

1 **Corporate environmental performance prediction in China: An empirical study of energy**
2 **service companies**

3

4 Saina Zheng¹, Chenhang He², and Shu-Chien Hsu^{3,*}, Joseph Sarkis⁴, and Jieh-Haur Chen⁵

5

6 **Abstract**

7 Businesses are constrained by and dependent upon nature and institutional context. The global
8 climate crisis has put pressure on and increased firm sensitivity to environmental issues.
9 Predicting corporate environmental performance can help plan for environmental impact
10 mitigation by adjusting organizational practices. Lack of environment-related information
11 makes it difficult to make such predictions. A theoretical framework informed by the natural-
12 resource-based view (NRBV) of the firm and institutional theory is used to identify variables
13 for predicting corporate environmental performance. Five dimensions including institutional
14 context, governance capability, information management capability, system capability, and
15 technology-related capability, populated with 14 variables are used to empirically investigate
16 the relationship of these variables with corporate environmental performance. Using 1100 data
17 points on energy service companies (ESCOs) from 2011 to 2015 in mainland China, the
18 Extreme Gradient Boosting (XGBoost) algorithm, a statistical nonlinear machine learning
19 approach, is utilized to predict corporate environmental performance. The results demonstrate
20 that the XGBoost model can be effective for ESCO environmental performance prediction,
21 with satisfactory prediction accuracy. This study also adopted the SHapley Additive

¹ Ph.D. candidate, Department of Civil and Environmental Engineering, The Hong Kong Polytechnic University.

² Mphil candidate, Department of Electronic and Information Engineering, The Hong Kong Polytechnic University.

^{3*} Corresponding author: Assistant Professor, Department of Civil and Environmental Engineering, The Hong Kong Polytechnic University.

E-mail: mark.hsu@polyu.edu.hk; Tel: +852-27666057

⁴ Professor, Foisie Business School, Worcester Polytechnic Institute.

⁵ Distinguished Professor, Department of Civil Engineering, National Central University.

22 exPlanations (SHAP) values for model interpretation, indicating that total assets, amount of
23 proactive environmental costs, proportion of technicians and number of patents contribute most
24 to corporate environmental performance. Several policies and environmental strategies for
25 improving corporate environmental performance in the ESCO industry are derived from this
26 analysis.

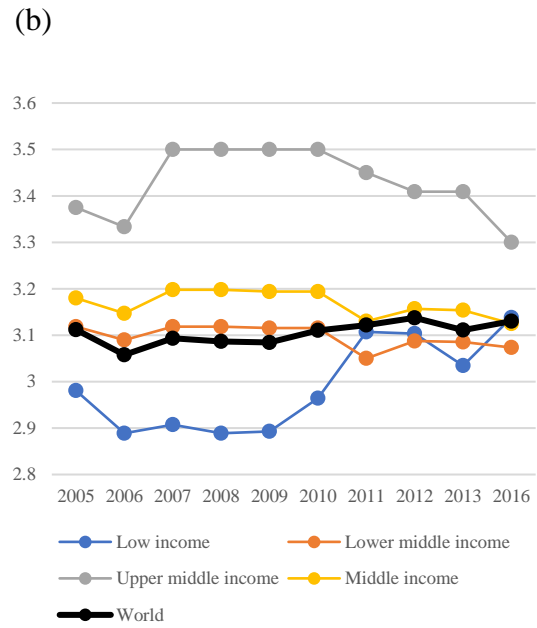
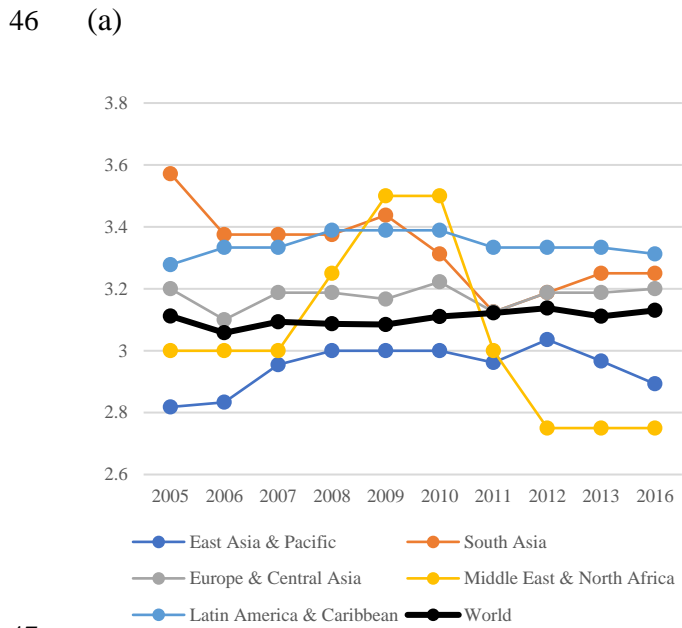
27

28 **Keywords:** Corporate environmental performance, Extreme Gradient Boosting (XGBoost),
29 Energy service company (ESCO).

30

31 **1 Introduction**

32 Growing climate and other environmental crises, such as resource depletion, have led to an
33 increased focus on shifting to a low-carbon world (Millar et al., 2018). Policies have focused
34 on attempts to cut or mitigate greenhouse gas (GHG) emissions, as policy measures are the
35 most direct way to reduce the risk of future climate change impacts (IPCC, 2015). Countries
36 are joining global environmental collaborative efforts including the Kyoto Protocol, the
37 Copenhagen Accord, and the Global Pact for the Environment. Several market-based
38 environmental instruments including green credits, green insurance, and pollution tax policies
39 have been adopted (Crowley, 2013; Garnaut, 2008; Neuhoff, 2011; Newell and Paterson, 2010;
40 Nyberg et al., 2013; Stern, 2008). According to the Country Policy and Institutional
41 Assessment (CPIA), ratings on policies and institutions for environmental sustainability of the
42 world have continued to rise since 2005, as shown in Figure 1(a). Figure 1(b) indicates that
43 countries with higher income tend to pay more attention to policies concerning environmental
44 sustainability. With its rapid economic development, China is expected to launch more
45 environmental policies.



48 Figure 1: CPIA ratings on policy and institutions for environmental sustainability (1=low to
 49 6=high). (a) CPIA ratings on policies and institutions for environmental sustainability by
 50 region; (b) CPIA ratings on policies and institutions for environmental sustainability according
 51 to country income level. (Data source: World Bank)

52 China became the world’s largest carbon emitter in 2006 and has increased attention to
 53 environmental issues arising from its population growth and economic development (Liu et al.,
 54 2016; Zhang et al., 2008). Past environmental policy in China focused on mandatory
 55 regulations. In recent years this role has shifted to market-oriented and voluntary approaches.
 56 Fiscal incentive policies, tax subsidies, a pollution levy system, and technology innovation
 57 support are being provided seeking to achieve economic and environmental protection ‘win-
 58 wins’ (Zhang et al., 2007). Energy performance contracting (EPC) is one of these instruments.
 59 EPC is a market-oriented approach in which the energy service companies (ESCOs) invest in
 60 implementing energy services for customers to improve energy efficiency, including energy
 61 savings guarantees, associated design, and installation services. ESCOs get paid annually from
 62 energy savings during the contract period (Deng et al., 2017; Zheng et al., 2018). ESCOs are

63 corporations which focus on improving energy efficiency and relieving climate change through
64 EPC (Liu et al., 2018; Xu et al., 2015; Xu and Chan, 2013).

65 Economic growth has been identified as the main driver for sharp CO₂ emissions
66 increases. Anthropogenic causes of climate change are intimately related to economic
67 behaviour, and the industry is increasingly being called upon to respond (Yeeles, 2018). The
68 sheer scale of the Chinese economy means that worldwide CO₂ emissions are strongly
69 determined there (Wiedenhofer et al., 2017). The enterprises are the primary damager of
70 environmental pollution and the major consumer of energy in China (Li et al., 2017). Nearly
71 two-thirds of China's groundwater was of poor quality, over 15% of China's soil and farmland
72 has been polluted, causing serious threat to food security and human health (Li et al., 2017;
73 Qiu, 2011). To reduce these impacts, regulating corporate environmental performance in China
74 is in urgent need.

75 The institutional theory stipulates that firms will respond to institutional pressures
76 (mainly regulatory policy pressures) to thrive and gain legitimacy (Meyer and Rowan, 1977;
77 Scott, 2013). Damaging corporate environmental impact has become increasingly part of the
78 public and social mindset. This concern over climate change and the policy regulations to
79 mitigate GHG emissions have exerted greater pressures on corporations to improve
80 environmental well-being rather than hastening its degradation. These forces have also
81 incentivized corporations to pursue good environmental performance (Hart, 2010).

82 The natural-resource-based view (NRBV) of the firm stipulates that a business will be
83 constrained by natural resources and firm competitiveness is related to natural resources (Hart,
84 1995). Many corporations now claim that developing a corporate culture that promotes social
85 and environmental sustainability can improve employee recruitment attraction, motivation, and
86 retention (Renwick et al., 2013). Corporations are also realizing the significance of climate
87 change and develop strategies for this environmental issue (Wright and Nyberg, 2015).

88 Business and industry play a dual role in climate politics. Firstly, corporations are the principal
89 agents producing CO₂ emissions; secondly, corporations can improve the environment and
90 reduce emissions through technological innovation. Better environmental performance reduces
91 the volatility of the firm's cash flows, decreases potential bankruptcy costs, and increases debt
92 capacity; all characteristics that can add resources for an organization to build competitive
93 advantage. Many theoretical and empirical studies also indicate that better environmental
94 performance boosts and is endogenously influenced by better financial performance (Dixon-
95 Fowler et al., 2013; King and Lenox, 2001; Stanwick and Stanwick, 1998).

96 Disclosure of corporate environmental performance has been advocated to achieve
97 better environmental performance. Scholars have argued that information about corporate
98 environmental performance disclosed to the public can play an important role in determining
99 business strategies, consumers' purchasing behaviour, and investors' financial investment
100 decisions (Meng et al., 2014; Rockness, 1985; Spicer, 1978). Environmental disclosure may
101 decrease the agency costs of debt and reduce estimation or information risk (Bansal and
102 Clelland, 2004; Gao and Connors, 2011). However, a vast majority of companies do not
103 produce corporate environmental reports or include environmental information in their annual
104 reports. This result may be due to environmental information disclosure resistance, a desire to
105 avoid additional costs, fear of threats to local employment, and concerns about reduced profits
106 (Wang et al., 2004).

107 Researchers have been investigating corporate environmental performance for decades.
108 They have mainly focused on evaluating environmental performance and environmental
109 management strategy using the environmental-related information (Bhatnagar, 1999; Delmas
110 and Blass, 2010; Ilinitich et al., 1998; Klassen and McLaughlin, 1996; Lober, 1996; Tyteca et
111 al., 2002; Zhang et al., 2008). In China, various corporations have made efforts toward
112 environmental protection and generating data related to environmental performance, especially

113 given governmental pressures for this type of information. This data provides information
114 which can be useful for practical governmental and organizational policy decisions, but also
115 for research purposes. However, this type of data is currently rare in China, resulting in some
116 hurdles in predicting corporate environmental performance using environment-related
117 information. Only a few studies emphasize predicting environmental performance, due to the
118 scant data available and lack of a detailed list of pollutants emitted by corporations (Delmas
119 and Blass, 2010). It is essential find ways to utilize corporate related information, which can
120 be accessed easily, for predicting corporate environmental performance. An application using
121 machine learning method-XGBoost can help in completing various predictive analyses for
122 multiple settings and purposes, especially when there exists the sparse and noise in the dataset.

123 In this study, a machine learning model for predicting corporate environmental
124 performance is constructed with two main functionalities: assessing the corporate
125 environmental performance of an unknown ESCO and calculating the future performance of
126 ESCOs. Predicting corporate environmental performance can be used to mitigate
127 environmental impacts through guiding organizational practices, and to improve a firm's
128 reputation. First, we combine the elements of Institutional Theory and NRBV to bring a fresh
129 perspective to environmental performance prediction research. Second, we explore the factors
130 in five domains: Institutional context, Governance capability, Information management
131 capabilities, Systems capability, and Technology-related investment, using 1100 data points on
132 corporate environmental performance from different industries in mainland from 2011 to 2015
133 in China. Lastly, we utilize a machine learning tool, XGBoost regression to predict future
134 performance and adopt Shapely to generate interpretations from the model. Conclusions and
135 future research finalize the paper.

