4
Environmental protection	1	1	7	10	1	7	7	0	0	0
Safety	0	37	3	3	5	10	5	2	1	0
Development	9	28	36	23	32	27	8	36	27	14
Ecology	0	0	0	0	0	0	0	0	0	0
Fairness	0	0	0	0	0	0	0	0	1	0

**Table 3.** Manual-checking for sample text mining results. Notes: Panel A reports the word counts under manual retrieval; and Panel B lists the word counts based on text mining techniques.

**Validation of the corporate sustainability disclosure index.** To further ascertain the validity of the corporate sustainability disclosure dataset, we follow the approach of Li *et al.* to examine the effectiveness of the disclosure indexes<sup>25</sup>. For the purposes of this study, we collect corporate real performance indicators for each of the three sustainability dimensions. To control for the potential effects of other factors on the regression results, we add a battery of firm-level control variables<sup>60–63</sup> when examining the relationship between each dimension disclosure index and the corresponding corporate real sustainability performance in that dimension. The following firm-level control variables are added: financial leverage (*leverage*), return on assets (*ROA*), shareholding of the largest shareholder (*Share1*), whether audited by a Big Four auditor company (*Big4*), firm size (*Size*), and firm age (*Listage*). The above data are all obtained from the CSMAR database<sup>64</sup> and are available when logged in. The validation results are described in detail below.

First, we test the validity of *Corporate Sustainability Disclosure Index\_Economic Dimension* (CSDI\_ECO). To verify the validity of CSDI\_ECO, we examine the relationship between CSDI\_ECO and the following corporate economic performance indicators: total factor productivity (*TFP*), growth in revenue (*Growth*), bankruptcy risk (*Oscore*), and financial constraints (*SA*). First, *TFP* reflects the efficiency of converting inputs into outputs in the production process. Thus, at the firm level, a firm's *TFP* can represent the firm's ability to become economically sustainable. We calculate *TFP* based on the method developed by Levinsohn and Petrin<sup>65</sup>. Second, we use revenue growth to reflect the economic growth ability of the firm<sup>66</sup>. Third, because financial constraints limit firms' investment in research and development (R&D) activities<sup>67</sup>, which will affect firms' long-term economic development, we select the *SA* index<sup>68</sup> to represent firms' actual financial constraints. The larger the *SA* index, the more severe the firms' financial constraints are, and so the less sustainable their economic development. Lastly, we select *Oscore*<sup>69</sup> to reflect corporate economic sustainability from the perspective of risk management, with a larger *Oscore* representing greater distress in business operations and so higher likelihood of being bankruptcy. The validation results are reported in Table 4. Overall, the associations of CSDI\_ECO with all of the four economic indicators are highly correlated. Thus, the validation results shown in Table 4 provide solid evidence that the CSDI\_ECO constructed in this study can reflect the actual performance of firms' economic sustainability.

Second, we examine the validity of *Corporate Sustainability Disclosure Index\_Environmental Dimension* (CSDI\_ENV). To verify the validity of the CSDI\_ENV, we test the relationship between CSDI\_ENV and corporate environmental performance indicators. The following indicators are selected: disclosure of dust control (*Dust\_control*), disclosure of wastewater discharge (*Wastewater*), disclosure of solid wastes utilisation and disposal (*Solid\_waste*), and disclosure of waste gas abatement and control (*Waste\_gas*)<sup>70</sup>. The validation results are presented in Table 5. The associations between CSDI\_ENV and all of the four corporate environmental indicators are highly correlated. The results document that CSDI\_ENV constructed in this study is positively associated with companies' environmental performance.

