

A Bayesian LSTM Model to Evaluate the Effects of Air Pollution Control Regulations in Beijing, China

Yang Han^a, Jacqueline CK Lam^{a, b*}, Victor OK Li^{a, b}, and David Reiner^b

^aDepartment of Electrical and Electronic Engineering, The University of Hong Kong, Pokfulam, Hong Kong

^bEnergy Policy Research Group, Judge Business School, The University of Cambridge, Trumpington Street, Cambridge, United Kingdom

Acknowledgement

This research is supported in part by the General Research Fund of the Research Grants Council of Hong Kong, under Grant No. 17620920. We gratefully acknowledge the PM_{2.5} data provided by Beijing US Embassy, China; and AOD data provided by NASA, USA. We also gratefully acknowledge the valuable comments and suggestions of our reviewers, and participants who attended our first HKU-Cambridge AI-WiSe workshop, titled “AI for Social Good”, held in the Department of Computer Science and Technology, Cambridge, the United Kingdom, in May 2019, during which an earlier version of this work was presented.

Declarations of interest: none

© 2020. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>

* Corresponding author at: Department of Electrical and Electronic Engineering, The University of Hong Kong, Pokfulam, Hong Kong. Tel.: (852) 3917 4843; E-mail address: jcklam@eee.hku.hk

A Bayesian LSTM Model to Evaluate the Effects of Air Pollution Control Regulations in Beijing, China

Abstract

Rapid socio-economic development and urbanization have resulted in serious deterioration in air-quality in many world cities, including Beijing, China. This study attempts to examine the effectiveness of air pollution control regulations implemented in Beijing during 2008 – 2019 through a data-driven regulatory intervention analysis. Our proposed Bayesian deep learning model utilizes proxy data including Aerosol Optical Depth (AOD) and meteorology as well as socio-economic data, while accounting for confounding effects via propensity score estimation. Our results show that air pollution control regulatory measures implemented in China and Beijing during 2008 – 2019 reduced PM_{2.5} pollution in Beijing by 11% on average. After the introduction of Action Plan for Clean Air in China and Beijing in late 2013, as compared to the hypothetical PM_{2.5} concentration (without any regulatory interventions), the estimated PM_{2.5} reduction increased dramatically from 15% in 2015 to 44% in 2018. Our results suggest that Beijing's air quality has improved gradually over the past decade, though the annual PM_{2.5} pollution still exceeds the WHO threshold. In this regard, the air pollution control regulations introduced in Beijing and China tend to become more effective after 2015, suggesting a 2-year time lag before the stringent air pollution control regulations starting from 2013 takes any strong positive effects. Moreover, as compared to the air pollution control regulations introduced before 2013, newly introduced policy-making governance, which couples the policy-makings of the local jurisdictions with that of the central government, and the new policy measures that tackle the vested interests of the local stakeholders in Beijing and its nearby cities, alongside with the stringent local and national air pollution control regulations and plans, should help reduce air pollution and promote healthy living in Beijing over the longer term.

Keywords: air pollution control regulations, effects of regulatory interventions, Bayesian LSTM, propensity score, counterfactual analysis, causal inference

Highlights

- Aerosol optical depth, meteorology, and socio-economic data are collected
- A Bayesian deep-learning approach is proposed for regulatory intervention analysis

- 34 • Confounding effects are addressed by the propensity score estimation
- 35 • Air pollution controls reduced PM_{2.5} in Beijing by 11% during 2008 – 2019

36

37 **1. Introduction**

38 Over the past few decades, rapid socio-economic development and urbanization have resulted
39 in serious deterioration in air quality in Beijing, China. Air pollutants, especially PM_{2.5}
40 (particulates smaller than 2.5 micrometers in diameter), can lead to extremely detrimental
41 health consequences, such as cancer, stroke, asthma, or heart disease (Pope III and Dockery,
42 2006; Pui et al., 2014). To provide in a timely manner, the critical health advice for Beijing’s
43 citizens based on scientific evidence, the introduction of real-time air pollution monitoring and
44 reporting system in China has become increasingly crucial. Since April 2008, the US Embassy
45 in Beijing has been publishing hourly PM_{2.5} readings based on its own monitors installed in the
46 embassy building. In January 2013, Beijing officially launched a new air quality monitoring
47 system. Since then, PM_{2.5} has been fully monitored by Beijing’s automatic monitoring network,
48 with hourly air pollution concentrations released by Beijing’s Environmental Monitoring
49 Center. A number of air pollution control regulations have been introduced by the government
50 in China to control air pollution, with increasing stringency over the last two decades. Using
51 Beijing as a case study, this study proposes a data-driven regulatory intervention analysis
52 framework to study the causal relationship between air pollution control regulations and city-
53 level PM_{2.5} pollution concentration, based on available monitored air pollution data, proxy data
54 including AOD and meteorology, and socio-economic data. The effects of air pollution control
55 regulations in Beijing during 2008 – 2019 are evaluated. Our current work is an extension of
56 an earlier work, which evaluates the effectiveness of air pollution control regulations in Beijing,
57 China, during 2013 – 2017 (Han et al., 2018), by (1) adding socio-economic statistics in the
58 input data, (2) taking account of the effective periods of air pollution control regulations, (3)
59 extending the period of study to 2008 – 2019, and (4) reducing the confounding biases via
60 propensity score estimation. The rest of the paper is organized as follows. Section 2 reviews
61 related works. Section 3 discusses our collected data and proposed the methods for regulatory
62 intervention analysis. Section 4 presents our experimental results, followed by discussions on
63 limitations of study and future directions. Section 5 highlights the policy implications. Section
64 6 concludes our study.

65

66 **2. Related Work**

67 **2.1 Machine-Learning for Causal Inference and Policy Evaluation**

68 Examining the causal effects from observational datasets has been a subject of serious attention
69 in the fields of social science, policy, or medical science (Athey, 2017; Imbens and Rubin,
70 2015). Numerous theories attempted to account for the cause of an outcome/event, with the
71 counter-factual framework being widely recognized and adopted (Rubin, 2005). Under such
72 framework, the causal relationship between X and Y can be re-formulated as a counter-factual
73 question. For example, in order to test whether there is a causal relationship between X and Y ,
74 a question such as “If X had not occurred, what would Y be?” is raised. However, it remains a
75 difficult challenge to determine the counter-factual outcome. For any unit at a given time point,
76 only the factual (instead of the counter-factual) outcome of a specific intervention could be
77 observed (Rubin, 2005). Confounding factors were considered as the major barrier for
78 determining with confidence whether causal relationships exist across a set of examined
79 variables in a big dataset (Pearl, 2018). A proper and rigorous solution is to resort to the
80 randomized control trials (RCT) and perform statistical adjustments to reduce the confounding
81 biases. RCT is the “gold standard” for evaluating the causal effects while controlling for any
82 confounding variables. However, in many cases, it was an infeasible task due to high research
83 cost and ethical constraints (Stolberg et al., 2004). Hence, traditional statistical techniques,
84 such as matching, re-weighting, and propensity score, were proposed to reduce confounding
85 bias in observational studies (Athey and Imbens, 2017). However, these traditional methods
86 were usually based on low-dimensional linear modelling, and failed to capture the complex
87 non-linear relationships identified from high dimensional datasets (Hartford et al., 2017).