136

137 **2 Theory foundation and hypotheses development**

138 Corporate environmental performance can be defined as the results of an organisation's
139 management of its environmental aspects or more precisely 'is the totality of a firm's behaviour
140 toward the natural environment (i.e. it's level of total resource consumption and emissions)'
141 (Tyteca et al., 2002). Corporations compete over limited natural resources, tend to take
142 strategies to use the resources more efficiently, relieve their impact on the natural environment,
143 and focus more effort on pollution control. Corporate environmental performance evaluation
144 has been proposed for self-assessment, benchmarking, and reporting (Delmas and Blass, 2010;
145 Gao and Connors, 2011; Ilinitch et al., 1998; Veleva and Ellenbecker, 2001).

146 Several theories have been used to investigate and explain corporate environmental
147 performance. These theories include the natural-resource-based view (NRBV) (Hart, 1995),
148 institutional theory (Colwell and Joshi, 2013; Jennings and Zandbergen, 1995), stakeholder
149 theory (Freeman, 1984), agency theory (Berrone and Gomez-Mejia, 2009; Friedman, 2007),
150 and transaction cost theory; to name a few organizational theories. Two of these are especially
151 popular and salient. One is a general external to the organization theory, institutional theory,
152 the other is an internal theory used to build competitive advantage, the NRBV. Together these
153 two theories provide a more complete picture of how organizations manage their environmental
154 performance.

155

156 *2.1 Combining Institutional Theory and the Natural-Resource-Based View*

157 The institutional theory posits that organizations enhance or seek to protect their legitimacy
158 (Scott, 2013) by conforming to the expectations of institutional norms and stakeholder
159 requirements (Aldrich and Fiol, 1994; DiMaggio and Powell, 2000). Concern over legitimacy
160 forces firms to adopt managerial practices that are expected to conform to social values and
161 expectations (Berrone and Gomez-Mejia, 2009). With the increasing importance of

162 environmental issues, institutional theory stipulates that companies under heavier institutional
163 pressure will gain legitimacy by exhibiting good environmental performance (Bansal and
164 Clelland, 2004; Bansal Pratima, 2005). Researchers have applied institutional theory in the
165 investigation of corporate environmental performance (Berrone and Gomez-Mejia, 2009;
166 Campbell, 2007; Gallego-Alvarez et al., 2017; Tashman and Rivera, 2016). The institutional
167 context has a significant influence on environmental performance and the adoption of
168 environmental strategies (Chang et al., 2015; Christmann, 2004; Russo and Fouts, 1997;
169 Sharfman et al., 2004; Wang et al., 2018). Under institutional pressures, firms have tended to
170 adopt appropriate strategies and firms with an environmental legacy has incurred less risk
171 (Bansal and Clelland, 2004).

172 NRBV (Hart, (1995) and holds that the business is constrained by and dependent upon
173 natural ecosystems. Organizational competitiveness relies on the capabilities which facilitate
174 environmentally sustainable economic activity. Many researchers have examined the
175 relationship between corporate environmental performance and financial performance
176 (McWilliams and Siegel, 2001; Stanwick and Stanwick, 1998). Al-Tuwaijri et al (2004)
177 provided an analysis of the interrelations between environmental performance and economic
178 performance, finding that good environmental performance is significantly associated with
179 good economic performance. Trumpp and Guenther (2017) build on the theory of a non-linear,
180 specifically a U-shaped, relationship between corporate environmental performance and
181 corporate financial performance.

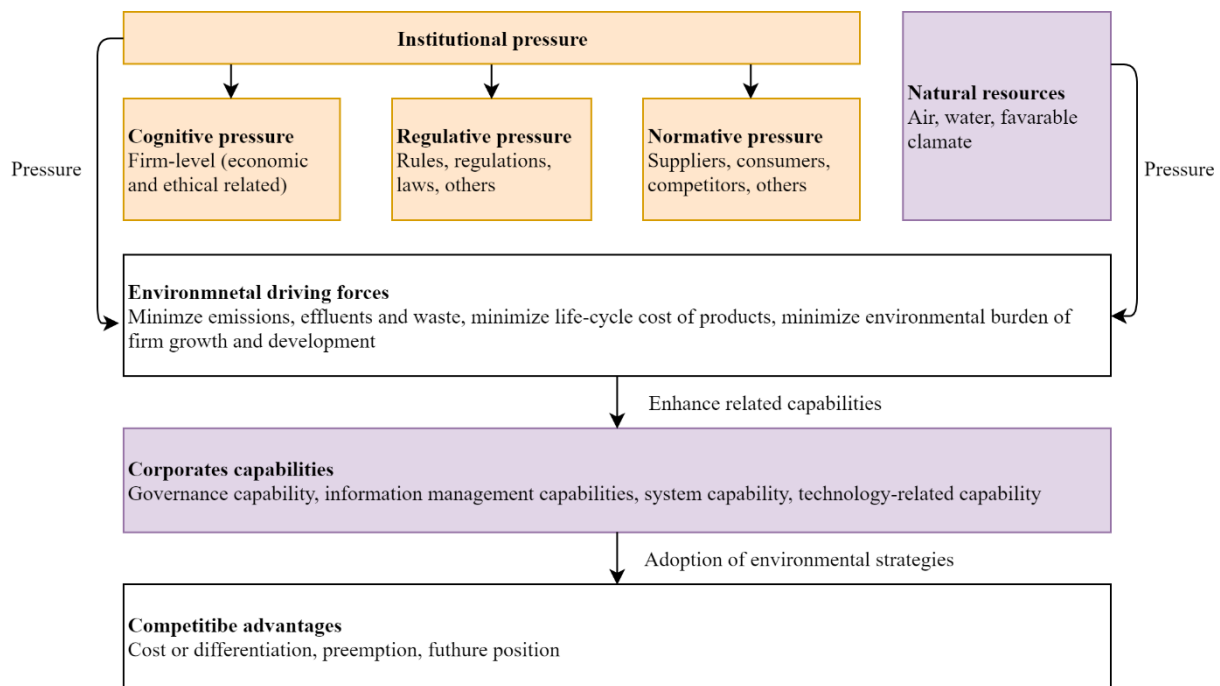
182 Empirical work by Buysse and Verbeke (2003) identified five essential resource
183 domains through which environmental proactiveness can be determined: strategic
184 environmental planning, formal routine-based environmental management, organizational
185 competencies in environmental management, employees' green skills, and conventional
186 technology-based green competencies. These determinants have also been categorized into

187 four resource domains: governance capability, information management capabilities, systems
188 capability, and technology-related investment (Backman et al., 2017). These elements will
189 prove helpful in our investigation using our machine learning models.

190 Scholars across a wide range of research areas and disciplines have focused on
191 examining the relationship between corporate environmental performance and other
192 organizational constructs or variables (Bansal and Gao, 2006; Trumpp et al., 2015). These
193 variables include external organizational factors, such as regulation (Camisón, 2010) or
194 stakeholder pressure (Ilinitch et al., 1998), as well as internal organizational factors, such as
195 different characteristics of the board (Post et al., 2015) or innovation (Hall and Wagner, 2012).
196 These findings corroborate the NRBV and institutional theory, dividing corporate
197 environmental performance into five categories: institutional context, governance capability,
198 information management capabilities, systems capability, and technology-related capability.
199 We conceptualize how the possible combinatorial configurations from institutional theory and
200 the NRBV relate to corporate environmental performance and try to understand how
201 corporations may achieve improved environmental performance.

202 Figure 2 summarizes the integration of institutional theory and NRBV for corporate
203 environmental performance. The figure depicts how corporations will choose to engage in
204 environmentally friendly behaviour due to limited natural resources in a given institutional
205 context. Institutional theory is adopted to explain how organizations react to institutional
206 pressures, while NRBV encompasses building corporate capabilities in such a way to gain a
207 competitive advantage in the market given natural resources consideration (Delmas and Toffel,
208 2004; Hart, 1995). Institutional theory categorizes the institutional pressure into three types,
209 namely, cognitive, regulative and normative pressure (Gao et al., 2019). The cognitive pressure
210 is related to the economic and ethical aspects which mainly refers to the environmental benefit
211 and ethical obligation (Gao et al., 2019). The regulative pressure comes from the regulations,

212 laws, rules and other formal instruments. The normative pressure is those pressure exerted by
213 the nongovernmental stakeholders, such as suppliers, consumers, competitors (DiMaggio and
214 Powell, 2000; Gao et al., 2019; Lee et al., 2018). Firms depend directly on natural capital and
215 ecosystem services (Pogutz & Winn, 2009; Starik & Rands, 1995). Without air, water, a
216 favourable climate, and a variety of natural resources, no organization can survive (Gladwin et
217 al., 1995). Key resources and capabilities also affect organizational ability to adopt competitive
218 environmental strategies (Hart, 1995). With the increasing consequences of climate change and
219 growing severity of resource scarcity, firms are facing loss of access to natural resources and
220 must adapt according to their dependence on nature. Both these institutional pressures and
221 natural resources pressure drive the organization to minimize emissions, effluents, waste, life-
222 cycle environmental costs of products, and environmental burden of firm growth and
223 development. That is, organizations adopt environmental strategies to achieve higher corporate
224 environmental performance. Corporate capabilities are the key factors which will affect the
225 adoption of environmental strategies. As stated by NRBV, these capabilities consist of
226 governance capability, information management, systems, and technology-related capabilities.
227 Corporates will adopt environmental strategies to pursue competitive advantages in the market
228 since researchers found that firms with better environmental performance have superior
229 financial performance (McWilliams et al., 2006). Environmental strategies can also lead to
230 reduced costs and improved environmental performance simultaneously (Lyon and Maxwell,
231 1999).



232

233 **Figure 2:** Framework of corporate environmental relationships and performance based on
 234 Institutional Theory and the Natural-Resource-Based View of the firm.
 235

236 *2.2 Institutional context and corporate environmental performance*

237 Institutional forces play a significant role in corporate environmental strategy adoption (Chang
 238 et al., 2015); impacting corporate environmental performance. High-income regions of the
 239 world have generally exerted the strongest regional or national institutional pressure for
 240 improved environmental performance. Institutional pressure has been found to be lower in
 241 middle-income and lessens in lowest-income regions (Luxmore et al., 2018)). The reason for
 242 this is an overwhelming need for economic development in some regions, where institutional
 243 measures from an environmental perspective may be lessened.

244 In China, developed regions and their governmental agencies, invest more in
 245 improving energy efficiency and require organizations to emphasize environmental
 246 performance (Zheng et al., 2018). In this study, the gross domestic product (GDP) and the
 247 population are proxies for institutional (government) pressure since this metric reflects the
 248 development of a region. Natural resource availability may also exert pressure on corporate

249 environmental performance from NRBV. Consumption of coal is used to reflect the level of
250 dependency on natural resources; renewable energies may also be dependent on natural
251 resources but are not as constrained due to the continuous sources (e.g. sunlight and wind
252 power). Using these perspectives, the following hypothesis is proposed for testing.

253 *Hypothesis 1 (H1):* Corporates facing greater institutional and natural resource pressures will
254 exhibit greater improvements in corporate environmental performance.

255 *2.3 Organizational characteristics and corporate environmental performance*

256 There are multiple levels of pressures and contexts. The first hypothesis focused on broader
257 social and natural resource considerations and relationships to environmental performance.
258 Organizational (corporate) contexts and characteristics will also relate to corporate
259 environmental performance. Corporate specific internal resources and capabilities are
260 particularly useful in generating unique, preventive and voluntary environmental actions to
261 reduce firms' environmental impacts (Hart, 1995). According to NRBV, corporate
262 characteristics can be divided into four domains: governance, information management, system,
263 and technology-related capabilities.

264 (1) Governance capability indicators

265 Governance capability refers to a strategic planning process reconfiguration ability and
266 integration of environmental issues into corporate policies and routines (Backman et al., 2017;
267 Walls et al., 2012). A measure of governance capability and policy focus includes
268 environmental costs which are the investments made in addressing pollution issues and
269 adopting environment strategies (Salo, 2008);. Another, proxy measure includes the number of
270 formal legal warnings a firm has received since its founding. This warnings measure indicates
271 how well the governance structure supports good or poor behaviour and can be closely linked
272 environmental policy; e.g. going beyond compliance (Li et al., 2017).