Third, we provide the validation of *Corporate Sustainability Disclosure Index\_Social Dimension* (CSDI\_SOCI). In verifying the validity of CSDI\_SOCI, we examine the relationship between CSDI\_SOCI and corporate

VARIABLES	TFP	Growth	SA	O'score
CSDI_ECO	0.384***	0.145***	-0.203***	-2.832***
	(3.87)	(2.83)	(-5.31)	(-4.83)
Leverage	0.930***	-0.033**	-0.174***	5.947***
	(28.14)	(-1.98)	(-14.54)	(32.31)
ROA	3.281***	0.538***	-0.020	-4.017***
	(33.11)	(10.38)	(-0.56)	(-6.54)
Share1	0.195***	0.007	0.193***	2.428***
	(6.51)	(0.38)	(14.73)	(12.25)
Big4	0.093***	0.025**	0.096***	0.563***
	(4.81)	(2.47)	(9.33)	(5.37)
Size	0.641***	0.004	0.008***	-1.118***
	(123.30)	(1.40)	(3.29)	(-39.68)
Listage	0.009	0.007*	-0.008***	-0.023
	(1.26)	(1.77)	(-2.58)	(-0.51)
Constant	-6.025***	-0.436***	-3.662***	14.408***
	(-50.98)	(-6.76)	(-72.36)	(22.45)
Observations	13,758	15,599	15,599	15,599
R-squared	0.800	0.183	0.253	0.187
Industry Fixed Effects	YES	YES	YES	YES
Year Fixed Effects	YES	YES	YES	YES

**Table 4.** Validation of CSDI\_ECO. Notes: \*\*\* $p < 0.01$ , \*\* $p < 0.05$ , \* $p < 0.1$ ; robust t-statistics in parentheses.

VARIABLES	Dust_control	Wastewater	Solid_waste	Waste_gas
CSDI_ENV	2.886***	3.144***	3.081***	3.781***
	(10.96)	(12.57)	(12.49)	(14.51)
Leverage	-0.159	-0.347**	-0.539***	-0.242*
	(-0.94)	(-2.52)	(-3.65)	(-1.70)
ROA	-1.197**	-0.013	0.783*	-0.448
	(-2.42)	(-0.03)	(1.85)	(-1.10)
Share1	0.417**	0.103	-0.142	0.214
	(2.47)	(0.73)	(-0.96)	(1.47)
Big4	0.281***	0.324***	0.393***	0.583***
	(2.96)	(3.69)	(4.56)	(6.67)
Size	0.506***	0.483***	0.520***	0.528***
	(19.63)	(21.42)	(22.18)	(22.78)
Listage	-0.121***	-0.049	-0.116***	-0.079**
	(-3.22)	(-1.54)	(-3.56)	(-2.47)
Constant	-13.908***	-13.154***	-14.659***	-14.438***
	(-21.06)	(-23.66)	(-21.62)	(-25.20)
Observations	15,139	15,529	15,501	15,516
Industry Fixed Effects	YES	YES	YES	YES
Year Fixed Effects	YES	YES	YES	YES

**Table 5.** Validation of CSDI\_ENV. Notes: \*\*\*:  $p < 0.01$ , \*\*:  $p < 0.05$ , \*:  $p < 0.1$ ; robust t-statistics in parentheses.

social performance indicators. To select social performance variables, we refer to the requirements of China Securities Regulatory Commission for listed firms to fulfil their corporate social responsibility. In this study, we follow Chen *et al.*'s definition of corporate social responsibility<sup>70</sup> and select the following four indicators: disclosure of public relations and public welfare (*Public*), disclosure of the protection of suppliers' rights and interests (*Suppliers*), disclosure of the protection of creditors' rights and interests (*Creditor*), disclosure of the protection of customers and consumers' rights and interests (*Custm\_consm*). The validation results presented in Table 6. The testing results indicate that the CSDI\_SOCI constructed in this study reflects the actual status of corporate CSR disclosure.