88 Advances in deep-learning have given rise to the remarkable success in overcoming
89 many computational challenges that involve non-linear modelling of high dimensional data,
90 such as natural language processing and computer vision (LeCun et al., 2015). However, these
91 deep-learning models were mostly trained on datasets that carried noisy and unrepresentative
92 big data (Caliskan et al., 2017), and often failed to account for the confounding effects when
93 making causal inference (Marcus, 2018). As a result, spurious causations might occur and
94 biased decisions made (Osoba and Welser IV, 2017). For supervised machine-learning
95 algorithms, a fundamental shift from correlation analysis to causality analysis is needed to fully
96 understand the causal relationship between an intervention (treatment or policy change) and an
97 outcome. Recently, there has been a growing interest in using deep-learning models for causal
98 inference and policy evaluation, based on techniques such as autoencoder (Atan et al., 2018)
99 or variational autoencoder (VAE) (Louizos et al., 2017), propensity dropout (Alaa et al., 2017),
100 propensity score estimation (Shi et al., 2019), domain adaptation (Shalit et al., 2017), multi-

101 task learning (Alaa and van der Schaar, 2017), and generative adversarial network (GAN)
102 (Yoon et al., 2018). These techniques aimed to improve the generalization ability of the models
103 beyond observational data, and to reduce the confounding biases in high dimensional data.
104 Moreover, since it is difficult to take into account all important confounders in counterfactual
105 modelling, some other techniques, including the instrument variable (IV) method, were applied
106 to the deep-learning models to control for any unobserved confounders, with additional
107 assumptions taken (Hartford et al., 2017).

108 **2.2 Evaluation of Air Pollution Regulatory Interventions**

109 Many studies examined the effect of regulatory interventions on pollution concentrations in
110 both the Chinese and the international context. Two major approaches, namely, (1) the
111 environmental engineering approach and (2) the environmental economic approach, were
112 adopted in these studies (Li et al., 2017d). The first approach provided an *ex ante* evaluation of
113 policy impacts, by forecasting air qualities under different policy scenarios or constructing
114 hypothetical air qualities in the absence of policy regulations, using physical and statistical
115 modelling (Liu et al., 2012). The second approach performed an *ex post* evaluation of the causal
116 effects of policy interventions, using experimental/quasi-experimental design and
117 observational data, and methods such as difference-in-differences estimation (Chen et al.,
118 2013), regression discontinuity design (RDD) (Li et al., 2017d), and panel data regression
119 (Zheng et al., 2015). However, both approaches had drawbacks. The first one was often
120 constrained by high computational costs, complex process modelling, and high uncertainties in
121 emission inventories (Li et al., 2017d; Liu et al., 2010). The second one often failed to model
122 the complex relationship between air pollution and other confounders such as meteorology and
123 time trends, account for the uncertainties in input data and model parameters, and establish the
124 causal relationship only after controlling for the confounders (Ferraro, 2009; Henneman et al.,
125 2017).

126 Rapid development in machine learning made the adoption of data-driven regulatory
127 analysis possible, with applications in resource allocation and causal inference (Athey, 2017).
128 Recently, deep-learning approaches achieved state-of-the-art performance in air pollution
129 estimation and forecasting (Li et al., 2017b; Li et al., 2017c; Ong et al., 2016), including PM_{2.5}
130 estimation, utilizing satellite-based Aerosol Optical Depth (AOD) as proxy data (Li et al.,
131 2017a). However, in studies such as Li et al. 2017a, the temporal correlation between PM_{2.5}
132 pollution concentration and AOD is yet to be fully exploited by the neural network structure.
133 Moreover, deep learning can still suffer from limited data source and low data quality when
134 compared to other machine-learning techniques. Incorporating the Bayesian approach into deep

135 learning can reduce network overfitting due to data sparsity and noise, and provide uncertainty
136 measure for the prediction (Gal, 2016). However, a data-driven approach is yet to be applied
137 to accurately estimate the counter-factual effects of air pollution regulatory interventions on air
138 pollution outcomes, while accounting for the confounding biases.

139

140 **3. Data and Method**

141 This study proposes a machine-learning framework to provide counter-factual inference and
142 evaluate the effects of air pollution regulatory interventions. The problem setup is similar to
143 previous studies where potential outcome frameworks were adopted for causal inference (Alaa
144 and van der Schaar, 2017; Atan et al., 2018). Differing from previous studies, our work focusses
145 more specifically on evaluating the aggregate effect of multiple air pollution regulatory
146 interventions. For each daily observation, there is a corresponding regulatory intervention state,
147 which falls into two potential outcomes: the first potential outcome is a regulatory state where
148 all regulatory interventions have been implemented as planned, whilst the second potential
149 outcome is a regulatory state where no regulatory intervention has been implemented. Our goal
150 is to learn how each feature-intervention pair is mapped to its corresponding factual outcomes,
151 based on the observational air pollution samples collected during the period of study. Once the
152 mapping model is trained, given an observed sample of air pollution outcomes after a group of
153 regulatory interventions has been implemented, the counter-factual outcomes can be estimated
154 for the scenario when the equivalent regulatory interventions are not implemented. Moreover,
155 to estimate the causal effects of regulatory interventions, we follow the un-confoundedness
156 assumption made in the potential outcome framework (Wooldridge, 2000). We assume that all
157 important confounders that can potentially affect the regulatory interventions and the air quality
158 outcomes have been taken into account in our model, and the confounding biases can be
159 addressed via the propensity score estimation. Our proposed Bayesian deep learning policy
160 intervention framework consists of four components, covering, data collection, data pre-
161 processing, model training, and regulatory intervention analysis (see Figure 1).

162

163 [Insert Figure 1 about here]

164

165 **3.1 Data Collection**

166 We collected data consisting of air quality, AOD, meteorology, socio-economic, and air
167 pollution regulatory measures from 2008 to 2019 (see Table 1 for a summary of data sources).

168

[Insert Table 1 about here]

169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201

3.1.1 Air Quality Data

We collected hourly PM_{2.5} concentration data recorded at the US Embassy, Beijing from 9 April 2008 to 31 December 2019 (US Department of State, 2020). Only PM_{2.5} observations with the quality control label “Valid” were included. Hourly PM_{2.5} data were aggregated to daily means. Given that official PM_{2.5} concentration data had not been available until 2013, air pollution observations at the US Embassy, Beijing for city-level PM_{2.5} pollution concentrations in Beijing during the study period were used as the ground truths. Existing studies showed that Beijing’s city-level PM_{2.5} concentrations was highly correlated with PM_{2.5} concentrations observed at the US Embassy, Beijing. Hence, it was reasonable to assume that the readings reported by the US Embassy can be used to represent the level of air-quality throughout the city in Beijing (Wang et al., 2013). To further examine the representativeness of the US Embassy PM_{2.5} data, we collected official hourly station-level PM_{2.5} concentration data from 1 January 2014 to 31 December 2019 using the data source provided in Zhang et al. (2019), and examined the correlation between the daily average PM_{2.5} concentrations measured at the US Embassy, Beijing and the daily city-level average PM_{2.5} concentrations measured at the 35 official stations in Beijing during 2014 – 2019. Result showed that the two measurements were highly correlated ($R^2=96.2\%$; see Figure 2).

[Insert Figure 2 about here]

3.1.2 Proxy Data (AOD and Meteorology)

Previous studies showed that AOD and meteorology data can be incorporated into the statistical modelling to examine the effects of regulatory interventions on air pollution concentrations in Beijing (Liu et al., 2012). Our study had incorporated the AOD data into our statistical modelling. AOD observations at the city level were collected from the NASA MODIS satellite database from 26 March 2008 to 21 March 2019 (US NASA, 2020). Five features were selected based on data availability during the period of study, including AOD at 1020 nm, AOD at 870 nm, AOD at 675 nm, AOD at 440 nm, and precipitable water. AOD data points observed each day were aggregated into daily means. In addition, hourly city-level meteorology data, including temperature, relative humidity, wind speed, wind bearing, and visibility, across the period from 1 January 2008 to 31 December 2019, were collected from a weather data

202 application program interface (API), based on the official data sources (Apple Inc., 2020).
203 Hourly meteorology data were aggregated to daily means.

204 **3.1.3 Socio-economic Data**

205 Previous studies showed that socio-economic data can be used as control variables to model
206 statistically the effects of regulatory interventions on air pollution concentrations at the
207 provincial-level in China (Zheng et al., 2015). In this study, we collected the yearly socio-
208 economic statistics including the percentage of GDP generated from the secondary sector, the
209 population density, and the number of vehicles, during the period of 2008 to 2019 (Beijing
210 Municipal Bureau of Statistics, 2020; Beijing Transport Institute, 2020).