273 Firms with a high level of environmental commitment and stronger governance policies
274 are more likely to regard environmental protection as their corporate social responsibility and
275 be eager to protect the environment, thus achieving higher corporate environmental
276 performance (Al-Tuwaijri et al., 2004; Muller and Kolk, 2010; Wang et al., 2018). Furthermore,
277 organizations with historically poor environmental records are often subjected to more scrutiny
278 by their local communities and regulators. Thus, organizations with poor environmental
279 records may try to build greater corporate environmental governance capabilities to achieve
280 higher environmental performance to gain more resources. Together, the following is expected.
281 *Hypothesis 2 (H2):* Greater organizational environmental governance capability, measured by
282 the combination of environmental cost and formal legal warning, relates to higher corporate
283 environmental performance.

284 (2) Information management capability

285 Information management capabilities mainly focus on formal management systems and
286 procedures of investment. Researchers found that effective information management
287 capabilities and corporate social responsibility are synergistically related, and can facilitate
288 transition to corporate sustainability (Gangi et al., 2019). Countries paying attention to climate
289 change mitigation tend to set develop stronger information management capability for
290 organizations (Backman et al., 2017). This information refers to environmental-related
291 information, such as the climate change impact mitigation and carbon footprint, denoting the
292 attitude towards sustainability. The work environment is considered to evaluate the
293 organizational culture regarding how a corporation views the importance of environment
294 (Bhatnagar, 1999). High environmental awareness can help firms to implement environmental
295 management practices smoothly and then help them improve environmental performance.
296 Based on this analysis, the following is hypothesized.

297 *Hypothesis 3 (H3):* Corporates with stronger information management capability tend to
298 achieve higher corporate environmental performance than corporates with weaker information
299 management capability.

300 (3) Systems capability

301 Systems capability covers investments in employee skills and organizational competencies,
302 such as research and development funding, finance and accounting, and storage and human
303 resources in environmental management (Backman et al., 2017; Buysse and Verbeke, 2003).
304 Previous research investigated the relationship between organizational characteristic variables
305 and environmental performance/environmental benefits, such as the top management's
306 leadership skills, human resources and organizational size (Etzion, 2007; Lee et al., 2018).
307 Kitada and Ölçer (2015) put forward that employee element is essential when considering
308 corporate social responsibility. It was found that there appears to be a positive relationship
309 between a firm's environmental performance and its financial performance (Dixon-Fowler et
310 al., 2013; Rockness, 1985; Spicer, 1978). Accordingly, we postulate the following.

311 *Hypothesis 4 (H4):* Corporations with stronger system capability tend to achieve higher
312 corporate environmental performance than corporates with weaker system capability.

313 (4) Technology-related capability

314 Technology-related capability covers the conventional green competencies related to green
315 product and manufacturing technologies. Technologies will affect corporate competitiveness
316 since the environmental problems arise increasing awareness (Shrivastava, 1995). Technology
317 in energy efficiency proved options and solutions for organizations to pursue better
318 environmental performance by implementing the energy efficiency retrofit projects (Kitada and
319 Ölçer, 2015). Benitez-Amado and Walczuch (2012) believed that technology-related
320 capabilities are a key enabler for organizations to achieve better environmental performance.
321 Environmental innovations contribute to corporate environmental performance since they can

322 improve energy efficiency and reduce pollution. Therefore, our last hypothesis is as following
323 (Kagan et al., 2003).

324 *Hypothesis 5 (H5):* Corporates with stronger technology-related capability tend to achieve
325 higher corporate environmental performance than corporates with weaker technology-related
326 capability.

327 **3 Research method and data processing**

328 *3.1 Sample and data*

329 The combination of institutional theory and NRBV identifies five domains for selecting the
330 index to predict corporate environmental performance. Considering the data availability and
331 referring to the previous research, 14 factors as shown in Table 1 were chosen to test the
332 hypotheses. The research data we employed was provided by the ESCO Committee of China
333 Energy Conservation Association (EMCA). The collected data covers 3225 ESCOs in 30
334 provinces in mainland China from 2011 to 2015 (Zheng et al., 2018). However, some ESCOs
335 were excluded for one or more of the following reasons: (i) data for environmental performance
336 is missing; (ii) data for more than 3 variables are missing. Thus, 1134 ESCO projects have
337 sufficient information for further analysis. The value of the corporate environmental
338 performance for most projects are between 0 and 1, however, the corporate environmental
339 performance of 34 projects (3% of total projects) is 0, meaning there is no environmental
340 income for these companies, which is not suitable for our research. Then, 1100 ESCO projects
341 are finally analysed to predict the corporate environmental performance, which is mainly
342 located in the Beijing, Shandong, and Guangdong provinces (see Fig. 3). Table 2 shows an
343 example of the detailed information for each project, including investment, number of formal
344 legal warnings since foundation, proportion of in-plant environmental, proportion of
345 technicians, assets, equity, environmental projects payback period, asset age, revenue, tax
346 bracket, and number of patents. All the variables in Table 1 can get or calculated based on

347 Table 2. The amount of proactive environmental costs is the investment for improving energy
 348 efficiency and reducing the impact of environment. The proportion of technician can be get
 349 using the number of technicians divided by number of employees. Information related to GDP,
 350 population, and consumption of coal was gained through the National Bureau of Statistics of
 351 China.

352

353 Table 1: Corporate environmental performance indicator system.

Destination layer	Standard layer	Index layer	Data source	References
Corporate environmental performance prediction indicator system	Institutional context	GDP (GDP)	National Bureau of Statistics of China	(Chan and Makino, 2007; Zheng et al., 2018)
		Population (PO)	National Bureau of Statistics of China	(Cui and Jiang, 2012)
		Consumption of coal (CC)	National Bureau of Statistics of China	(Zheng et al., 2018)
	Governance capability	Amount of proactive environmental costs (PEC)	EMCA	(Fu et al., 2017; Salo, 2008)
		Number of formal legal warnings since firm founding (FLW)	EMCA	(Li et al., 2017; Yoon et al., 2006)
	Information management capability	Proportion of In-plant environment (PIE)	EMCA	(Bhatnagar, 1999)
	Systems capability	Proportion of technicians (PT)	EMCA, $\frac{\text{number of technicians}}{\text{number of employees}}$	(Etzion, 2007; Lee et al., 2018)
		Total assets (TA)	EMCA	(Backman et al., 2017; Buysse and

				Verbeke, 2003)
		Equity (EQ)	EMCA	(Backman et al., 2017; Buysse and Verbeke, 2003)
		Environmental projects payback period (PP)	EMCA	(Dibrell et al., 2011)
		Asset age (AA)	EMCA	(Li et al., 2017)
		Revenue (RE)	EMCA	(Orlitzky et al., 2003; Russo and Fouts, 1997)
		Tax bracket (TB)	EMCA	(Hoi et al., 2013)
	Technology-related capability	Number of patents (PA)	EMCA	(Benitez-Amado and Walczuch, 2012)

354

355

356 Table 2: Example of detailed information about ESCO

Liaoning Nengfaweiye Energy Technology Co., Ltd.	Region	Number of Employees	Number of Technicians	Number of Patents	Investment (million yuan)	Assets
	Liaoning	450	68	13	27.53	240.23
	Equity	Payback Period	Asset Age	Ratepaying Credit Grade	Number of Penalties Received	Environmental Performance (Environmental income per unit of an asset)
	20.205	0.8	6	A	0	0.354502

357

358



359

360 Figure 3: Distribution of Sampled ESCOs

361

362 These samples cover all kinds of firms, including state-owned enterprises, corporations,
 363 general partnership firms, private enterprises, foreign-owned enterprises, and others, with their
 364 assets varying from 0 to more than 1 trillion yuan (shown in Table 3).

365

366 Table 3: Sampled firms by business type and asset size

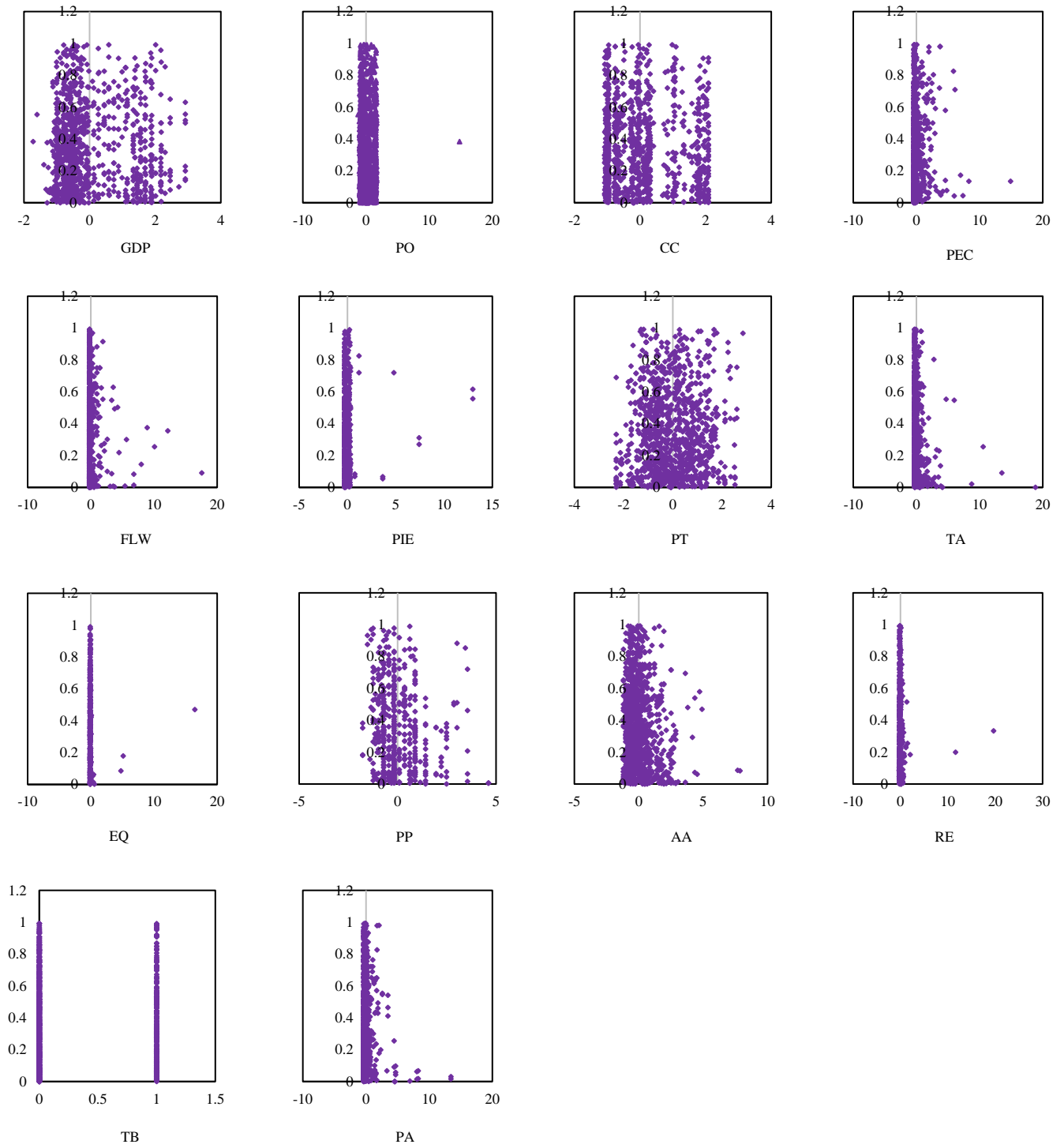
Business type	Number	Assets (in yuan)	Number
State-owned enterprise	157	0-500	20
Corporation	128	500-1000	202
General Partnership	54	1000-5000	531
Private enterprise	736	5000-10000	171
Foreign-owned enterprise	14	10000-100000	171
Other	11	≥100000	5

367

368 Fig. 4 displays the scatter plots for each of the (normalized) input variables and output variables.

369 These scatter plots show that none of the functional relationships between the input variables

370 and the output variables are trivial.