Last but more important, we test the validity of *Corporate Sustainability Disclosure Index* (CSDI). To measure the validity of CSDI, we construct aggregate corporate sustainability performance indicators and examine the relationship between CSDI and these indicators. First, we standardise the above-used twelve corporate sustainability performance indicators. Second, to ensure the robustness of the results, we adopt different weighting



VARIABLES	Public	Creditor	Supplier	Custm_consm
CSDI_SOCI	9.050*** (16.70)	6.686*** (14.28)	4.731*** (9.84)	4.679*** (9.65)
Leverage	-0.340** (-2.55)	-0.019 (-0.16)	-0.538*** (-4.38)	-0.241** (-2.00)
ROA	1.452*** (3.64)	0.532 (1.52)	0.096 (0.27)	0.522 (1.47)
Share1	-0.221 (-1.58)	-0.256** (-2.05)	-0.143 (-1.12)	0.039 (0.30)
Big4	0.580*** (4.99)	-0.329*** (-4.01)	0.246*** (3.05)	0.435*** (4.70)
Size	0.668*** (28.01)	0.156*** (8.06)	0.492*** (24.09)	0.522*** (25.06)
Listage	0.128*** (3.95)	0.189*** (6.79)	0.141*** (5.01)	0.135*** (4.70)
Constant	-16.223*** (-28.78)	-5.680*** (-12.42)	-12.743*** (-25.46)	-12.826*** (-25.53)
Observations	15,555	15,575	15,548	15,548
Industry Fixed Effects	YES	YES	YES	YES
Year Fixed Effects	YES	YES	YES	YES

**Table 6.** Validation of CSDI\_SOCI. Notes: \*\*\* $p < 0.01$ , \*\* $p < 0.05$ , \* $p < 0.1$ ; robust t-statistics in parentheses.

VARIABLES	Panel A – Score1				Panel B – Score2				Panel C – Score3			
	Score1	Score1	Score1	Score1	Score2	Score2	Score2	Score2	Score3	Score3	Score3	Score3
CSDI	0.468*** (11.50)				1.066*** (11.64)				0.685*** (12.22)			
Leverage	-0.024 (-0.86)	-0.015 (-0.53)	-0.021 (-0.78)	-0.008 (-0.30)	-0.105* (-1.69)	-0.085 (-1.36)	-0.105* (-1.69)	-0.074 (-1.19)	-0.122*** (-3.24)	-0.110*** (-2.88)	-0.121*** (-3.21)	-0.102*** (-2.68)
ROA	-0.021 (-0.26)	-0.035 (-0.43)	-0.037 (-0.46)	-0.061 (-0.75)	-0.255 (-1.36)	-0.286 (-1.51)	-0.287 (-1.53)	-0.331* (-1.76)	-0.168 (-1.47)	-0.186 (-1.62)	-0.190* (-1.66)	-0.220* (-1.92)
Share1	0.120*** (4.13)	0.117*** (4.01)	0.122*** (4.20)	0.115*** (3.97)	0.218*** (3.33)	0.211*** (3.21)	0.226*** (3.45)	0.210*** (3.20)	0.047 (1.20)	0.043 (1.07)	0.052 (1.30)	0.042 (1.05)
Big4	0.079*** (3.93)	0.073*** (3.65)	0.078*** (3.89)	0.077*** (3.84)	0.183*** (3.87)	0.170*** (3.60)	0.183*** (3.88)	0.175*** (3.70)	0.066** (2.38)	0.058** (2.08)	0.065** (2.37)	0.062** (2.22)
Size	0.125*** (25.98)	0.126*** (26.06)	0.124*** (25.56)	0.130*** (26.64)	0.289*** (25.91)	0.291*** (25.97)	0.284*** (25.42)	0.295*** (26.24)	0.149*** (22.44)	0.151*** (22.52)	0.146*** (21.95)	0.154*** (22.93)
Listage	0.015** (2.27)	0.013* (1.94)	0.014** (2.00)	0.015** (2.14)	-0.008 (-0.53)	-0.013 (-0.83)	-0.012 (-0.76)	-0.012 (-0.77)	0.005 (0.58)	0.002 (0.25)	0.003 (0.32)	0.003 (0.35)
CSDI_ECO		0.247*** (2.72)				0.609*** (2.98)				0.428*** (3.42)		
CSDI_ENV			0.500*** (8.85)				1.408*** (10.90)				0.840*** (10.67)	
CSDI_SOCI				1.004*** (9.62)				1.201*** (5.12)				0.974*** (6.85)
Constant	-3.148*** (-27.89)	-3.087*** (-27.01)	-3.026*** (-26.83)	-3.153*** (-27.87)	-7.063*** (-27.57)	-6.932*** (-26.68)	-6.773*** (-26.48)	-6.965*** (-27.04)	-2.945*** (-19.28)	-2.866*** (-18.48)	-2.762*** (-18.09)	-2.903*** (-18.92)
Observations	12,672	12,672	12,672	12,672	12,672	12,672	12,672	12,672	12,672	12,672	12,672	12,672
R-squared	0.206	0.197	0.202	0.202	0.209	0.201	0.209	0.202	0.186	0.177	0.185	0.179
Industry Fixed Effects	YES	YES	YES	YES	YES	YES	YES	YES	YES	YES	YES	YES
Year Fixed Effects	YES	YES	YES	YES	YES	YES	YES	YES	YES	YES	YES	YES