211 **3.1.4 Regulatory Measures Data**

212 We identified major air pollution control regulations at the city- or the national-level during
213 the period of 2008 to 2019 (DieselNet, n.d.; Lam et al., 2019; Zhang et al., 2016). These
214 regulations were directly responsible for air pollution prevention and control in Beijing/China,
215 with a strong focus on the energy and transportation sectors, including emission controls on
216 the coal-fired power plants and the industrial facilities and vehicles, emission standards on cars
217 and light trucks, optimization of energy structures and traffic systems, technological
218 innovations of clean environment, emergency plans for high pollution episodes, and legal
219 responsibilities. Some were updated during the period of study, including, the Air Pollution
220 Prevention and Control Law in China (see Figure 3).

221

222 [Insert Figure 3 about here]

223

224 **3.2 Data Preprocessing**

225 The daily air quality data and the daily proxy data were combined to generate a tabular dataset,
226 ranging from 9 April 2008 to 21 March 2019. The dataset was pre-processed for model training,
227 validation, and evaluation. The data pre-processing procedure was listed as follows. First, a
228 random 80/10/10 split of the data was used as the training set, the validation set, and the test
229 set. Second, each dataset was converted into the input/output pairs. The input data consisted of
230 two parts: a vector representing the historical daily proxy data (including AOD and
231 meteorology) of the current day and the previous days, and a binary vector representing the
232 current status of regulatory interventions. To account for the socio-economic variation, the
233 corresponding yearly statistics were included in the input vector. To account for the unobserved
234 time trends and recurrent effects, the month and the day of the week were included as the
235 categorical features in the input vector. The output was a continuous value representing the

236 corresponding city-level daily PM_{2.5} concentration (i.e., PM_{2.5} concentration observed at the US
 237 Embassy, Beijing). Then, missing data was filled in via iterative imputation (Buuren and
 238 Groothuis-Oudshoorn, 2010). To avoid any potential information leakages, the iterative
 239 imputer was constructed based on the training set, which was subsequently used to impute the
 240 validation and test datasets. Next, each input feature (except for the time trend) in the training
 241 set was standardized according to its mean and standard deviation; which were then used to
 242 standardize the corresponding feature of the validation and the test datasets. Finally, the data
 243 pre-processing procedure was repeated five times. Eventually, five datasets for model training,
 244 validation, and test were constructed from five random data splits.

245 **3.3 Model Training**

246 The pre-processed data was fed into a Bayesian deep learning model for training. During the
 247 period of study, the covariate data at day t was denoted as x_t . The input data for day t consisted
 248 of the observations over the past $L + 1$ days (including the current day t) and the time trend:
 249 $X_t = \{x_{t-L}, \dots, x_t, \text{Month}_t, \text{Day of week}_t\}$. The regulatory status vector at day t consisted of
 250 the status of K regulatory interventions $I_t = \{I_t^1, \dots, I_t^K\}$, e.g., {Regulation 1 is implemented,
 251 Regulation 2 is not implemented, ..., Regulation K is not implemented} (see Figure 3 for the
 252 effective periods). We used zero or one to indicate the status of a particular regulatory
 253 intervention I_t^k , namely, one for “is implemented” and zero for “is not implemented”. The
 254 output y_t was the observed city-level air quality (i.e. PM_{2.5} concentrations observed at the US
 255 Embassy, Beijing). The proposed framework had two potential outputs, the first one
 256 corresponded to I_t , where all regulatory interventions are implemented as planned, while the
 257 second one corresponded to a regulatory state where no regulatory interventions are
 258 implemented. A Bayesian deep-learning model with network structure f and parameters θ was
 259 denoted as f_θ . During the study period of length T , given the input X_t , the regulatory
 260 intervention status I_t , and the output $y_t = f(X_t, I_t)$, we aimed to estimate the counter-factual
 261 output $\tilde{y}_t = f(X_t, \mathbf{0})$. The model f_θ aimed to find the optimal posterior distribution of the
 262 network weight parameters θ , given the observed tuples $\{(X_t, I_t, y_t)\}_{t=1}^T$. To better address the
 263 confounding effects, a shared representation layer was used (1) to predict air quality based on
 264 the covariate and the regulatory intervention status and (2) to predict regulatory intervention
 265 status from the covariates (i.e., the propensity score estimation). By incorporating the
 266 propensity score estimation model into the proposed framework, the input features relevant for
 267 confounding effects could be distilled automatically (Shi et al., 2019). More specifically, we
 268 focussed on the Bayesian RNN, which is a particular type of Bayesian deep learning model

269 capable of modelling time-series data (Fortunato et al., 2017). We used LSTM as the recurrent
 270 unit of the network. A Bayesian embedding layer was used to map the time trend vector into a
 271 vector of continuous values (Yi et al., 2018). Another Bayesian embedding layer was used to
 272 map the regulatory status vector to a vector of continuous values (Pham et al., 2017). Two
 273 Bayesian fully connected linear layers were utilized. One Bayesian fully connected linear layer
 274 was used to predict y_t , while the other Bayesian fully connected linear layer followed by a
 275 sigmoid function, was used to predict I_t (Shi et al., 2019). Both of them were based on the
 276 shared representation, which consisted of three parts, the final hidden state of Bayesian LSTM
 277 (h_t), the embedded time trend (e_t^1), and the embedded regulatory intervention status (e_t^2).
 278 Conceptually, our proposed model was as follows:

279

$$280 \quad h_t = \text{Bayesian-LSTM}(x_t, h_{t-1}) \quad (1)$$

$$281 \quad e_t^1 = \text{Bayesian-Embedding}(\text{Month}_t, \text{Day of week}_t) \quad (2)$$

$$282 \quad e_t^2 = \text{Bayesian-Embedding}(I_t) \quad (3)$$

$$283 \quad y_t = \text{Bayesian-Linear}(h_t, e_t^1, e_t^2) \quad (4)$$

$$284 \quad I_t = \text{Sigmoid}(\text{Bayesian-Linear-Propensity-Score}(h_t, e_t^1, e_t^2)) \quad (5)$$

285

286 To train our proposed model, we followed the work done by Blundell et al. (2015); Fortunato
 287 et al. (2017). In the network, each weight parameter was a random variable with a Gaussian
 288 mixture prior, and the weight at each time step had the same distribution. A diagonal Gaussian
 289 distribution was used as the variational posterior distribution, which is often computationally
 290 tractable and numerically stable, assuming that the network weights were uncorrelated. The
 291 loss function of the proposed model consisted of three components. The first part was the Mean
 292 Squared Error (MSE) loss calculated by the predicted and the observed air quality values,
 293 which is the most commonly used loss for predicting continuous values. The second part was
 294 the Binary Cross Entropy (BCE) loss calculated by the predicted and the observed regulatory
 295 intervention status, which is often used for multi-label classification. The BCE loss enforced
 296 the learned shared representation layer to account for the propensity score estimation, in order
 297 to address the confounding effects. The third part was the Kullback–Leibler (KL) divergence
 298 between the posterior and the prior distribution, which is a regularization term to penalize
 299 model overfitting. Bayes by Backprop was adopted to update the weight parameters of the
 300 network while minimizing the loss function, given the observed inputs (see Algorithm 1). The

301 proposed model was trained via the shuffled mini batches, using a stochastic gradient descent
302 (SGD) optimizer.