371

372 Figure 4: Scatter plots of the relationships between each input variable and output

373 This suggests that we can reasonably accept that classical learners such as linear regression
374 may fail to find an accurate mapping of the input variables to the output variables. Therefore,
375 these plots intuitively justify the need to experiment with more complicated learners such as
376 machine learning methods. However, the machine learning methods are mainly used for
377 prediction and classification, without the ability to interpret the relationship between variables.
378 Recently, SHAP (SHapley Additive exPlanations) was developed to interpret the variables'
379 impact on the model's prediction (Lundberg and Lee, 2017). A SHAP value for a feature of a
380 specific prediction represents how much the model prediction changes when we observe that
381 feature. This SHAP figure not only indicates which features are most important but also their
382 range of effects over the dataset, revealing the relationship between variables and model output.

383 *3.2 Machine learning method-XGBoost*

384 In recent years, machine learning has been generating a lot of curiosity for its superior
385 performance compared to other more traditional statistical techniques. Numerous machine
386 learning models like Linear/Logistic regression, Support Vector Machines, Neural Networks,
387 and Tree-based models are being tried and applied in analysis and prediction (Gumus and Kiran,
388 2017). Tso and Yau (2007) predicted electricity energy consumption adopting the decision tree
389 and neural network models. Lee (2007) applied support vector machines to suggest a new
390 model for corporate credit ratings with better explanatory power and stability. Among these
391 methods, Extreme Gradient Boosting, also known as XGBoost (Chen and Guestrin, 2016), is
392 a model that has a high success rate in the majority of machine learning competitions and has
393 proven to be efficient for predictive modeling.

394 XGBoost has algorithms that can deeply explore data-label correlations by adaptively
395 fitting large-scale data via tree boosting. Compared to conventional regression approaches such
396 as logistic regression and SVM regression, XGBoost's tree-ensemble approaches can easily
397 handle data with missing values (Torlay et al., 2017). Third, XGBoost penalizes the complexity

398 of an individual tree as a regularization term, which has better generalization ability compared
399 to other MART (multiple additive regress trees) methods. Ajit and Punnoose R (2016) applied
400 XGBoost to predict employee turnover within an organization, addressing the prevalence of
401 noise in data to reduce overfitting and improve accuracy. XGBoost is suitable for our case since
402 there exist sparse data and noisy data in the realm of corporate environmental performance.
403 Furthermore, the tree-ensemble algorithm provides strong interpretability of the model. By
404 constructing the model, we can visualize the tree's structure and explore implicitly how the
405 model makes decisions and which attributes are dominant.

406 XGBoost is a typical tree-ensemble model related to CART (Classification And
407 Regression Trees), which grows the tree in a top-down manner. Each tree consists of internal
408 (or split) nodes and terminal (or leaf) nodes. Each split node will make a binary decision and
409 the final decision is made based on the terminal node reached by the input feature. Tree-
410 ensemble methods regard different decision trees as weak learners, and then construct a strong
411 learner by either bagging or boosting. Bagging, also known as bootstrap aggregating (Breiman,
412 1996), is used to reduce the variance of the model. Multiple random subsets of the dataset with
413 replacements are first selected, one for training an individual sub-model. Then an average
414 prediction from these sub-models is calculated. Random Forest (Liaw and Wiener, 2002)
415 extends the bagging by exploiting a small tweak that reduces the correlation between the
416 bagged trees. For the boosting algorithm, the boosted tree (strong learner) is regarded as a
417 combination of the single trees (weak learners). The weight of the combination is updated
418 adaptively according to the different designs of the objective function and optimization
419 methods. AdaBoost (Freund and Schapire, 1997) is the first version of the boosting method, in
420 which the weak learners are iteratively trained on a weighted dataset by minimizing the
421 exponential loss. XGBoost extends to more general loss function via gradient boosting
422 optimization and learns a model with an additive training trick.

423 The objective of XGBoost is to learn a model with good variance-bias balance. In other
 424 words, the model should have strong predictive power but also large variance to be generalized
 425 on the extra data. This can be represented with the following objective function with respect to
 426 model parameter θ :

$$427 \quad \text{obj}(\theta) = L(\theta) + \Omega(\theta)$$

428 where the first term is the loss function which should be minimized, and the second term is a
 429 regularization term of the model's complexity to prevent it from over-fitting. Considering a
 430 tree-ensemble model where the overall prediction is the summation of K predictive values
 431 across all the trees $f_k(x_i)$,

$$432 \quad p_i = \sum_{k=1}^K f_k(x_i),$$

433 the objective function can be written as:

$$434 \quad \text{obj}(\theta) = \sum_i^n l(p_i, t_i) + \sum_{k=1}^K \Omega(f_k),$$

435 where $l(p_i, t_i)$ is the mean-squared loss imposed on each sample i regarding its predictive
 436 value p_i and the label t_i , and $\Omega(f_k)$ is the regularization constraint imposed on each tree.
 437 XGBoost applies an efficient additive training algorithm to optimize such an objective
 438 function. This algorithm will learn one tree at each step, then add a new tree by fixing what it
 439 has learned, mathematically,

$$440 \quad p_i^{(0)} = 0,$$

$$441 \quad p_i^{(1)} = f_1(x_i) = p_i^{(0)} + f_1(x_i),$$

$$442 \quad \dots \dots$$

$$443 \quad p_i^{(t)} = \sum_{k=1}^t f_k(x_i) = p_i^{(t-1)} + f_t(x_i).$$

444 Thus, the objective at step t becomes,

$$\begin{aligned}
445 \quad \text{obj}^{(t)} &= \sum_{i=1}^n \left(t_i - \left(p_i^{(t-1)} + f_t(x_i) \right) \right)^2 + \sum_{i=1}^t \Omega(f_i) \\
446 \quad &= \sum_{i=1}^n \left[2 \left(p_i^{(t-1)} - t_i \right) f_t(x_i) + f_t(x_i)^2 \right] + \Omega(f_t) + \text{constant}
\end{aligned}$$

447 This can be easily optimized with second-order Tylor expansion, considering the first and
448 second-order gradients, $g_i = \partial_{p_i^{(t-1)}} l(t_i, p_i^{(t-1)})$ and $h_i = \partial_{p_i^{(t-1)}}^2 l(t_i, p_i^{(t-1)})$ respectively, with
449 the objective function at step t now becoming,

$$450 \quad \text{obj}^{(t)} \approx \sum_{i=1}^n \left[g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \Omega(f_t)$$

451 To this end, we introduced how to efficiently train the boosted trees with an additive
452 strategy, i.e. training a new tree at a step t by optimizing above step-based objective function.
453 One of the merits of this definition is that the objective value only depends on the g_i and h_i ,
454 which allows using custom loss function. $\Omega(f_t)$ is the regularization term, which controls the
455 complexity of the model. Now, we re-define the tree by a vector of prediction score in leaves,

$$456 \quad f_t(x) = w_{q(x)}, w \in \mathbb{R}^T$$

457 where $q(x_i)$ is a mapping function that maps a training instance to a leaf. Based on this re-
458 defined formulation, $\Omega(f)$ can be heuristically defined as

$$459 \quad \Omega(f) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2,$$

460 where T is the number of leaves of the tree and w_j is the prediction score in each leaf. By re-
461 grouping the training samples on each leaf j , the objective function can hence be reformed as

$$\begin{aligned}
462 \quad \text{obj}^{(t)} &\approx \sum_{i=1}^n \left[g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \\
463 \quad &= \sum_{j=1}^T \left[\left(\sum_{i \in I_j} g_i \right) w_j + \frac{1}{2} \left(\sum_{i \in I_j} h_i + \lambda \right) w_j^2 \right] + \gamma T,
\end{aligned}$$

464 where $I_j = \{i | q(x_i) = j\}$ is the indices of training instances which reach the j th leaf. We use
 465 $G_j = \sum_{i \in I_j} g_i$ and $H_j = \sum_{i \in I_j} h_i$ to express the summation of first/second ordered gradients
 466 across leaves. Thus, the objective function can then be further simplified as

$$467 \quad obj^{(t)} \approx \sum_{j=1}^T \left[G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2 \right] + \gamma T$$

468 Note that w_j are independent with each other, thus the equation has a quadratic form, the
 469 solution for the above equation is

$$470 \quad w_j^* = -\frac{G_j}{H_j + \lambda},$$

471 and the resulting objective value is

$$472 \quad obj^* = -\frac{1}{2} \sum_{j=1}^T \frac{G_j^2}{H_j + \lambda} + \gamma T.$$

473 This equation measures how good a tree structure $q(x_i)$ is for a certain training instance.
 474 Based on this property, one can grow a tree greedily using the information gain. To specify this
 475 information gain, we consider the gradients flow before and after splitting,

$$476 \quad G_L = \sum_{j \in T_L} g_j, \quad G_R = \sum_{j \in T_R} g_j$$

$$477 \quad H_L = \sum_{j \in T_L} h_j, \quad H_R = \sum_{j \in T_R} h_j$$

478 where T_L and T_R are the indices of left and right leaves respectively. Before splitting, the tree's
 479 complexity is

$$480 \quad \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} + \gamma.$$

481 After splitting, the tree has complexity,

$$482 \quad \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} + 2\gamma,$$

483 Then the information gain of a splitting tree can be calculated as

484
$$Gain = \frac{1}{2} \left[\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma$$

485 As a result, we can outline the *XGBoost* algorithm as an iteration process. For each iteration,
486 we perform the following operations: 1) Grow the tree to the maximum depth by finding the
487 best splitting points via information gain. 2) Assign prediction score to the two new leaves. 3)
488 Prune the tree by deleting the nodes with negative gain.

489 **4 Implementation of XGBoost – a reliable prediction model**

490 There exists some sparse data in our experimental dataset that needs the adoption of XGBoost.
491 The dataset was arbitrarily split into two subsets; 75% of the data was used as a training set
492 and 25% as a validation set. All the training data for Xbgoost was used to construct the model.
493 The validation data was used to test the results with the data that was not utilized to develop
494 the model. In order to improve the calculation efficiency, and prevent individual data from
495 overflowing during the calculation, input and output parameters were normalized. In addition,
496 all 14 variables show independence from each other after doing correlation analysis, which
497 indicates these 14 variables can be used for predicting the environmental performance in a
498 model.

499 PyCharm was adopted to train and develop the XGBoost model for corporate environmental
500 performance. A statistical package scikit-learn in python was used to implement the XGBoost.
501 To determine the hyper-parameters of the model, we applied a brute force grid search with 5-
502 fold cross-validation. In order to achieve optimal parameter setting, we needed to initialize the
503 searching with some prior knowledge of the parameters' ranges. For example, the learning rate
504 for XGBoost is usually 0.05, and the maximum depth is usually 6, 7, or 8. Other parameters,
505 such as 'min_child_weight', 'subsample', and 'colsample_bytree' need to be carefully tuned
506 since they greatly affect the model's generalizability. Thus, we applied different seed during
507 the searching to increase the variance of the model. The boosting iterations were determined

508 using early stopping, and mean squared error was applied as the evaluation metrics during the
 509 searching. Table 4 shows the finally determined values for the hyper-parameters of the
 510 XGBoost model which achieve the best performance.

511 Table 4: Values Determined for the Hyper-parameters of the XGBoost Model

	Description	Value
'eta'	Boosting learning rate	0.03
'subsample'	Subsample ratio of the training instance	0.8
'colsample_bytree'	Subsample ratio of columns when constructing each tree	0.8
'objective'	Specify the learning task and the corresponding learning objective	'linear'
'max_depth'	Maximum tree depth for base learners	7
'min_child_weight'	Minimum sum of instance weight(hessian) needed in a child	0.5
'num_boost_round'	Number of boosting iterations	1000

512

513 4.1 Evaluation Criteria for model

514 The performance evaluation indices for the models tested in this paper are Mean Absolute Error
 515 (MAE), Root Mean Square Error (RMSE), Correlation Coefficient (CC) and Coefficient of
 516 Determination (R-square, R^2), which are defined as follows:

$$517 \quad MAE = \frac{1}{m} \sum_{i=1}^m |y_i - \hat{y}_i|$$

$$518 \quad RMSE = \sqrt{\frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2}$$

$$519 \quad CC = \frac{\sum_{i=1}^m (x_i - \bar{x}_i)(y_i - \bar{y}_i)}{\sqrt{\sum_{i=1}^m (x_i - \bar{x}_i)^2 \sum_{i=1}^m (y_i - \bar{y}_i)^2}}$$

$$520 \quad R^2 = 1 - \frac{\sum_{i=1}^m (y_i - \hat{y}_i)^2}{\sum_{i=1}^m (y_i - \bar{y}_i)^2}$$

521 where y_i is the observed value for parameter y , \hat{y}_i is the predicted value and \bar{y}_i is the mean of
 522 observed values.