**Table 7.** Validation of CSDI. Notes: \*\*\* $p < 0.01$ , \*\* $p < 0.05$ , \* $p < 0.1$ ; robust t-statistics in parentheses; Score1: Composite sustainability score calculated based on equal-weighting scheme; Score2: Composite sustainability score calculated based on the principal component analysis (PCA); Score3: Composite sustainability score calculated based on the entropy method.

methods to assign weights to the twelve indicators and then calculate the composite sustainability indicator. The first composite indicator is obtained through an equally-weighted average of the standardised twelve sustainability indicators, denoted as *Score1*. The second composite indicator is calculated based on the application of the principal component analysis (PCA) method to the twelve sustainability indicators<sup>71,72</sup>, denoted as *Score2*. The third composite score is obtained based on the entropy method<sup>73</sup>, denoted as *Score3*. These validation results are presented in Table 7, respectively. The first columns in all three panels report the association of the composite score with CSDI, while the remaining three columns report that with the three dimension indexes (i.e., CSDI\_ECO, CSDI\_ENV, and CSDI\_SOC1).

Overall, the associations between CSDI and all composite corporate sustainability scores are significantly correlated. Similar results are reached for the three dimension scores. Accordingly, the validation results document that the CSDI constructed in this study can reflect the actual performance of firms' sustainable development.

Taken together, these validation results provide solid evidence that the constructed *Corporate Sustainability Disclosure Index* (CSDI), *CSDI\_Economic Dimension* (CSDI\_ECO), *CSDI\_Environmental Dimension* (CSDI\_ENV), and *CSDI\_Social Dimension* (CSDI\_SOC1) are valid and can reflect the actual disclosure of corporate sustainability performance.

**Limitations.** Our dataset has the following limitations. First, limited by the availability of latest firm-level data, the dataset constructed in this study does not take into account the impact of the COVID-19 pandemic on corporate sustainability. Second, comparisons between different industries are beyond the scope of current study and can be further explored in future studies. Third, the sustainability information disclosure index constructed in this study does not fully consider the context of the words. Although negative tones are very rare in the MD&A documents, a lack of context consideration might still generate some (albeit very limited) influences<sup>25,29</sup>. Future studies can make attempts to improve machine learning techniques to take into account the context where key words are located when constructing text-based quantitative measures. Lastly, as market conditions in developed countries are different from those in developing countries such as China, future studies can build on the spirit of this study to investigate the progress of corporate sustainability disclosure in developed countries.

### Code availability

The codes used for calculation and analysis in this study are available in *figshare*<sup>74–78</sup>.