303

304 **Algorithm 1.** Bayesian LSTM Model Training via Bayes by Backprop

305 **Require:** training data $D = \{(X_t, I_t, y_t)\}_{t=1}^{t=T}$, epoch size E , batch size B , and learning rate α

306 **For epoch from 1 to E**

307 **Repeat**

308 1. Sample a mini batch of size B from the training data D without replacement

309 2. Sample $\varepsilon \sim \text{Gaussian}(0, I)$, where I is the identity matrix

310 3. Set network parameters $\theta = \mu + \sigma\varepsilon$, where μ and σ are the mean and
311 standard deviation, respectively

312 4. Compute the gradients of MSE loss plus BCE loss

313 with respect to θ using normal back-propagation: g_θ^L

314 5. Compute the gradients of $F(\mu, \sigma, \theta) = \log \text{Gaussian}(\mu, \sigma^2) - \log p(\theta)$ with
315 respect to μ, σ, θ : $g_\mu^F, g_\sigma^F, g_\theta^F$, where $p(\theta)$ is the Gaussian mixture prior

316 6. Update $\mu = \mu - \alpha \frac{g_\theta^L + g_\theta^F + g_\mu^F}{B}$

317 7. Update $\sigma = \sigma - \alpha \frac{g_\theta^L \varepsilon + g_\theta^F \varepsilon + g_\sigma^F}{B}$

318 **Until all mini-batches are sampled**

319 **End**

320 **Return** fitted network model f_θ

321

322 During the model training, the tuning hyper-parameters took into account the number of lagged
323 observations (0 or 7; 0 indicated that no lagged observations were used, while 7 indicated that
324 the past one week data was used for prediction), the embedding dimension of the regulatory
325 intervention status vector (3 or 5), the number of hidden units used in the neural network (128
326 or 256), the batch size (32 or 64). For each data split (including the training, the validation, and
327 the test dataset), the best hyper-parameters were selected based on the validation MSE.
328 Moreover, the fixed hyper-parameters included the number of training epochs (30), the learning
329 rate (0.01), the number of recurrent layers (1), the embedding dimension of the time trend
330 vector (3; based on a configuration used by Yi et al. (2018)), the prior distribution of the
331 Bayesian deep-learning model ($\pi = 0.25$, $-\log\sigma_1 = 0$, and $-\log\sigma_2 = 6$; based on a
332 configuration used by Blundell et al. (2015)).

333 3.4 Regulatory Intervention Analysis

334 After the model training, counter-factual outcomes, in the absence of regulatory interventions,
335 were predicted to quantify the net effects of regulatory intervention, based on the fitted model
336 f_θ . More specifically, for each data split j , the regulatory intervention analysis was performed
337 according to the following steps. First, a random sample was drawn from the posterior of the
338 network weight parameters to obtain a model $f_{\theta_{i,j}}$. Next, the corresponding regulatory status
339 vector was constructed with the hypothesis that no regulatory intervention was implemented
340 and represented by a vector of zeros. Such hypothetical regulatory intervention status vector,
341 after combining with the covariate data X_t , were used to re-estimate PM_{2.5} concentration using
342 model $f_{\theta_{i,j}}$. This was repeated N times, such that the mean of PM_{2.5} re-estimations could be
343 calculated to account for the uncertainties of the model parameters (Kendall and Gal, 2017).
344 During the study period of length T (2008 – 2019 or a particular year such as 2017), the final
345 estimation of PM_{2.5} concentrations on day t and the average regulatory effect (ARE) were
346 calculated by the following equations:

347

$$348 \quad \tilde{y}_t^{i,j} = E[f_{\theta_{i,j}}(X_t, \mathbf{0})] \quad (6)$$

$$349 \quad \text{ARE}_{i,j} = \tilde{y} - y = E_{t \in T}[\tilde{y}_t^{i,j}] - E_{t \in T}[y_t] \quad (7)$$

350

351 where $\theta_{i,j}$ was the i th sample from the network weights posterior trained on the j th data split,
352 i ranged from 1 to N , j ranged from 1 to 5, T was in the *ex post* evaluation period, y_t and \tilde{y}_t
353 were the observed and counter-factual air quality, respectively, and X_t is the covariate data.
354 The number of posterior samples N was set to 100, in order to obtain a reasonable estimation
355 of the re-estimated air quality values. Note that the regulatory intervention analysis was
356 performed for five times, using the models trained across different data splits. Finally, given
357 that the air quality values may not follow a Gaussian distribution (see Figure 2), the final
358 estimation of ARE with 95% confidence interval (CI) was calculated based on bootstrapping.
359 More specifically, a list of ARE values of length $5 * N$ was resampled from
360 $\{\text{ARE}_{1,1}, \text{ARE}_{2,1}, \dots, \text{ARE}_{N,5}\}$ with replacement, the resampled mean was subsequently
361 calculated. This was repeated 10,000 times, and the 250 percentiles and the 9,750 percentiles
362 of the resampled means were selected as the lower and the upper bound of the ARE during the
363 study period, respectively.

364

365 4. Results

366 4.1 Baseline Selection and Model Evaluation

367 Previous research suggested that non-linear relationship might exist between PM_{2.5} pollution
368 concentration and other covariates data (Han et al., 2018). Hence, in our experiment, two non-
369 linear machine-learning models, namely, Support Vector Regression (SVR) and Random
370 Forest (RF), were selected as the baseline models. We used Mean Absolute Error (MAE) and
371 Mean Absolute Percentage Error (MAPE) for model evaluation and comparison. For Bayesian
372 LSTM, we fine-tuned the hyper-parameters as listed in Section 3.3. For SVR, we fine-tuned
373 three hyper-parameters, including the lagged observations (0 or 7), the kernel function
374 (polynomial function or radial basis function) and the penalty parameter of the error term (0.1,
375 1, or 10). For RF, we fine-tuned four hyper-parameters, including the lagged observations (0
376 or 7), the number of estimators (10 or 100), the maximum depth of the tree (1, 16, or 32), and
377 the maximum number of features (n , \sqrt{n} , or $\log_2(n)$, where n is the number of features).
378 Finally, the models with the lowest MSE on the validation set were selected as the final models
379 for further analysis.

380 The performance of the proposed model and the baseline models are shown in Table 2.
381 Results have clearly revealed that the Bayesian LSTM model outperforms the baseline models.
382 On the test set, the mean MAE of the proposed model is 20.3, while the mean MAE of the SVR
383 and RF model are 22.1 and 22.4, respectively. The mean MAPE of the proposed model is
384 36.8%, while the MAPE of the SVR and RF model are 38.8% and 46.9%, respectively.
385 Moreover, the standard deviation of the proposed model's performance is also the lowest as
386 compared to the baseline models. This suggests that our proposed model can give a much better
387 prediction of the out-of-sample data as compared to traditional machine-learning techniques,
388 across different training/validation/test data splits. Note that the absolute/relative error rates of
389 the proposed model remain high, partly due to the fact that some features inputs (which were
390 irrelevant to the causal relationships according to the propensity score estimation) were
391 considered as noise for air quality estimation. However, the causal estimation of ARE can be
392 improved through such a trade-off between predictive accuracy and propensity score estimation
393 (Shi et al., 2019).

394

395 [Insert Table 2 about here]

396

397 4.2 Regulatory Intervention Analysis

398 We used the final fitted Bayesian LSTM models across different data splits to estimate/simulate
399 the average counter-factual air quality, based on the assumption that all regulatory interventions
400 were not implemented, in order to examine the ARE of all regulatory interventions during 2008
401 – 2019. Figure 4 shows that the observed monthly average daily air quality and simulated
402 monthly average daily air quality without any regulatory interventions during 2008 – 2019.
403 The average of observed daily $PM_{2.5}$ concentration was $86 \mu g/m^3$ during 2008 – 2019. Had the
404 same set of regulatory interventions not been implemented before 2008, the hypothetical
405 average daily $PM_{2.5}$ pollution would be $97 \mu g/m^3$ (95% CI: $96 \mu g/m^3$ to $99 \mu g/m^3$). The average
406 intervention effect of all regulatory interventions was $11 \mu g/m^3$ (95% CI: $10 \mu g/m^3$ to $13 \mu g/m^3$).
407 This implies that the aggregate effect of all air pollution regulatory interventions implemented
408 during this period can lead to a 11% reduction in $PM_{2.5}$ pollution concentration on average.
409 Based on Eq. (7), the relative reduction was calculated as ARE / \tilde{y} , where \tilde{y} is the hypothetical
410 average daily $PM_{2.5}$ pollution.