523

524 *4.2 Reliability of XGBoost model*

525 Random Forest (RF) and Support Vector Machines for regression (SVMreg) are commonly
526 adopted machine learning methods when dealing with prediction problems (Chaudhuri and De,
527 2011; Lee, 2007; Pan, 2018; Tsanas and Xifara, 2012). In this research, RF and SVMreg
528 prediction methods were implemented and compared with an XGBoost model.

529 Fig. 5 presents the initial data curve and relative error curve of the training set and testing
530 data. For the training course curve and testing course curve, a dot was extracted from the curve
531 every 10 samples, 88 samples in total. And for the training error and testing error curves, a dot
532 was extracted from each curve every 3 samples, 73 samples each in total. It can be seen that
533 the prediction relative errors of the training samples under the XGBoost model are nearly
534 0.04%, exhibiting much better performance compared to SVMReg and RF. This demonstrates
535 that the developed XGBoost model can more precisely describe the complex relationship
536 between corporate environmental performance and explanatory variables. The predicted
537 environmental performance on validation data by the three models and the relative errors
538 between the predicted value and real value are illustrated in Fig. 5(b), Fig. 5(d), and Fig. 5(f).
539 The MAE, RMSE, CC, and R^2 of the testing samples under the three models are compared in
540 Table 5.

541 Fig. 5 and Table 5 also show that using the XGBoost method to predict corporate
542 environmental performance is better than using RF and SVMreg. The XGBoost method is more
543 efficient and is a reliable alternative for corporate environmental performance prediction.

544

545



546

547 Figure 5: Course curves and relative error curves. (a) training course curve of initial data and
 548 three models, (b) testing course curve of initial data and three models, (c) training error curve
 549 of three models, (d) testing error curve of three models, (e) XGBoost error curve on training
 550 data, and (f) XGBoost error curve on testing data.

551

552 Table 5: Comparison of the prediction accuracy of SVMReg, RF, and XGBoost

Method	MAE	RMSE	CC	R2
SVMreg	0.19304	0.23527	0.39971	0.15951
RF	0.16578	0.20429	0.61295	0.36630
XGBoost	0.14546	0.18336	0.70244	0.48952

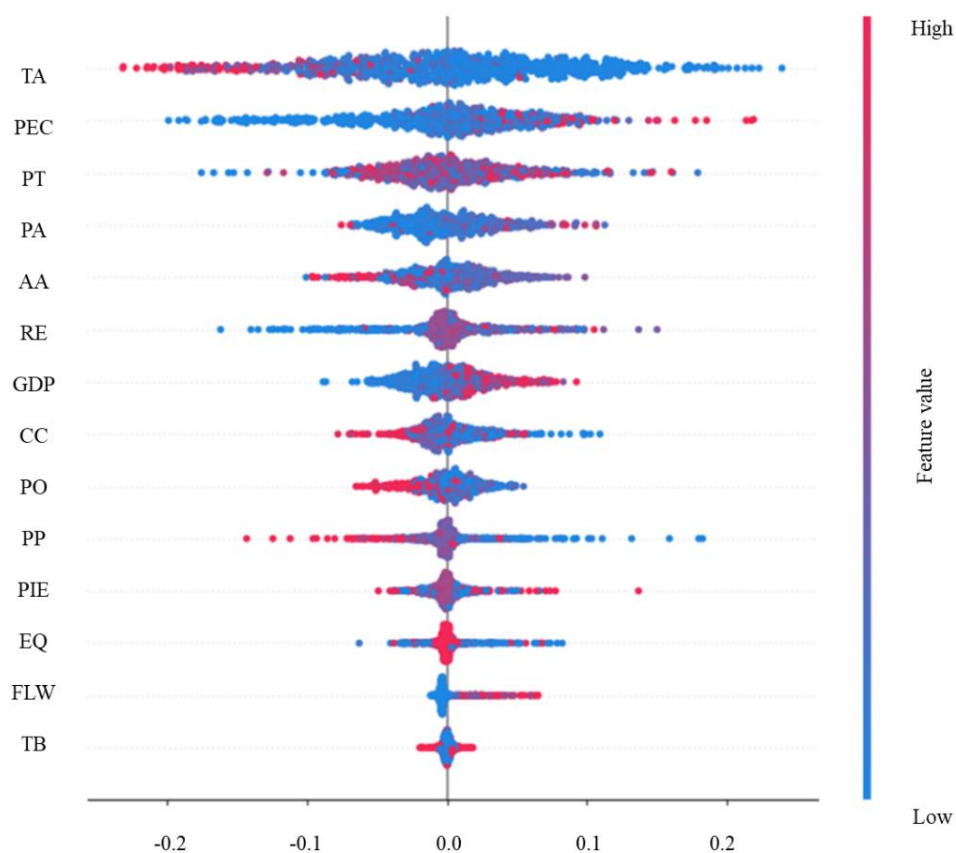
553

554 5 Empirical results and discussions

555 A SHAP value for a feature of a specific prediction represents how much the model prediction
 556 changes when we observe that feature. In the summary plot below (Fig. 6), all the SHAP values

557 for a single feature on a row are drawn, where the x-axis is the SHAP value (which for this
558 model is in units of log odds of corporate environmental performance).

559 Fig. 6 indicates that total asset (TA), amount of proactive environmental costs (PEC),
560 proportion of technicians (PT) and number of patents (PA) were more important in this model
561 while tax bracket (TB), formal legal warning since firm founding (FLW), equity (EQ),
562 proportion of in-plant environment (PIE) and environmental projects payback period (PP) were
563 relatively less important.



564

565 Figure 6: Summary of SHAP values for 14 variables (impact on model output)
566

567 This SHAP figure not only indicates which features are most important but also their range of
568 effects over the dataset. Each dot is coloured by the value of that feature from high to low. For
569 example, as shown in Fig. 6, red dots for the rows of 'total asset' tend to appear on the left side.
570 This means that high values for total asset lead to lower corporate environmental performance,

571 or in other words, total asset exhibit a positive relationship to corporate environmental
572 performance. As shown in Fig. 6, GDP shows a positive relationship with corporate
573 environmental performance, while population and consumption of coal have a negative effect
574 on corporate environmental performance. Based on the analysis above, Hypothesis 1 cannot be
575 supported. Amount of proactive environmental costs and number of formal legal warnings
576 since firm founding reflect a positive relationship with corporate environmental performance.
577 These results provide support for Hypothesis 2. It is revealed that the higher value for
578 proportion of in-plant environment, the better environmental performance corporates will get.
579 As such, Hypothesis 3 is supported. There indicates a negative effect on corporate
580 environmental performance for total assets, environmental projects payback period and asset
581 age. In contrast, a positive relationship exists between proportion of technicians, equity,
582 revenue, and corporate environmental performance. Fewer total assets correlated with higher
583 corporate environmental performance. As for tax bracket, the relationship reveals unclear. Thus,
584 Hypothesis 4 cannot be supported. Fig. 6 clearly indicates that number of patents shows a
585 positive relationship with corporate environmental performance, verifying Hypothesis 5.

586 Total assets (TA) represents the size and ability of a firm, which is highly related to
587 corporate environmental performance (Trumpp Christoph and Guenther Thomas, 2017; Zhang
588 et al., 2008). Amount of proactive environmental costs (PEC) is a direct reflection of
589 investment in the environmental strategy of a firm. The proportion of technicians (PT) and
590 number of patents (PA) show the technological innovation ability of a firm. Advanced
591 technology can reduce the environmental impact of firms, improve energy efficiency, and
592 increase corporate environmental performance (Dietz and Rosa, 1994; Wang et al., 2013). The
593 variable explanation rankings show the firm characteristic variables explain more about the
594 prediction model, indicating that Natural-resource-based view is better to study the corporate
595 environmental performance. As researchers pointed out before, the external environment,

596 including normative and regulative environmental, remains undeveloped and fragmented in
597 China (Gao et al., 2019). Thus, the Institutional theory explained less about corporate
598 environmental performance in China.

599 There are massive reasons leading to the complex relationship between variables and
600 corporate environmental performance. Amount of proactive environmental costs is the direct
601 investment in environmental strategy, and the result is similar with the previous research which
602 indicates that greater investment leads to higher corporate environmental performance (Fu et
603 al., 2017). The negative relationship between formal legal warnings since firm founding could
604 due to that the corporates need to keep its positive image. If a firm receives a legal warning
605 that damages its social image, it may adopt measures to mitigate this effect, such as developing
606 and implementing environmental strategies.

607 The negative relationship between total assets and corporate environmental performance
608 stands in contrast to previous research findings (Al-Tuwaijri et al., 2004). This may be because
609 the proportion of assets dedicated to environmental investment by large firms is relatively low
610 although the large firms care more about social responsibility (Udayasankar, 2008). Firms
611 survive based on their profitability, which enforces firms to invest in profitable projects. As for
612 the payback period, firms tend to invest in projects with short payback period to avoid the risks.
613 If the payback period is too long, firms will engage in these projects, leading to less
614 environmental improvement projects and poorer environmental performance. Revenue can
615 show the profitability of a firm and firms with strong profitability tend to pay more attention
616 to environmental issues (Orlitzky et al., 2003; Russo and Fouts, 1997). Regions with higher
617 GDP are usually developed regions in China, such as Beijing, Shanghai, Jiangsu, and
618 Guangdong. These regions are stricter about environmental protection and have inaugurated
619 several policies regarding environmental sustainability (Zheng et al., 2018).

620

621 Based on the findings of positive and negative relationships between variables and
622 corporate environmental performance, several policies can be put in place to improve corporate
623 environmental performance. To increase firms' contributions to corporate proactive
624 environmental costs, the government should provide more financial incentives for
625 environmental protection, including tax benefits, green loans, and environmental subsidies.
626 When providing these incentives, the payback period should be taken into consideration since
627 longer periods entail more risks. Currently many corporations have spent hundreds of millions
628 of dollars on environmental projects (Berry and Rondinelli, 1998). Fines for noncompliance
629 need to be increased and enforcement of environmental regulations should be strengthened,
630 making business executives and owners liable for environmental pollution. The number of fines
631 and intensity of enforcement also need to be applied in accordance with the size of the
632 corporation. Corporate environmental issues need to receive more attention in regions with
633 lower GDP. The local governments of these regions can learn from the experiences of more
634 developed regions. At present, the incentives are primarily provided to the larger-scale
635 corporations since they demonstrate better financial performance. However, as knowledge,
636 practices, systems, and routines at the business and natural environment interface become more
637 widely dispersed, smaller companies may also begin to adopt voluntary niche environmental
638 strategies (Sharma, 2000). Governments ought to promote the environmental strategies of
639 small companies and develop some targeted incentives tailored to them.

640 Whatever policies the government may implement, corporations could internally choose
641 to direct more investment toward environmental prevention and minimization. Introducing
642 advanced technologies, employing more technicians, reusing materials, and adopting an
643 environmental corporate culture are other advisable measures. Full-cost accounting is
644 suggested for adoption when managing the corporations, considering direct costs (labour,
645 capital, and raw materials), hidden costs (monitoring and reporting;), contingent liability costs

646 (fines and remedial action), and less tangible costs (public relations and goodwill).
647 Corporations can use full-cost accounting to choose the most eco-effective projects and
648 improve corporate environmental performance.

649

650 **6 Conclusion**

651 In this study, we identified the relationship between institutional context, corporate
652 environmental performance with corporate environmental performance based on a
653 combination of Institutional Theory and the Natural-Resource-Based View. We presented an
654 approach to conducting this identification by predicting corporate environmental performance
655 with machine learning methods. The key challenge of dealing with noise in the data from
656 ESCOs that compromises the accuracy of these predictive models was also highlighted. In this
657 study, a newly introduced machine learning algorithm, XGBoost, was applied to predict
658 corporate environmental performance. Data from 1100 projects for ESCOs in the time period
659 between 2011 and 2015 was analyzed to explore the statistical relationship between 14 input
660 variables (GDP, population, consumption of coal, amount of proactive environmental costs,
661 number of formal legal warnings since a firm's founding, proportion of in-plant area,
662 proportion of technicians, total worth, equity, environmental projects payback period, asset age,
663 revenue, tax bracket, number of patents) and the output variable, corporate environmental
664 performance. The results indicate that XGBoost achieved higher accuracy than other learning
665 algorithms and was reliable to test the relationship.