Received: 14 November 2022; Accepted: 20 March 2023;

Published online: 31 March 2023

### References

- Bansal, P. Evolving sustainably: A longitudinal study of corporate sustainable development. *Strategic Manage J.* **26**, 197–218 (2005).
- Leisinger, K. Business needs to embrace sustainability targets. *Nature.* **528**, 165–165 (2015).
- Hart, S. L. Beyond greening: Strategies for a sustainable world. *Harvard Bus Rev.* **75**, 66–77 (1997).
- Eccles, R. G., Ioannou, I. & Serafeim, G. The impact of corporate sustainability on organizational processes and performance. *Manage Sci.* **60**, 2835–2857 (2014).
- Glaser, G. Base sustainable development goals on science. *Nature.* **491**, 35–35 (2012).
- Maclaren, V. Urban sustainability reporting. *J Am Plann Assoc.* **62**, 184–202 (1996).
- Scipioni, A., Mazzi, A., Mason, M. & Manzardo, A. The dashboard of sustainability to measure the local urban sustainable development: The case study of Padua Municipality. *Ecol Indic.* **9**, 364–380 (2009).
- Chung, M. G., Frank, K. A., Pokhrel, Y., Dietz, T. & Liu, J. Natural infrastructure in sustaining global urban freshwater ecosystem services. *Nat Sustain.* **4**, 1068–1075 (2021).
- Leal, Filho, W. et al. A. L. Towards symbiotic approaches between universities, sustainable development, and cities. *Sci Rep.* **12**, 1–8 (2022).
- Searcy, C. Corporate sustainability performance measurement systems: A review and research agenda. *J Bus Ethics.* **107**, 239–253 (2012).
- Hahn, R. & Lülfs, R. Legitimizing negative aspects in GRI-oriented sustainability reporting: A qualitative analysis of corporate disclosure strategies. *J Bus Ethics.* **123**, 401–420 (2014).
- Adler, R., Mansi, M. & Pandey, R. “Biodiversity and threatened species reporting by the top fortune global companies”. *Account Audit Accoun.* **31**, 787–825 (2018).
- Milne, M. J. & Gray, R. W. (h) ither ecology? The triple bottom line, the global reporting initiative, and corporate sustainability reporting. *J Bus Ethics.* **118**, 13–29 (2013).
- Larrinaga, C. & Bebbington, J. “The pre-history of sustainability reporting: a constructivist reading”. *Account Audit Accoun.* **34**, 162–181 (2021).
- Rezaee, Z. & Tuo, L. Are the quantity and quality of sustainability disclosures associated with the innate and discretionary earnings quality? *J Bus Ethics.* **155**, 763–786 (2019).
- Adams, C. A. The ethical, social and environmental reporting performance portrayal gap. *Account Audit Accoun.* **17**, 731–757 (2004).
- Clarkson, P. M., Kao, J. L. & Richardson, G. D. Evidence that management discussion and analysis (MD&A) is a part of a firm's overall disclosure package. *Contemp Account Res.* **16**, 111–134 (1999).
- Mayew, W. J., Sethuraman, M. & Venkatachalam, M. MD&A disclosure and the firm's ability to continue as a going concern. *Account Rev.* **90**, 1621–1651 (2015).
- Luo, Y. Determinants of entry in an emerging economy: A multilevel approach. *J Manage Stud.* **38**, 443–472 (2001).
- Guan, D. et al. Structural decline in China's CO2 emissions through transitions in industry and energy systems. *Nat. Geosci.* **11**, 551–555 (2018).
- Qian, H. et al. China industrial environmental database 1998–2015. *Sci Data.* **9**, 1–13 (2022).
- Shan, Y., Huang, Q., Guan, D. & Hubacek, K. China CO2 emission accounts 2016–2017. *Sci Data.* **7**, 1–9 (2020).
- Shan, Y., Liu, J., Liu, Z., Shao, S. & Guan, D. An emissions-socioeconomic inventory of Chinese cities. *Sci Data.* **6**, 1–10 (2019).
- Wang, W. & Zhang, Y. J. Does China's carbon emissions trading scheme affect the market power of high-carbon enterprises? *Energy Econ.* **108**, 105906 (2022).
- Li, K., Mai, F., Shen, R. & Yan, X. Measuring corporate culture using machine learning. *Rev Financ Stud.* **34**, 3265–3315 (2021).
- Zhang, G. et al. China's environmental policy intensity for 1978–2019. *Sci Data.* **9**, 1–10 (2022).