411

[Insert Figure 4 about here]

413

414 Table 3 shows the observed yearly average daily air quality and simulated yearly average daily
415 air quality without any regulatory interventions during the period of study. Results have shown
416 that the estimated $PM_{2.5}$ reduction due to the implementation of the set of air pollution
417 regulatory interventions implemented during the 2008 – 2019 period on average was not as
418 significant as expected on average (11%), even when a series of stringent air pollution control
419 regulations and plans have been introduced during the period (see Figure 2). However, after
420 the introduction of Action Plan for Clean Air in China and Beijing in late 2013, the estimated
421 $PM_{2.5}$ reduction increased dramatically from 2% in 2014 to 15% in 2015. After 2015, the
422 estimated $PM_{2.5}$ reduction increased up to 44% in 2018, and dropped to 37%¹ in 2019.

423

424

[Insert Table 3 about here]

425

426 **4.3 Limitations of Study and Future Work**

427 There are some limitations in our current study. First, the interpretability of the proposed
428 Bayesian deep-learning framework for policy evaluation can be improved. Although the

¹ This only covers the data in the first quarter of 2019. We expect the average improvement would be changed (would likely be increased) when the full year data is incorporated into our model.

429 confounding variables have been addressed in the proposed model by incorporating a
430 propensity score estimation layer, it remains difficult to understand which variables have
431 contributed most to the confounding biases and the causal relationships, and when/where the
432 proposed model works better as compared to the traditional statistical methods for policy
433 evaluation (such as propensity score estimation using a logistic linear regression model). Future
434 work can focus on an interpretable machine-learning framework for policy evaluation. Second,
435 our study only examines the aggregate effect of air pollution control regulations and plans
436 during the period of study. More sophisticated analysis is needed to understand the individual
437 effect of a particular regulatory intervention on air quality, and over a particular sector. Third,
438 the proxy data is still very limited. Additional data, such as satellite images and industrial
439 outputs published by the government's statistical bureau, can be included in the regulatory
440 analysis to improve the accuracy of policy evaluation. Finally, this study only uses a single-
441 point PM_{2.5} monitor data. Given that the air quality can vary across different parts of Beijing,
442 in future work, more fine-grained air quality data obtained from the 35 official stations can be
443 used to evaluate the effects of air pollution regulatory interventions since 2013.

444

445 **5. Policy Implications**

446 Evaluating the effects of air pollution control regulations has significant implications for
447 environmental policy-makings in China and the rest of the world. We have identified two major
448 policy implications with regard to our proposed Bayesian deep-learning policy intervention
449 study methodology and results. First, our proposed data-driven regulatory analysis
450 methodology can estimate the aggregate effects of air pollution control regulations and plans
451 with reduced confounding biases and higher accuracies, when compared to other machine-
452 learning techniques. Hence, our model can provide the needed evidence to support evidence-
453 based air pollution policy-makings. For instance, the governments can perform *ex post*
454 evaluation on air pollution control regulations to test the effectiveness of the regulations they
455 implemented based on our model. Second, though the annual PM_{2.5} pollution concentration in
456 Beijing remains far beyond the WHO threshold (10 μ g/m³), our results suggest that Beijing's
457 air quality has been improved gradually over the past decade (11% improvement on average;
458 see Table 3). The air pollution control regulations implemented during 2008 – 2019 tend to be
459 more effective after 2015, i.e., after the air pollution control laws in Beijing/China have been
460 further revised and stringent air pollution control action plans have been implemented in
461 Beijing/China since 2013 (see Figure 3 and Table 3). This suggests that there is a 2-year time
462 lag before the stringent air pollution control regulations in Beijing/China taken any strong

463 positive effects. As compared to the regulatory interventions introduced before 2013, policy-
464 makings that coordinate that of the local jurisdictions and the central governments (such as the
465 guidelines on air quality monitoring and law enforcement introduced by provincial authorities,
466 effective in 2016), and laws and policies that tackle the vested interests of the local stakeholders
467 in Beijing and neighbouring cities (such as the joint action plan for air pollution control in
468 Beijing-Tianjin-Hebei Region, effective in 2013), alongside with the stringent air pollution
469 control regulations and plans, can help reduce air pollution and promote healthy living in
470 Beijing over the longer term (Lam et al., 2019).

471

472 **6. Conclusion**

473 This study extends our previous work on modelling the effects of air pollution control
474 regulations (Han et al., 2018), to investigate the effectiveness of existing and newly introduced
475 air pollution control regulations in Beijing, China during 2008 – 2019, using a Bayesian deep-
476 learning approach. Our approach can model the complex relationship between PM_{2.5} pollution
477 concentrations and other confounding factors that potentially affect PM_{2.5} pollution
478 concentrations, better address the confounding effects in policy evaluation, and predict the
479 hypothetical PM_{2.5} pollution concentrations in the absence of any regulatory interventions
480 (MAE=20.3; MAPE=36.8%). Results of our novel Bayesian deep learning regulatory
481 intervention analysis show that the PM_{2.5} pollution concentrations in Beijing were reduced by
482 11% on average, due to the aggregate effects of all regulatory interventions implemented
483 during the period of 2008 – 2019. Moreover, after the introduction of Action Plan for Clean
484 Air in China and Beijing in late 2013, as compared to the hypothetical PM_{2.5} concentration
485 (without any regulatory interventions), the estimated PM_{2.5} reduction increased dramatically
486 from 15% in 2015 to 44% in 2018. In the future, more relevant data should be collected, and
487 more advanced machine-learning methods can be used to improve the interpretability of our
488 proposed model and provide more fine-grained estimation of the regulatory effects in China
489 and elsewhere.

490

491 **Acknowledgement**

492 This research is supported in part by the General Research Fund of the Research Grants Council
493 of Hong Kong, under Grant No. 17620920. We gratefully acknowledge the PM_{2.5} data provided
494 by Beijing US Embassy, China; and AOD data provided by NASA, USA. We also gratefully
495 acknowledge the valuable comments and suggestions of our reviewers, and participants who
496 attended our first HKU-Cambridge AI-WiSe workshop, titled “AI for Social Good”, held in

497 the Department of Computer Science and Technology, Cambridge, the United Kingdom, in
498 May 2019, during which an earlier version of this work was presented.

499

500 **Declarations of interest:** none

501

502 **References**

503 Alaa, A.M., van der Schaar, M., 2017. Bayesian inference of individualized treatment effects
504 using multi-task Gaussian processes, *Advances in Neural Information Processing Systems*, pp.
505 3424-3432.

506 Alaa, A.M., Weisz, M., Van Der Schaar, M., 2017. Deep counterfactual networks with
507 propensity-dropout. arXiv preprint arXiv:1706.05966.

508 Apple Inc., 2020. Dark Sky API, <https://darksky.net/dev>

509 Atan, O., Jordon, J., van der Schaar, M., 2018. Deep-treat: Learning optimal personalized
510 treatments from observational data using neural networks, *Thirty-Second AAAI Conference*
511 *on Artificial Intelligence*.

512 Athey, S., 2017. Beyond prediction: Using big data for policy problems. *Science* 355, 483-485.

513 Athey, S., Imbens, G.W., 2017. The state of applied econometrics: Causality and policy
514 evaluation. *Journal of Economic Perspectives* 31, 3-32.

515 Beijing Municipal Bureau of Statistics, 2020. Beijing Statistical Year Book [Webpage; in
516 Chinese], <http://nj.tjj.beijing.gov.cn/nj/main/2019-tjnj/zk/indexch.htm>

517 Beijing Transport Institute, 2020. Annual Report for Transport Development in Beijing [PDF;
518 in Chinese], <http://www.bjtrc.org.cn/List/index/cid/7/p/1.html>

519 Blundell, C., Cornebise, J., Kavukcuoglu, K., Wierstra, D., 2015. Weight uncertainty in neural
520 networks. arXiv preprint arXiv:1505.05424.