666 The findings of this research agree with those in the machine learning literature strongly
667 endorsing the use of XGBoost in complex applications (Gumus and Kiran, 2017; Pan, 2018).
668 Applying SHAP in XGBoost model interpretation enables the impact of input variables on the
669 output to be determined. In the model, total assets (TA), amount of proactive environmental
670 costs (PEC), proportion of technicians (PT) and number of patents (PA) are found to contribute

671 the most to corporate environmental performance. Also, the impacts each feature has on the
672 model output was obtained through SHAP summary plotting. Amount of proactive
673 environmental costs (PEC), Revenue (RE), GDP, and number of formal legal warnings since
674 the firm's founding (FLW) show a positive relationship with corporate environmental
675 performance, while total assets (TA) and environmental projects payback period (PP) show a
676 negative relationship. Based on the SHAP findings, several policy recommendations and
677 environmental strategies for governments and corporations to carry out are proposed to
678 improve corporate environmental performance. Corporates with stronger governance
679 capability, information management capability and technology-related capability will perform
680 better corporate environmental performance.

681 Although this paper contributes to corporate environmental performance, there are still
682 some research limitations. First, the prediction accuracy for all observations is relatively low
683 due to the result of noisy data and the limited input gaps between machine learning and social
684 sciences (CHEN et al., 2018). The rate may increase if more information about corporate
685 environmental performance is considered. Second, the data used are only from the ESCO
686 industry in China. It could add more value if the model can be tested in other industries and in
687 other countries.

688 In future studies, more variables and more data should be introduced to achieve greater
689 accuracy in predicting corporate environmental performance. Since corporates in China are
690 becoming increasingly aware of the environmental performance, along with increasing national
691 policies regarding corporate environmental performance, more variables from the perspective
692 of institutional theory can be taken into consideration. In addition, due to the contrast results
693 with previous research about the relationship between total assets and corporate environmental
694 performance, the total assets could be considered to have a nonlinear relationship as the

695 moderators when investigating the relationship between corporate financial performance and
696 corporate environmental performance.

697

698 **References**

699 Ajit, P., Punnoose R, 2016. Prediction of Employee Turnover in Organizations using Machine
700 Learning Algorithms. *Int. J. Adv. Res. Artif. Intell.* 5, 22–26.

701 Aldrich, H.E., Fiol, C.M., 1994. Fools Rush in? The Institutional Context of Industry Creation. *Acad.*
702 *Manage. Rev.* 19, 645–670. <https://doi.org/10.5465/AMR.1994.9412190214>

703 Al-Tuwaijri, S.A., Christensen, T.E., Hughes, K.E., 2004. The relations among environmental
704 disclosure, environmental performance, and economic performance: a simultaneous equations
705 approach. *Account. Organ. Soc.* 29, 447–471. [https://doi.org/10.1016/S0361-3682\(03\)00032-1](https://doi.org/10.1016/S0361-3682(03)00032-1)

706 Backman, C.A., Verbeke, A., Schulz, R.A., 2017. The Drivers of Corporate Climate Change
707 Strategies and Public Policy: A New Resource-Based View Perspective. *Bus. Soc.* 56, 545–575.
708 <https://doi.org/10.1177/0007650315578450>

709 Bansal, P., Clelland, I., 2004. Talking Trash: Legitimacy, Impression Management, and Unsystematic
710 Risk in the Context of the Natural Environment. *Acad. Manage. J.* 47, 93–103.
711 <https://doi.org/10.2307/20159562>

712 Bansal, P., Gao, J., 2006. Building the Future by Looking to the Past: Examining Research Published
713 on Organizations and Environment. *Organ. Environ.* 19, 458–478.
714 <https://doi.org/10.1177/1086026606294957>

715 Bansal Pratima, 2005. Evolving sustainably: a longitudinal study of corporate sustainable
716 development. *Strateg. Manag. J.* 26, 197–218. <https://doi.org/10.1002/smj.441>

717 Benitez-Amado, J., Walczuch, R.M., 2012. Information technology, the organizational capability of
718 proactive corporate environmental strategy and firm performance: a resource-based analysis. *Eur. J.*
719 *Inf. Syst.* 21, 664–679. <https://doi.org/10.1057/ejis.2012.14>

720 Berrone, P., Gomez-Mejia, L.R., 2009. Environmental Performance and Executive Compensation: An
721 Integrated Agency-Institutional Perspective. *Acad. Manage. J.* 52, 103–126.

722 Berry, M.A., Rondinelli, D.A., 1998. Proactive Corporate Environmental Management: A New
723 Industrial Revolution. *Acad. Manag. Exec.* 1993-2005 12, 38–50.

724 Bhatnagar, V., 1999. Evaluating corporate environmental performance in developing countries.
725 *Sustain. Meas. Eval. Report. Environ. Soc. Perform.* Greenleaf Publ. Sheffield S3 8GG UK.

726 Breiman, L., 1996. Bagging predictors. *Mach. Learn.* 24, 123–140.

727 Buysse, K., Verbeke, A., 2003. Proactive environmental strategies: a stakeholder management
728 perspective. *Strateg. Manag. J.* 24, 453–470. <https://doi.org/10.1002/smj.299>

729 Camisón, C., 2010. Effects of coercive regulation versus voluntary and cooperative auto-regulation on
730 environmental adaptation and performance: Empirical evidence in Spain. *Eur. Manag. J.* 28, 346–361.
731 <https://doi.org/10.1016/j.emj.2010.03.001>

732 Campbell, J.L., 2007. Why Would Corporations Behave in Socially Responsible Ways? An
733 Institutional Theory of Corporate Social Responsibility. *Acad. Manage. Rev.* 32, 946–967.
734 <https://doi.org/10.2307/20159343>

735 Chan, C.M., Makino, S., 2007. Legitimacy and multi-level institutional environments: implications
736 for foreign subsidiary ownership structure. *J. Int. Bus. Stud.* 38, 621–638.
737 <https://doi.org/10.1057/palgrave.jibs.8400283>

738 Chang, L., Li, W., Lu, X., 2015. Government Engagement, Environmental Policy, and Environmental
739 Performance: Evidence from the Most Polluting Chinese Listed Firms. *Bus. Strategy Environ.* 24, 1–
740 19. <https://doi.org/10.1002/bse.1802>

741 Chaudhuri, A., De, K., 2011. Fuzzy Support Vector Machine for bankruptcy prediction. *Appl. Soft
742 Comput., The Impact of Soft Computing for the Progress of Artificial Intelligence* 11, 2472–2486.
743 <https://doi.org/10.1016/j.asoc.2010.10.003>

744 CHEN, N.-C., DROUHARD, M., KOCIELNIK, R., SUH, J., ARAGON, C.R., 2018. Using Machine
745 Learning to Support Qualitative Coding in Social Science: Shifting The Focus to Ambiguity.

746 Chen, T., Guestrin, C., 2016. XGBoost: A Scalable Tree Boosting System, in: *Proceedings of the
747 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16.*
748 ACM, New York, NY, USA, pp. 785–794. <https://doi.org/10.1145/2939672.2939785>

749 Christmann, P., 2004. Multinational Companies and the Natural Environment: Determinants of Global
750 Environmental Policy. *Acad. Manage. J.* 47, 747–760. <https://doi.org/10.2307/20159616>

751 Colwell, S.R., Joshi, A.W., 2013. Corporate Ecological Responsiveness: Antecedent Effects of
752 Institutional Pressure and Top Management Commitment and Their Impact on Organizational
753 Performance. *Bus. Strategy Environ.* 22, 73–91. <https://doi.org/10.1002/bse.732>

754 Crowley, K., 2013. Irresistible force? Achieving carbon pricing in Australia. *Aust. J. Polit. Hist.* 59,
755 368–381.

756 Cui, L., Jiang, F., 2012. State ownership effect on firms' FDI ownership decisions under institutional
757 pressure: a study of Chinese outward-investing firms. *J. Int. Bus. Stud.* 43, 264–284.
758 <https://doi.org/10.1057/jibs.2012.1>

759 Delmas, M., Blass, V.D., 2010. Measuring corporate environmental performance: the trade-offs of
760 sustainability ratings. *Bus. Strategy Environ.* 19, 245–260. <https://doi.org/10.1002/bse.676>

761 Delmas, M., Toffel, M.W., 2004. Stakeholders and environmental management practices: an
762 institutional framework. *Bus. Strategy Environ.* 13, 209–222. <https://doi.org/10.1002/bse.409>

763 Deng, X., Zheng, S., Xu, P., Zhang, X., 2017. Study on dissipative structure of China's building
764 energy service industry system based on brusselator model. *J. Clean. Prod.* 150, 112–122.
765 <https://doi.org/10.1016/j.jclepro.2017.02.198>

766 Dibrell, C., Craig, J., Hansen, E., 2011. Natural Environment, Market Orientation, and Firm
767 Innovativeness: An Organizational Life Cycle Perspective. *J. Small Bus. Manag.* 49, 467–489.
768 <https://doi.org/10.1111/j.1540-627X.2011.00333.x>

769 Dietz, T., Rosa, E.A., 1994. Rethinking the Environmental Impacts of Population, Affluence and
770 Technology. *Hum. Ecol. Rev.* 1, 277–300.

771 DiMaggio, P.J., Powell, W.W., 2000. The iron cage revisited institutional isomorphism and collective
772 rationality in organizational fields, in: *Economics Meets Sociology in Strategic Management*. Emerald
773 Group Publishing Limited, pp. 143–166.

774 Dixon-Fowler, H.R., Slater, D.J., Johnson, J.L., Ellstrand, A.E., Romi, A.M., 2013. Beyond “Does it
775 Pay to be Green?” A Meta-Analysis of Moderators of the CEP—CFP Relationship. *J. Bus. Ethics* 112,
776 353–366.

777 Etzion, D., 2007. Research on Organizations and the Natural Environment, 1992-Present: A Review.
778 *J. Manag.* 33, 637–664. <https://doi.org/10.1177/0149206307302553>

779 Freeman, R.E., 1984. *Strategic management: a stakeholder approach*. Pitman, Boston.

780 Freund, Y., Schapire, R.E., 1997. A decision-theoretic generalization of on-line learning and an
781 application to boosting. *J. Comput. Syst. Sci.* 55, 119–139.

782 Friedman, M., 2007. The Social Responsibility of Business Is to Increase Its Profits, in: Zimmerli,
783 W.C., Holzinger, M., Richter, K. (Eds.), *Corporate Ethics and Corporate Governance*. Springer Berlin
784 Heidelberg, Berlin, Heidelberg, pp. 173–178. https://doi.org/10.1007/978-3-540-70818-6_14

785 Fu, X., Lahr, M., Yaxiong, Z., Meng, B., 2017. Actions on climate change, Intended Reducing carbon
786 emissions in China via optimal industry shifts: Toward hi-tech industries, cleaner resources and
787 higher carbon shares in less-develop regions. *Energy Policy*.
788 <https://doi.org/10.1016/j.enpol.2016.10.038>

789 Gallego-Alvarez, I., Ortas, E., Vicente-Villardón, J.L., Etxeberria, I.Á., 2017. Institutional
790 Constraints, Stakeholder Pressure and Corporate Environmental Reporting Policies. *Bus. Strategy*
791 *Environ.* 26, 807–825. <https://doi.org/10.1002/bse.1952>

792 Gangi, F., Mustilli, M., Varrone, N., 2019. The impact of corporate social responsibility (CSR)
793 knowledge on corporate financial performance: evidence from the European banking industry. *J.*
794 *Knowl. Manag.* <https://doi.org/10.1108/JKM-04-2018-0267>

795 Gao, L.S., Connors, E., 2011. Corporate Environmental Performance, Disclosure and Leverage: An
796 Integrated Approach. *Int. Rev. Account. Bank. Finance* 1.