27. Vuong, Q. H. The (ir) rational consideration of the cost of science in transition economies. *Nat Hum Behav.* **2**, 5–5 (2018).
28. Vuong, Q. H., Nguyen, H. T. T., Pham, T. H., Ho, M. T. & Nguyen, M. H. Assessing the ideological homogeneity in entrepreneurial finance research by highly cited publications. *Hum Soc Sci Commun.* **8**, 1–11 (2021).
29. Henry, E. & Leone, A. J. Measuring qualitative information in capital markets research: comparison of alternative methodologies to measure disclosure tone. *Account Rev.* **91**, 153–178 (2016).
30. Vuong, Q. H. *et al.* Covid-19 vaccines production and societal immunization under the serendipity-mindsponge-3D knowledge management theory and conceptual framework. *Humanit Soc Sci Commun.* **9**, 1–12 (2022).
31. Vuong, Q. H. The semiconducting principle of monetary and environmental values exchange. *Econ Bus Lett.* **10**, 284–290 (2021).
32. Vuong, Q. H. *et al.* An open database of productivity in Vietnam's social sciences and humanities for public use. *Sci Data.* **5**, 1–15 (2018).
33. China Research Data Service Platform <https://www.cnrds.com/Home/Index#/FeaturedDatabase/DB/CMDA> (2022).
34. Huang, Y. S. & Wang, C. J. Corporate governance and risk-taking of Chinese firms: the role of board size. *Int Rev Econ Financ.* **37**, 96–113 (2015).
35. Kothari, S. P., Li, X. & Short, J. E. The effect of disclosures by management, analysts, and business press on cost of capital, return volatility, and analyst forecasts: A study using content analysis. *Account Rev.* **84**, 1639–1670 (2009).
36. Wilbur, W. J. & Sirotkin, K. The automatic identification of stop words. *Inf Sci.* **18**, 45–55 (1992).
37. Fan, C. *et al.* China's Gridded Manufacturing Dataset. *Sci Data.* **9**, 1–14 (2022).
38. Lian, Y. *et al.* Cyber violence caused by the disclosure of route information during the COVID-19 pandemic. *Humanit Soc Sci Commun.* **9**, 417 (2022).
39. Loughran, T. & McDonald, B. When is a liability not a liability? Textual analysis, dictionaries, and 10-Ks. *J Financ.* **66**, 35–65 (2011).
40. Scherer, A. G. & Palazzo, G. The new political role of business in a globalized world: A review of a new perspective on CSR and its implications for the firm, governance, and democracy. *J Manage Stud.* **48**, 899–931 (2011).
41. Holmes, S. L. Corporate social performance: past and present areas of commitment. *Acad Manage J.* **20**, 433–438 (1977).
42. Abbott, W. F. & Monsen, R. J. On the measurement of corporate social responsibility: Self-reported disclosure as a method of measuring corporate social investment. *Acad Manage J.* **22**, 501–515 (1979).
43. McGuire, J. B., Sundgren, A. & Schneeweis, T. Corporate social responsibility and firm financial performance. *Acad Manage J.* **31**, 854–872 (1988).
44. Fryxell, G. & Wang, J. The fortune corporate reputation index: reputation for what? *J Manage.* **20**, 1–14 (1994).
45. Power, M. The risk management of nothing. *Account, Org Soc.* **34**, 849–855 (2009).
46. Ball, R., Gerakos, J., Linnainmaa, J. T. & Nikolaev, V. V. Deflating profitability. *J Financ Econ.* **117**, 225–248 (2015).
47. Wartick, S. L. & Cochran, P. L. The evolution of the corporate social performance model. *Acad Manage Rev.* **10**, 758–769 (1985).
48. Bebchuk, L., Cohen, A. & Ferrell, A. What matters in corporate governance? *Rev Financ Stud.* **22**, 783–827 (2009).
49. Sharma, S. & Henriques, I. Stakeholder influences on sustainability practices in the Canadian forest products industry. *Strategic Manage J.* **26**, 159–180 (2005).
50. Chan, K. Mass communication and pro-environmental behaviour: Waste recycling in Hong Kong. *J Environ Manage.* **52**, 317–325 (1998).
51. Cochran, P. L. & Wood, R. A. Corporate social responsibility and financial performance. *Acad Manage J.* **27**, 42–56 (1984).
52. Guiso, L., Sapienza, P. & Zingales, L. The value of corporate culture. *J Financ Econ.* **117**, 60–76 (2015).
53. Church, K. W. Word2Vec. *Nat Lang Eng.* **23**, 155–162 (2017).
54. Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., Dean, J. Distributed representations of words and phrases and their compositionality. *Adv Neural Inf Process Syst.* **26** (2013).
55. Hawn, O., Chatterji, A. K. & Mitchell, W. Do investors actually value sustainability? New evidence from investor reactions to the Dow Jones Sustainability Index (DJSI). *Strategic Manage J.* **39**, 949–976 (2018).
56. Flammer, C., Hong, B. & Minor, D. Corporate governance and the rise of integrating corporate social responsibility criteria in executive compensation: effectiveness and implications for firm outcomes. *Strategic Manage J.* **40**, 1097–1122 (2019).
57. Tian, J., Cheng, Q., Xue, R., Han, Y. & Shan, Y. Information disclosure index of sustainable development of listed firms in China 2010–2019. *figshare* <https://doi.org/10.6084/m9.figshare.21550155.v3> (2022).
58. Tian, J., Cheng, Q., Xue, R., Han, Y. & Shan, Y. The dictionary of corporate sustainable development. *figshare* <https://doi.org/10.6084/m9.figshare.21550230.v5> (2022).
59. Tian, J., Cheng, Q., Xue, R., Han, Y. & Shan, Y. Variables used in the validation section. *figshare*. <https://doi.org/10.6084/m9.figshare.21923322.v2> (2023).
60. Kim, Y., Park, M. S. & Wier, B. Is earnings quality associated with corporate social responsibility? *Account Rev.* **87**, 761–796 (2012).
61. Gerged, A. M., Beddewela, E. & Cowton, C. J. Is corporate environmental disclosure associated with firm value? A multicountry study of gulf cooperation council firms. *Bus Strateg Environ.* **30**, 185–203 (2021).
62. Anderson, R. C. & Reeb, D. M. Founding-family ownership and firm performance: Evidence from the S&P 500. *J Financ.* **58**, 1301–1328 (2003).
63. Roychowdhury, S. Earnings management through real activities manipulation. *J Account Econ.* **42**, 335–370 (2006).
64. China Stock Market & Accounting Research Database <https://www.gtarsc.com> (2022).
65. Levinsohn, J. & Petrin, A. Estimating production functions using inputs to control for unobservables. *Rev Econ Stud.* **70**, 317–341 (2003).
66. Semrau, T. & Sigmund, S. Networking ability and the financial performance of new ventures: A mediation analysis among younger and more mature firms. *Strateg Entrep J.* **6**, 335–354 (2012).
67. Chang, K., Zeng, Y., Wang, W. & Wu, X. The effects of credit policy and financial constraints on tangible and research & development investment: firm-level evidence from china's renewable energy industry. *Energy Policy.* **130**, 438–447 (2019).
68. Fee, C. E., Hadlock, C. J. & Pierce, J. R. Investment, financing constraints, and internal capital markets: Evidence from the advertising expenditures of multinational firms. *Rev Financ Stud.* **22**, 2361–2392 (2009).
69. Ohlson, J. A. Financial ratios and the probabilistic prediction of bankruptcy. *J Account Res.* **18**, 109–131 (1980).
70. Chen, Y. C., Hung, M. & Wang, Y. The effect of mandatory CSR disclosure on firm profitability and social externalities: Evidence from China. *J Account Econ.* **65**, 169–190 (2018).
71. Tan, F. & Lu, Z. Assessing regional sustainable development through an integration of nonlinear principal component analysis and Gram Schmidt orthogonalization. *Ecol Indic.* **63**, 71–81 (2016).
72. Lever, J., Krzywinski, M. & Altman, N. Points of significance: Principal component analysis. *Nat Methods.* **14**, 641–643 (2017).
73. Zhao, J., Ji, G., Tian, Y., Chen, Y. & Wang, Z. Environmental vulnerability assessment for mainland China based on entropy method. *Ecol Indic.* **91**, 410–422 (2018).
74. Tian, J., Cheng, Q., Xue, R., Han, Y. & Shan, Y. MD&A text preprocessing code. *figshare*. <https://doi.org/10.6084/m9.figshare.21923382.v1> (2023).
75. Tian, J., Cheng, Q., Xue, R., Han, Y. & Shan, Y. Code for calculating similar words. *figshare*. <https://doi.org/10.6084/m9.figshare.21923502.v1> (2023).
76. Tian, J., Cheng, Q., Xue, R., Han, Y. & Shan, Y. Code for calculating tf.idf. *figshare*. <https://doi.org/10.6084/m9.figshare.21923514.v1> (2023).