521 Buuren, S.v., Groothuis-Oudshoorn, K., 2010. mice: Multivariate imputation by chained
522 equations in R. *Journal of Statistical Software*, 1-68.

523 Caliskan, A., Bryson, J.J., Narayanan, A., 2017. Semantics derived automatically from
524 language corpora contain human-like biases. *Science* 356, 183-186.

525 Chen, Y., Jin, G.Z., Kumar, N., Shi, G., 2013. The promise of Beijing: Evaluating the impact
526 of the 2008 Olympic Games on air quality. *Journal of Environmental Economics and*
527 *Management* 66, 424-443.

528 DieselNet, n.d. China: Cars and Light Trucks, <https://dieselnet.com/standards/cn/ld.php>

529 Ferraro, P.J., 2009. Counterfactual thinking and impact evaluation in environmental policy.
530 *New Directions for Evaluation* 2009, 75-84.

531 Fortunato, M., Blundell, C., Vinyals, O., 2017. Bayesian recurrent neural networks. arXiv
532 preprint arXiv:1704.02798.

533 Gal, Y., 2016. Uncertainty in deep learning. PhD thesis, University of Cambridge.

534 Han, Y., Lam, J.C.K., Li, V.O.K., 2018. A Bayesian LSTM Model to Evaluate the Effects of
535 Air Pollution Control Regulations in China, 2018 IEEE International Conference on Big Data
536 (Big Data). IEEE, pp. 4465-4468.

537 Hartford, J., Lewis, G., Leyton-Brown, K., Taddy, M., 2017. Deep IV: A flexible approach for
538 counterfactual prediction, Proceedings of the 34th International Conference on Machine
539 Learning-Volume 70. JMLR. org, pp. 1414-1423.

540 Henneman, L.R., Liu, C., Mulholland, J.A., Russell, A.G., 2017. Evaluating the effectiveness
541 of air quality regulations: A review of accountability studies and frameworks. Journal of the
542 Air & Waste Management Association 67, 144-172.

543 Imbens, G.W., Rubin, D.B., 2015. Causal inference in statistics, social, and biomedical
544 sciences. Cambridge University Press.

545 Kendall, A., Gal, Y., 2017. What uncertainties do we need in Bayesian deep learning for
546 computer vision?, Advances in Neural Information Processing Systems, pp. 5574-5584.

547 Lam, J.C.K., Han, Y., Wang, S., Li, V.O.K., Pollitt, M., Warde, P., 2019. A comparative study
548 of air pollution trends in historical London and contemporary Beijing, In Search of Good
549 Energy Policy. Cambridge University Press.

550 LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. Nature 521, 436-444.

551 Li, T., Shen, H., Yuan, Q., Zhang, X., Zhang, L., 2017a. Estimating ground-level PM2.5 by
552 fusing satellite and station observations: A geo - intelligent deep learning approach.
553 Geophysical Research Letters 44, 11,985-911,993.

554 Li, V.O.K., Lam, J.C.K., Chen, Y., Gu, J., 2017b. Deep learning model to estimate air pollution
555 using M-BP to fill in missing proxy urban data, GLOBECOM 2017-2017 IEEE Global
556 Communications Conference. IEEE, pp. 1-6.

557 Li, X., Peng, L., Yao, X., Cui, S., Hu, Y., You, C., Chi, T., 2017c. Long short-term memory
558 neural network for air pollutant concentration predictions: Method development and evaluation.
559 Environmental Pollution 231, 997-1004.

560 Li, X., Qiao, Y., Zhu, J., Shi, L., Wang, Y., 2017d. The “APEC blue” endeavor: Causal effects
561 of air pollution regulation on air quality in China. Journal of Cleaner Production 168, 1381-
562 1388.

563 Liu, X.-H., Zhang, Y., Cheng, S.-H., Xing, J., Zhang, Q., Streets, D.G., Jang, C., Wang, W.-
564 X., Hao, J.-M., 2010. Understanding of regional air pollution over China using CMAQ, part I
565 performance evaluation and seasonal variation. *Atmospheric Environment* 44, 2415-2426.

566 Liu, Y., He, K., Li, S., Wang, Z., Christiani, D.C., Koutrakis, P., 2012. A statistical model to
567 evaluate the effectiveness of PM_{2.5} emissions control during the Beijing 2008 Olympic Games.
568 *Environment International* 44, 100-105.

569 Louizos, C., Shalit, U., Mooij, J.M., Sontag, D., Zemel, R., Welling, M., 2017. Causal effect
570 inference with deep latent-variable models, *Advances in Neural Information Processing*
571 *Systems*, pp. 6446-6456.

572 Marcus, G., 2018. Deep learning: A critical appraisal. arXiv preprint arXiv:1801.00631.

573 Ong, B.T., Sugiura, K., Zettsu, K., 2016. Dynamically pre-trained deep recurrent neural
574 networks using environmental monitoring data for predicting PM_{2.5}. *Neural Computing and*
575 *Applications* 27, 1553-1566.

576 Osoba, O.A., Welser IV, W., 2017. An intelligence in our image: The risks of bias and errors
577 in artificial intelligence. Rand Corporation.

578 Pearl, J., 2018. Theoretical impediments to machine learning with seven sparks from the causal
579 revolution. arXiv preprint arXiv:1801.04016.

580 Pham, T., Tran, T., Phung, D., Venkatesh, S., 2017. Predicting healthcare trajectories from
581 medical records: A deep learning approach. *Journal of Biomedical Informatics* 69, 218-229.

582 Pope III, C.A., Dockery, D.W., 2006. Health effects of fine particulate air pollution: lines that
583 connect. *Journal of the Air & Waste Management Association* 56, 709-742.

584 Pui, D.Y., Chen, S.-C., Zuo, Z., 2014. PM_{2.5} in China: Measurements, sources, visibility and
585 health effects, and mitigation. *Particuology* 13, 1-26.

586 Rubin, D.B., 2005. Causal inference using potential outcomes: Design, modeling, decisions.
587 *Journal of the American Statistical Association* 100, 322-331.

588 Shalit, U., Johansson, F.D., Sontag, D., 2017. Estimating individual treatment effect:
589 generalization bounds and algorithms, *Proceedings of the 34th International Conference on*
590 *Machine Learning*-Volume 70. JMLR. org, pp. 3076-3085.

591 Shi, C., Blei, D., Veitch, V., 2019. Adapting neural networks for the estimation of treatment
592 effects, *Advances in Neural Information Processing Systems*, pp. 2507-2517.

593 Stolberg, H.O., Norman, G., Trop, I., 2004. Randomized controlled trials. *American Journal of*
594 *Roentgenology* 183, 1539-1544.