797 Gao, Y., Gu, J., Liu, H., 2019. Interactive effects of various institutional pressures on corporate
798 environmental responsibility: Institutional theory and multilevel analysis. *Bus. Strategy Environ.* 28,
799 724–736. <https://doi.org/10.1002/bse.2276>

800 Garnaut, R., 2008. *The Garnaut climate change review*. Camb. Camb.

801 Gumus, M., Kiran, M.S., 2017. Crude oil price forecasting using XGBoost, in: *2017 International*
802 *Conference on Computer Science and Engineering (UBMK)*. Presented at the 2017 International

803 Conference on Computer Science and Engineering (UBMK), pp. 1100–1103.
804 <https://doi.org/10.1109/UBMK.2017.8093500>

805 Hall, J., Wagner, M., 2012. Integrating Sustainability into Firms' Processes: Performance Effects and
806 the Moderating Role of Business Models and Innovation. *Bus. Strategy Environ.* 21, 183–196.
807 <https://doi.org/10.1002/bse.728>

808 Hart, S.L., 2010. Capitalism at the crossroads: Next generation business strategies for a post-crisis
809 world. FT Press.

810 Hart, S.L., 1995. A Natural-Resource-Based View of the Firm. *Acad. Manage. Rev.* 20, 986–1014.
811 <https://doi.org/10.5465/AMR.1995.9512280033>

812 Hoi, C.K., Wu, Q., Zhang, H., 2013. Is Corporate Social Responsibility (CSR) Associated with Tax
813 Avoidance? Evidence from Irresponsible CSR Activities. *Account. Rev.* 88, 2025–2059.
814 <https://doi.org/10.2308/accr-50544>

815 Ilinitich, A.Y., Soderstrom, N.S., E. Thomas, T., 1998. Measuring corporate environmental
816 performance. *J. Account. Public Policy* 17, 383–408. [https://doi.org/10.1016/S0278-4254\(98\)10012-1](https://doi.org/10.1016/S0278-4254(98)10012-1)

817 IPCC, 2015. Climate change 2014: mitigation of climate change. Cambridge University Press.

818 Jennings, P.D., Zandbergen, P.A., 1995. Ecologically Sustainable Organizations: An Institutional
819 Approach. *Acad. Manage. Rev.* 20, 1015–1052. <https://doi.org/10.2307/258964>

820 Kagan, R.A., Gunningham, N., Thornton, D., 2003. Explaining Corporate Environmental
821 Performance: How Does Regulation Matter? *Law Soc. Rev.* 37, 51–90. <https://doi.org/10.1111/1540-5893.3701002>

822

823 King, A.A., Lenox, M.J., 2001. Does It Really Pay to Be Green? An Empirical Study of Firm
824 Environmental and Financial Performance: An Empirical Study of Firm Environmental and Financial
825 Performance. *J. Ind. Ecol.* 5, 105–116. <https://doi.org/10.1162/108819801753358526>

826 Kitada, M., Ölçer, A., 2015. Managing people and technology: The challenges in CSR and energy
827 efficient shipping. *Res. Transp. Bus. Manag., Energy Efficiency in Maritime Logistics Chains* 17, 36–
828 40. <https://doi.org/10.1016/j.rtbm.2015.10.002>

829 Klassen, R.D., McLaughlin, C.P., 1996. The Impact of Environmental Management on Firm
830 Performance. *Manag. Sci.* 42, 1199–1214. <https://doi.org/10.1287/mnsc.42.8.1199>

831 Lee, J.W., Kim, Y.M., Kim, Y.E., 2018. Antecedents of Adopting Corporate Environmental
832 Responsibility and Green Practices. *J. Bus. Ethics* 148, 397–409. <https://doi.org/10.1007/s10551-016-3024-y>

833

834 Lee, Y.-C., 2007. Application of support vector machines to corporate credit rating prediction. *Expert*
835 *Syst. Appl.* 33, 67–74. <https://doi.org/10.1016/j.eswa.2006.04.018>

836 Li, D., Zhao, Y., Sun, Y., Yin, D., 2017. Corporate environmental performance, environmental
837 information disclosure, and financial performance: Evidence from China. *Hum. Ecol. Risk Assess.*
838 *Int. J.* 23, 323–339. <https://doi.org/10.1080/10807039.2016.1247256>

839 Liaw, A., Wiener, M., 2002. Classification and regression by randomForest. *R News* 2, 18–22.

840 Liu, G., Zheng, S., Xu, P., Zhuang, T., 2018. An ANP-SWOT approach for ESCOs industry strategies
841 in Chinese building sectors. *Renew. Sustain. Energy Rev.* 93, 90–99.
842 <https://doi.org/10.1016/j.rser.2018.03.090>

843 Liu, Z., Davis, S.J., Feng, K., Hubacek, K., Liang, S., Anadon, L.D., Chen, B., Liu, J., Yan, J., Guan,
844 D., 2016. Targeted opportunities to address the climate–trade dilemma in China. *Nat. Clim. Change* 6,
845 201. <https://doi.org/10.1038/nclimate2800>

846 Lober, D.J., 1996. Evaluating The Environmental Performance Of Corporations. *J. Manag. Issues* 8,
847 184–205.

848 Lundberg, S.M., Lee, S.-I., 2017. A Unified Approach to Interpreting Model Predictions, in: Guyon,
849 I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (Eds.),
850 *Advances in Neural Information Processing Systems* 30. Curran Associates, Inc., pp. 4765–4774.

851 Luxmore, S.R., Hull, C.E., Tang, Z., 2018. Institutional Determinants of Environmental Corporate
852 Social Responsibility: Are Multinational Entities Taking Advantage of Weak Environmental
853 Enforcement in Lower-Income Nations? *Bus. Soc. Rev.* 123, 151–179.
854 <https://doi.org/10.1111/basr.12138>

855 Lyon, T.P., Maxwell, J.W., 1999. Corporate environmental strategies as tools to influence regulation.
856 *Bus. Strategy Environ.* 8, 189–196. [https://doi.org/10.1002/\(SICI\)1099-0836\(199905/06\)8:3<189::AID-BSE194>3.0.CO;2-0](https://doi.org/10.1002/(SICI)1099-0836(199905/06)8:3<189::AID-BSE194>3.0.CO;2-0)

858 McWilliams, A., Siegel, D., 2001. Corporate social responsibility: A theory of the firm perspective.
859 *Acad. Manage. Rev.* 26, 117–127.

860 McWilliams, A., Siegel, D.S., Wright, P.M., 2006. Corporate Social Responsibility: Strategic
861 Implications*. *J. Manag. Stud.* 43, 1–18. <https://doi.org/10.1111/j.1467-6486.2006.00580.x>

862 Meng, X.H., Zeng, S.X., Shi, J.J., Qi, G.Y., Zhang, Z.B., 2014. The relationship between corporate
863 environmental performance and environmental disclosure: An empirical study in China. *J. Environ.
864 Manage.* 145, 357–367. <https://doi.org/10.1016/j.jenvman.2014.07.009>

865 Meyer, J.W., Rowan, B., 1977. Institutionalized Organizations: Formal Structure as Myth and
866 Ceremony. *Am. J. Sociol.* 83, 340–363.

867 Millar, R.J., Hepburn, C., Beddington, J., Allen, M.R., 2018. Principles to guide investment towards a
868 stable climate. *Nat. Clim. Change* 8, 2. <https://doi.org/10.1038/s41558-017-0042-4>

869 Muller, A., Kolk, A., 2010. Extrinsic and Intrinsic Drivers of Corporate Social Performance: Evidence
870 from Foreign and Domestic Firms in Mexico. *J. Manag. Stud.* 47, 1–26.
871 <https://doi.org/10.1111/j.1467-6486.2009.00855.x>

872 Neuhoff, K., 2011. *Climate policy after Copenhagen: the role of carbon pricing.* Cambridge
873 University Press.

874 Newell, P., Paterson, M., 2010. *Climate capitalism: global warming and the transformation of the
875 global economy.* Cambridge University Press.

876 Nyberg, D., Spicer, A., Wright, C., 2013. Incorporating citizens: corporate political engagement with
877 climate change in Australia. *Organization* 20, 433–453.

878 Orlitzky, M., Schmidt, F.L., Rynes, S.L., 2003. Corporate Social and Financial Performance: A Meta-
879 Analysis. *Organ. Stud.* 24, 403–441. <https://doi.org/10.1177/0170840603024003910>

880 Pan, B., 2018. Application of XGBoost algorithm in hourly PM2.5 concentration prediction. *IOP*
881 *Conf. Ser. Earth Environ. Sci.* 113, 012127. <https://doi.org/10.1088/1755-1315/113/1/012127>

882 Post, C., Rahman, N., McQuillen, C., 2015. From Board Composition to Corporate Environmental
883 Performance Through Sustainability-Themed Alliances. *J. Bus. Ethics* 130, 423–435.
884 <https://doi.org/10.1007/s10551-014-2231-7>

885 Qiu, J., 2011. China to Spend Billions Cleaning Up Groundwater. *Science* 334, 745–745.
886 <https://doi.org/10.1126/science.334.6057.745>

887 Renwick, D.W., Redman, T., Maguire, S., 2013. Green human resource management: A review and
888 research agenda. *Int. J. Manag. Rev.* 15, 1–14.

889 Rockness, J.W., 1985. An Assessment of the Relationship Between Us Corporate Environmental
890 Performance and Disclosure. *J. Bus. Finance Account.* 12, 339–354. <https://doi.org/10.1111/j.1468-5957.1985.tb00838.x>

892 Russo, M.V., Fouts, P.A., 1997. A Resource-Based Perspective On Corporate Environmental
893 Performance And Profitability. *Acad. Manage. J.* 40, 534–559. <https://doi.org/10.2307/257052>

894 Salo, J., 2008. Corporate Governance and Environmental Performance: Industry and Country Effects.
895 *Compet. Change* 12, 328–354. <https://doi.org/10.1179/102452908X357293>

896 Scott, W.R., 2013. *Institutions and organizations: Ideas, interests, and identities.* Sage Publications.

897 Sharfman, M.P., Shaft, T.M., Tihanyi, L., 2004. A Model of the Global and Institutional Antecedents
898 of High-Level Corporate Environmental Performance. *Bus. Soc.* 43, 6–36.
899 <https://doi.org/10.1177/0007650304262962>

900 Sharma, S., 2000. Managerial Interpretations and Organizational Context as Predictors of Corporate
901 Choice of Environmental Strategy. *Acad. Manage. J.* 43, 681–697. <https://doi.org/10.5465/1556361>

902 Shrivastava, P., 1995. Environmental technologies and competitive advantage. *Strateg. Manag. J.* 16,
903 183–200. <https://doi.org/10.1002/smj.4250160923>

904 Spicer, B.H., 1978. Investors, Corporate Social Performance and Information Disclosure: An
905 Empirical Study. *Account. Rev.* 53, 94–111.

906 Stanwick, P.A., Stanwick, S.D., 1998. The Relationship between Corporate Social Performance, and
907 Organizational Size, Financial Performance, and Environmental Performance: An Empirical
908 Examination. *J. Bus. Ethics* 17, 195–204.