77. Tian, J., Cheng, Q., Xue, R., Han, Y. & Shan, Y. Code for calculating word frequency. *figshare*. <https://doi.org/10.6084/m9.figshare.21923532.v2> (2023).
78. Tian, J., Cheng, Q., Xue, R., Han, Y. & Shan, Y. Index validation code. *figshare*. <https://doi.org/10.6084/m9.figshare.21929394.v1> (2023).
79. Nini, G., Smith, D. C. & Sufi, A. Creditor control rights, corporate governance, and firm value. *Rev Financ Stud.* **25**, 1713–1761 (2012).
80. Grunig, J. E. A new measure of public opinion on corporate social responsibility. *Acad Manage J.* **22**, 738–764 (1979).
81. Prasad, P. & Elmes, M. In the name of the practical: Unearthing the hegemony of pragmatics in the discourse of environmental management. *J Manage Stud.* **42**, 845–866 (2005).
82. Chan, R. Y. K. Does the natural-resource-based view of the firm apply in an emerging economy? A survey of foreign invested enterprises in China. *J Manage Stud.* **42**, 625–672 (2005).
83. Sonenshein, S. How corporations overcome issue illegitimacy and issue equivocality to address social welfare: The role of the social change agent. *Acad Manage Rev.* **41**, 349–366 (2016).

## Acknowledgements

This research was supported by the National Social Science Foundation of China (Grant Number: 21BTJ019).

## Author contributions

**Jinfang Tian:** Conceptualization, Funding acquisition, Project administration, Writing – original draft. **Qian Cheng:** Data curation, Investigation, Methodology, Writing – original draft. **Rui Xue:** Conceptualization, Formal analysis, Validation, Writing – review & editing. **Yilong Han:** Data curation, Writing – review & editing. **Yuli Shan:** Conceptualization, Validation, Writing – review & editing.

## Competing interests

The authors declare no competing interests.

## Additional information

**Correspondence** and requests for materials should be addressed to R.X. or Y.S.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023