595 US Department of State, 2020. Beijing US Embassy Air Quality Data [CSV file],
596 [https://www.airnow.gov/international/us-embassies-and-consulates/#China\\$Beijing](https://www.airnow.gov/international/us-embassies-and-consulates/#China$Beijing)

597 US NASA, 2020. AERONET Data Download Tool [CSV file],
598 [https://aeronet.gsfc.nasa.gov/cgi-](https://aeronet.gsfc.nasa.gov/cgi-bin/webtool_aod_v3?stage=3&place_code=10®ion=Asia&state=China&site=Beijing&submit=Get+Download+Form)
599 [bin/webtool_aod_v3?stage=3&place_code=10®ion=Asia&state=China&site=Beijing&su](https://aeronet.gsfc.nasa.gov/cgi-bin/webtool_aod_v3?stage=3&place_code=10®ion=Asia&state=China&site=Beijing&submit=Get+Download+Form)
600 [bmit=Get+Download+Form](https://aeronet.gsfc.nasa.gov/cgi-bin/webtool_aod_v3?stage=3&place_code=10®ion=Asia&state=China&site=Beijing&submit=Get+Download+Form)
601 US NOAA, 2020. Hourly/Sub-Hourly Observational Data Map,
602 <https://gis.ncdc.noaa.gov/maps/ncei/cdo/hourly>
603 Wang, J.-F., Hu, M.-G., Xu, C.-D., Christakos, G., Zhao, Y., 2013. Estimation of citywide air
604 pollution in Beijing. PloS One 8, e53400.
605 Wooldridge, J.M., 2000. Econometric analysis of cross section and panel data. MIT Press.
606 Yi, X., Zhang, J., Wang, Z., Li, T., Zheng, Y., 2018. Deep distributed fusion network for air
607 quality prediction, Proceedings of the 24th ACM SIGKDD International Conference on
608 Knowledge Discovery & Data Mining, pp. 965-973.
609 Yoon, J., Jordon, J., van der Schaar, M., 2018. GANITE: Estimation of individualized
610 treatment effects using generative adversarial nets. International Conference on Learning
611 Representations.
612 Zhang, H., Wang, S., Hao, J., Wang, X., Wang, S., Chai, F., Li, M., 2016. Air pollution and
613 control action in Beijing. Journal of Cleaner Production 112, 1519-1527.
614 Zhang, Q., Zheng, Y., Tong, D., Shao, M., Wang, S., Zhang, Y., Xu, X., Wang, J., He, H., Liu,
615 W., 2019. Drivers of improved PM_{2.5} air quality in China from 2013 to 2017. Proceedings of
616 the National Academy of Sciences 116, 24463-24469.
617 Zheng, S., Yi, H., Li, H., 2015. The impacts of provincial energy and environmental policies
618 on air pollution control in China. Renewable and Sustainable Energy Reviews 49, 386-394.

619

620 **Table 1.** Data source

Data	Resolution	Variable	Source
Air quality	Hourly station-level (aggregated into daily means)	PM _{2.5} concentrations (ug/m ³)	US Department of State (2020)
Meteorology	Hourly city-level (aggregated into daily means)	Temperature, relative humidity, air pressure, wind speed, wind bearing, and visibility	Apple Inc. (2020) ¹

AOD	All city-level observations per day (aggregated into daily means)	AOD at 1020 nm, AOD at 870 nm, AOD at 675 nm, AOD at 440 nm, and precipitable water (cm)	US NASA (2020)
Socio-economic	Yearly	Population density (population per km ²), percentage of GDP generated from the secondary sector, and the number of vehicles	Beijing Municipal Bureau of Statistics (2020), Beijing Transport Institute (2020)
Notes			
1. The weather data application program interface (API) no longer accepts new signups (Apple Inc., 2020). The historical meteorology data in Beijing can also be downloaded from the US's National Climatic Data Center (US NOAA, 2020).			

621

622 **Table 2.** Comparison of the performance between Bayesian deep-learning and other baseline
623 air pollution regulatory intervention models based on the test set

Model	MAE ¹	MAPE ¹
SVR	22.1 (1.6)	38.8% (3.4%)
RF	22.4 (1.7)	46.9% (4.9%)
Bayesian LSTM	20.3 (0.6)	36.8% (1.8%)
Notes		
1. Standard deviation is shown in parenthesis.		

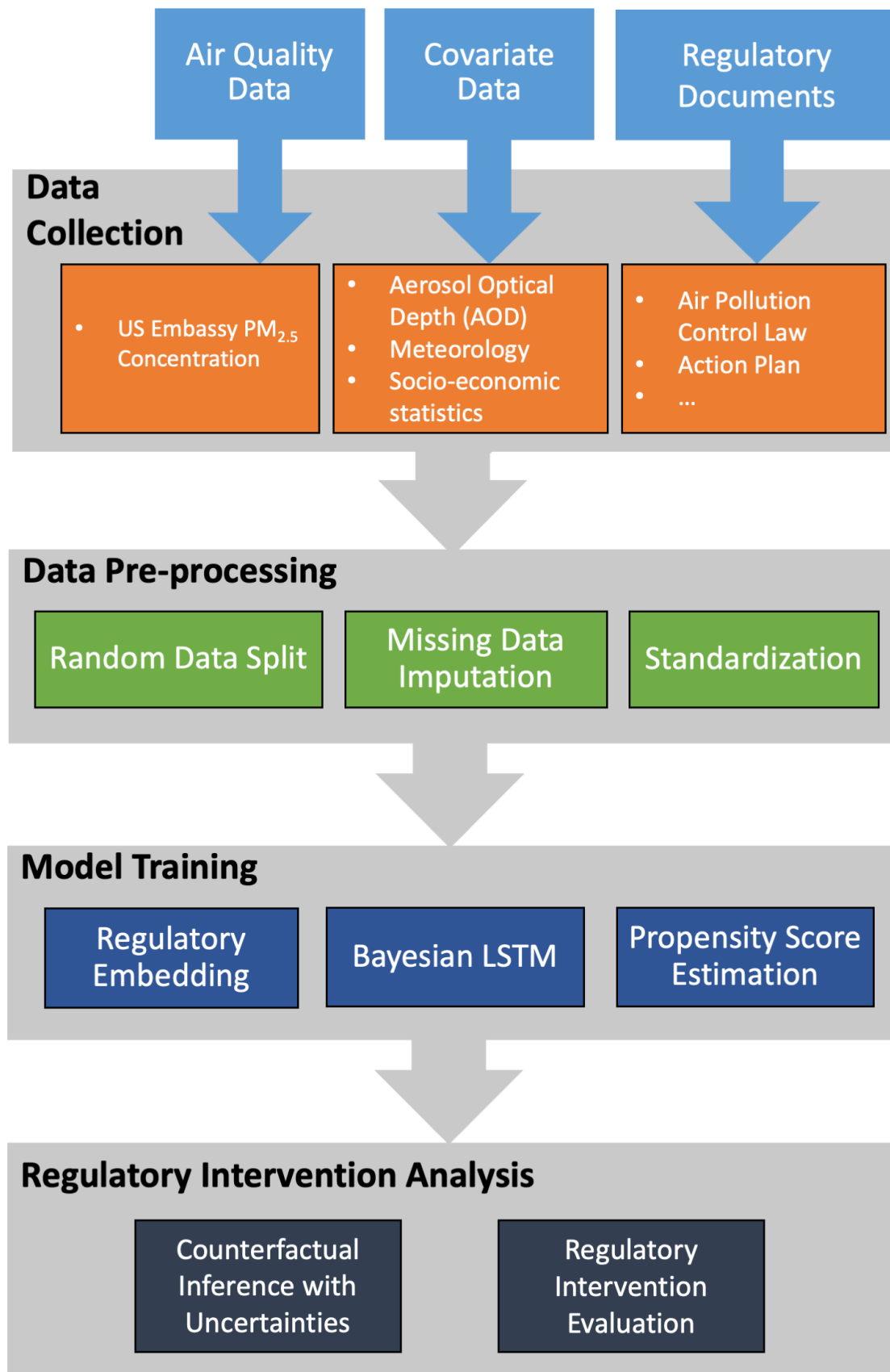
624

625 **Table 3.** Annual PM_{2.5} reduction due to local and national air pollution control regulations
626 implemented in Beijing, China

Year	Observed PM _{2.5} (µg/m ³)	Simulated PM _{2.5} (µg/m ³)	PM _{2.5} reduction (µg/m ³)	Relative reduction of PM _{2.5}
2008	92	96	4	4%
2009	102	99	-3	-3%
2010	104	100	-4	-4%
2011	98	99	1	1%
2012	91	98	7	7%

2013	101	101	0	0%
2014	98	100	2	2%
2015	82	97	15	15%
2016	73	96	23	24%
2017	59	94	35	37%
2018	51	91	40	44%
2019	58	92	34	37%
2008 – 2019	86	97	11	11%