909 Stern, N., 2008. The economics of climate change. *Am. Econ. Rev.* 98, 1–37.

910 Tashman, P., Rivera, J., 2016. Ecological uncertainty, adaptation, and mitigation in the U.S. ski resort
911 industry: Managing resource dependence and institutional pressures. *Strateg. Manag. J.* 37, 1507–
912 1525. <https://doi.org/10.1002/smj.2384>

913 Torlay, L., Perrone-Bertolotti, M., Thomas, E., Baciù, M., 2017. Machine learning–XGBoost analysis
914 of language networks to classify patients with epilepsy. *Brain Inform.* 4, 159–169.
915 <https://doi.org/10.1007/s40708-017-0065-7>

916 Trumpp, C., Endrikat, J., Zopf, C., Guenther, E., 2015. Definition, Conceptualization, and
917 Measurement of Corporate Environmental Performance: A Critical Examination of a
918 Multidimensional Construct. *J. Bus. Ethics* 126, 185–204. <https://doi.org/10.1007/s10551-013-1931-8>

919 Trumpp Christoph, Guenther Thomas, 2017. Too Little or too much? Exploring U-shaped
920 Relationships between Corporate Environmental Performance and Corporate Financial Performance.
921 *Bus. Strategy Environ.* 26, 49–68. <https://doi.org/10.1002/bse.1900>

922 Tsanas, A., Xifara, A., 2012. Accurate quantitative estimation of energy performance of residential
923 buildings using statistical machine learning tools. *Energy Build.* 49, 560–567.
924 <https://doi.org/10.1016/j.enbuild.2012.03.003>

925 Tso, G.K.F., Yau, K.K.W., 2007. Predicting electricity energy consumption: A comparison of
926 regression analysis, decision tree and neural networks. *Energy* 32, 1761–1768.
927 <https://doi.org/10.1016/j.energy.2006.11.010>

928 Tyteca, D., Carlens, J., Berkhout, F., Hertin, J., Wehrmeyer, W., Wagner, M., 2002. Corporate
929 environmental performance evaluation: evidence from the MEPI project. *Bus. Strategy Environ.* 11,
930 1–13. <https://doi.org/10.1002/bse.312>

931 Udayasankar, K., 2008. Corporate Social Responsibility and Firm Size. *J. Bus. Ethics* 83, 167–175.
932 <https://doi.org/10.1007/s10551-007-9609-8>

933 Veleva, V., Ellenbecker, M., 2001. Indicators of sustainable production: framework and methodology.
934 *J. Clean. Prod.* 9, 519–549. [https://doi.org/10.1016/S0959-6526\(01\)00010-5](https://doi.org/10.1016/S0959-6526(01)00010-5)

935 Walls, J.L., Berrone, P., Phan, P.H., 2012. Corporate governance and environmental performance: is
936 there really a link? *Strateg. Manag. J.* 33, 885–913. <https://doi.org/10.1002/smj.1952>

937 Wang, H., Bi, J., Wheeler, D., Wang, J., Cao, D., Lu, G., Wang, Y., 2004. Environmental
938 performance rating and disclosure: China’s GreenWatch program. *J. Environ. Manage.* 71, 123–133.
939 <https://doi.org/10.1016/j.jenvman.2004.01.007>

940 Wang, Q., Zhao, Z., Zhou, P., Zhou, D., 2013. Energy efficiency and production technology
941 heterogeneity in China: A meta-frontier DEA approach. *Econ. Model.* 35, 283–289.
942 <https://doi.org/10.1016/j.econmod.2013.07.017>

943 Wang, S., Li, J., Zhao, D., 2018. Institutional Pressures and Environmental Management Practices:
944 The Moderating Effects of Environmental Commitment and Resource Availability. *Bus. Strategy
945 Environ.* 27, 52–69. <https://doi.org/10.1002/bse.1983>

946 Wiedenhofer, D., Guan, D., Liu, Z., Meng, J., Zhang, N., Wei, Y.-M., 2017. Unequal household
947 carbon footprints in China. *Nat. Clim. Change* 7, 75–80. <https://doi.org/10.1038/nclimate3165>

948 Xu, P., Chan, E.H.W., 2013. ANP model for sustainable Building Energy Efficiency Retrofit (BEER)
949 using Energy Performance Contracting (EPC) for hotel buildings in China. *Habitat Int., Low-Carbon
950 Cities and Institutional Response* 37, 104–112. <https://doi.org/10.1016/j.habitatint.2011.12.004>

951 Xu, P., Chan, E.H.W., Visscher, H.J., Zhang, X., Wu, Z., 2015. Sustainable building energy efficiency
952 retrofit for hotel buildings using EPC mechanism in China: analytic Network Process (ANP)
953 approach. J. Clean. Prod. 107, 378–388. <https://doi.org/10.1016/j.jclepro.2014.12.101>

954 Yeeles, A., 2018. Business as usual. Nat. Clim. Change 8, 13. <https://doi.org/10.1038/s41558-017-0051-3>

956 Yoon, Y., Gürhan-Canli, Z., Schwarz, N., 2006. The Effect of Corporate Social Responsibility (CSR)
957 Activities on Companies With Bad Reputations. J. Consum. Psychol. 16, 377–390.
958 https://doi.org/10.1207/s15327663jcp1604_9

959 Zhang, B., Bi, J., Yuan, Z., Ge, J., Liu, B., Bu, M., 2008. Why do firms engage in environmental
960 management? An empirical study in China. J. Clean. Prod. 16, 1036–1045.
961 <https://doi.org/10.1016/j.jclepro.2007.06.016>

962 Zhang, K., Wen, Z., Peng, L., 2007. Environmental Policies in China: Evolvement, Features and
963 Evaluation. China Popul. Resour. Environ. 17, 1–7. [https://doi.org/10.1016/S1872-583X\(07\)60006-0](https://doi.org/10.1016/S1872-583X(07)60006-0)

964 Zheng, S., Alvarado, V., Xu, P., Leu, S.-Y., Hsu, S.-C., 2018. Exploring spatial patterns of carbon
965 dioxide emission abatement via energy service companies in China. Resour. Conserv. Recycl. 137,
966 145–155. <https://doi.org/10.1016/j.resconrec.2018.06.004>

967

968 **Appendix-XGBoost algorithm**

969

```

970 from XGBoost import plot_tree
971 import matplotlib.pyplot as plt
972 import numpy as np
973 import pandas as pd
974 from pandas import read_csv, read_excel
975 import XGBoost as xgb
976 from sklearn.model_selection import train_test_split
977 from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error
978 from sklearn.ensemble import RandomForestRegressor
979 from sklearn.preprocessing import Imputer, StandardScaler
980 from statsmodels.stats.outliers_influence import variance_inflation_factor
981 from sklearn.base import BaseEstimator, TransformerMixin
982 import matplotlib.pyplot as plt
983 from sklearn import svm
984 import shap
985
986 data= pd.read_excel("data/111.xlsx")
987 data = data.drop(['NO.'],axis=1)
988 label = data.pop('Y2')
989
990
991
992 def main():
993     # split data into train and test sets
994     seed = 7
995     test_size = .25

```

```

996
997     X_train, X_test, y_train, y_test = train_test_split(data, label, test_size=test_size,
998 random_state=seed)
999     original_col = X_train.columns
1000     imp = Imputer(missing_values='NaN', strategy='mean', axis=0)
1001     imp.fit(X_train)
1002     X_train = imp.transform(X_train)
1003     X_test = imp.transform(X_test)
1004
1005
1006     # random forest algorithm
1007     regr_rf = RandomForestRegressor(max_depth=30, random_state=2)
1008     regr_rf.fit(X_train, y_train)
1009     y_pred_train1 = regr_rf.predict(X_train)
1010     y_pred1 = regr_rf.predict(X_test)
1011     # random forest end
1012
1013     # XGBoost algorithm
1014     xgdmatrix = xgb.DMatrix(X_train, y_train)
1015     our_params = {'eta':.03, 'seed':0, 'subsample':0.8, \
1016                  'colsample_bytree':0.8, 'objective':'reg:linear', \
1017                  'max_depth':7, 'min_child_weight':.5}
1018
1019     # train the model
1020     final_gb = xgb.train(our_params, xgdmatrix, num_boost_round=1500)
1021
1022     testmat = xgb.DMatrix(X_test)
1023     trainmat = xgb.DMatrix(X_train)
1024     y_pred2 = final_gb.predict(testmat)
1025     y_pred_train2 = final_gb.predict(trainmat)
1026     # XGBoost end
1027
1028     # svm regression
1029     clf = svm.SVR(kernel='rbf', degree = 3, gamma = 'auto', coef0=0.0, tol=0.1, C=1.0, epsilon=0.1,
1030 shrinking = True, cache_size=200, verbose=False, max_iter=-1)
1031     clf.fit(X_train, y_train)
1032     y_pred_train3 = clf.predict(X_train)
1033     y_pred3 = clf.predict(X_test)
1034     # end svm
1035
1036     #random forest
1037     mae = mean_absolute_error(y_test.values, y_pred1)
1038     print("MAE: %.5f" % mae)
1039     rmse = np.sqrt(mean_squared_error(y_test.values, y_pred1))
1040     print("RMSE: %.5f" % rmse)
1041     R = np.corrcoef(y_test.values, y_pred1)
1042
1043     print("Correlation Coef: %.5f" % R[0,1])
1044     r2 = r2_score(y_test.values, y_pred1)
1045     print("r2 score: %.5f" % r2)
1046
1047     #XGBoost
1048     mae = mean_absolute_error(y_test.values, y_pred2)
1049     print("MAE: %.5f" % mae)
1050     rmse = np.sqrt(mean_squared_error(y_test.values, y_pred2))

```

```

1051 print("RMSE: %.5f" % rmse)
1052 R = np.corrcoef(y_test.values,y_pred2)
1053
1054 print("Correlation Coef: %.5f" % R[0,1])
1055 r2 = r2_score(y_test.values,y_pred2)
1056 print("r2 score: %.5f" % r2)
1057
1058 #svm
1059 mae = mean_absolute_error(y_test.values, y_pred3)
1060 print("MAE: %.5f" % mae)
1061 rmse = np.sqrt(mean_squared_error(y_test.values, y_pred3))
1062 print("RMSE: %.5f" % rmse)
1063 R = np.corrcoef(y_test.values,y_pred3)
1064
1065 print("Correlation Coef: %.5f" % R[0,1])
1066 r2 = r2_score(y_test.values,y_pred3)
1067 print("r2 score: %.5f" % r2)
1068
1069 # #plot predict error
1070 plt.gcf().set_size_inches((10, 4))
1071
1072 plt.plot(((y_pred1-y_test.values)/y_test.values)[:8], color='g', marker='*', label='random forest')
1073 plt.plot(((y_pred2-y_test.values)/y_test.values)[:8], color='c', marker='s', markerfacecolor='none',
1074 label='XGBoost')
1075 plt.plot(((y_pred3-y_test.values)/y_test.values)[:8], color='y', marker='o',
1076 markerfacecolor='none', label='SVM')
1077 # plt.gca().legend()
1078 plt.legend(loc='upper right')
1079 plt.savefig('junk.jpg')
1080
1081 # plot training error
1082 plt.gcf().set_size_inches((10, 4))
1083 plt.plot(((y_pred_train1-y_train.values)/y_train.values)[:20], color='g', marker='*',
1084 label='random forest')
1085 plt.plot(((y_pred_train2-y_train.values)/y_train.values)[:20], color='c', marker='s',
1086 markerfacecolor='none', label='XGBoost')
1087 plt.plot(((y_pred_train3-y_train.values)/y_train.values)[:20],color='y', marker='o',
1088 markerfacecolor='none', label='SVM')
1089 # plt.gca().legend()
1090 plt.legend(loc='upper right')
1091 plt.savefig('junk.jpg')
1092
1093
1094 # plot predict test
1095 plt.gcf().set_size_inches((10, 4))
1096 plt.plot(y_test.values[:3], color='b', label='value')
1097 plt.plot(y_pred1[:3], color='g', marker='*', markerfacecolor='none', label='random
1098 forest',linestyle='None')
1099 plt.plot(y_pred2[:3], color='c', marker='s', markerfacecolor='none',
1100 label='XGBoost',linestyle='None')
1101 plt.plot(y_pred3[:3], color='y', marker='o', markerfacecolor='none',
1102 label='SVM',linestyle='None')
1103 # plt.gca().legend()
1104 plt.legend(loc='upper right')
1105 plt.savefig('junk.jpg')

```

```

1106
1107 #plot training data
1108     plt.gcf().set_size_inches((10, 4))
1109     plt.plot(y_train.values[::10], color='b', label='value')
1110     plt.plot(y_pred_train1[::10], color='g', marker='*', markerfacecolor='none', label='random
1111 forest',linestyle='None')
1112     plt.plot(y_pred_train2[::10], color='c', marker='s', markerfacecolor='none',
1113 label='XGBoost',linestyle='None')
1114     plt.plot(y_pred_train3[::10], color='y', marker='o', markerfacecolor='none',
1115 label='SVM',linestyle='None')
1116     # plt.gca().legend()
1117     plt.legend(loc='upper right')
1118     plt.savefig('junk2.jpg')
1119
1120 # shap value
1121     shap.initjs()
1122     shap_values = shap.TreeExplainer(final_gb).shap_values(X_train)
1123     X_train = pd.DataFrame(data=X_train, columns=original_col)
1124     X_train = X_train.rename(columns={
1125         "X2": "X7", "X3":
1126 "X6", "X4": "X14", "X5": "X4", "X6": "X8", "X7": "X9", "X8": "X10", "X9": "X12",
1127         "X10": "X11", "X11": "X13", "X12": "X5", "X13": "X1", "X14": "X2", "X15": "X3"})
1128     shap.summary_plot(shap_values, X_train)
1129
1130 main()

```