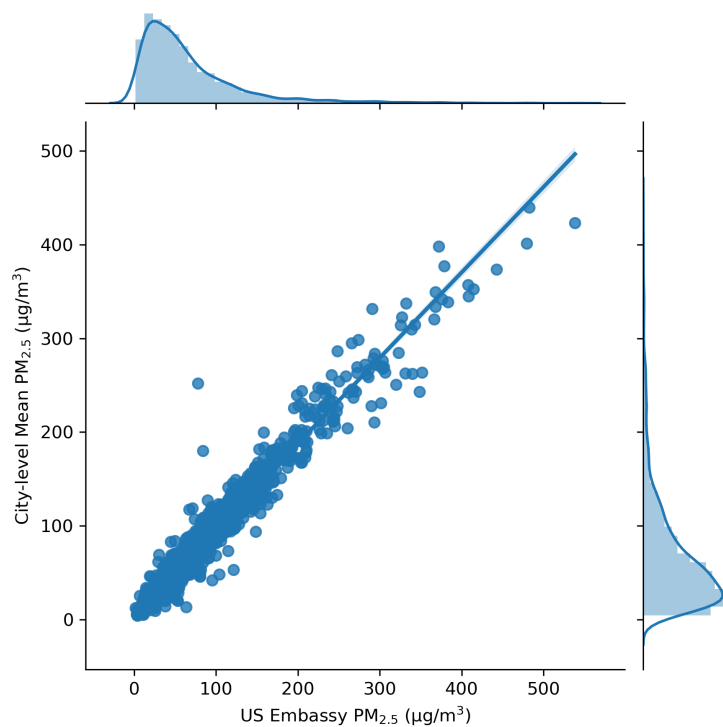
627



628

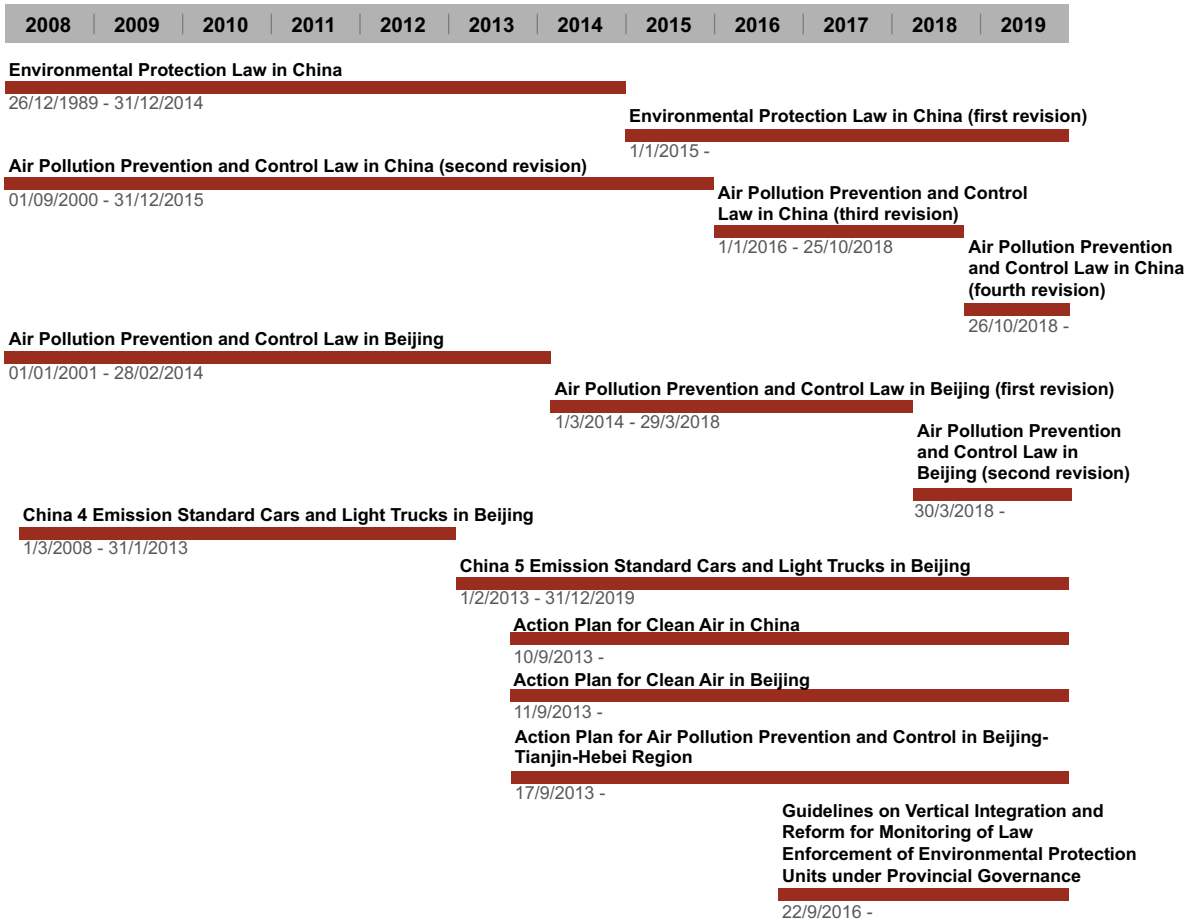
629 **Figure 1.** The overall framework of our proposed Bayesian deep-learning regulatory
 630 intervention analysis

631



632

633 **Figure 2.** Correlation between the daily PM_{2.5} concentrations monitored at the US Embassy,
634 Beijing and the daily city-level average PM_{2.5} concentrations monitored at the 35 official
635 stations in Beijing, 2014 – 2019

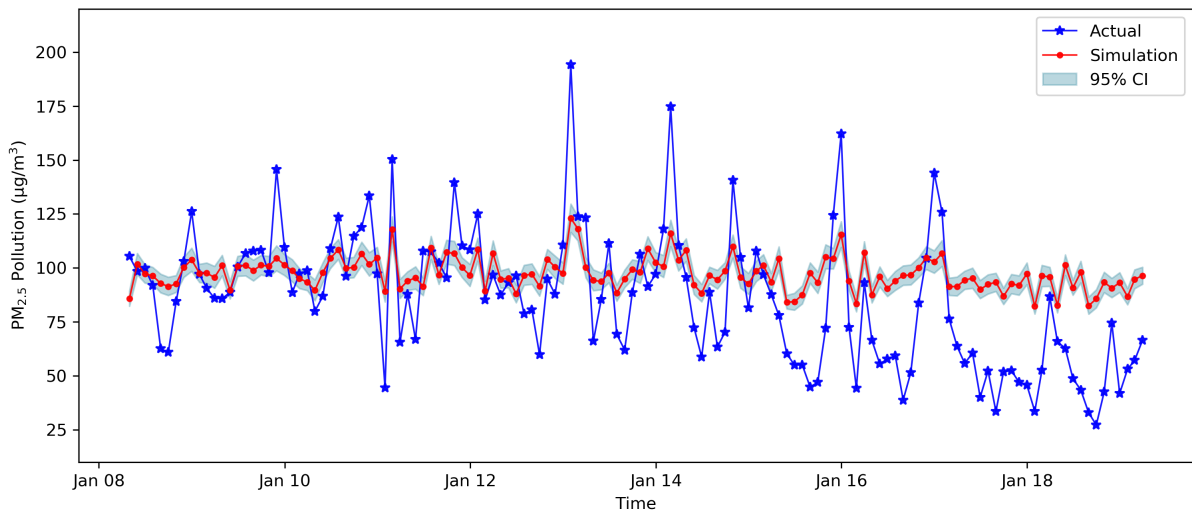


636

637 **Figure 3.** Timeline of air pollution control regulations implemented in Beijing/China during

638 2008 – 2019

639



640

641 **Figure 4.** The monthly trend of observed PM_{2.5} pollution concentrations (with regulatory

642 interventions) and simulated PM_{2.5} pollution concentrations (without any regulatory

643 interventions) during 2008 – 2019

644