

Estimating construction waste generation in the Greater Bay Area, China using machine learning

Weisheng Lu^a, Jinfeng Lou^{a,1}, Chris Webster^b, Fan Xue^a, Zhikang Bao^a, and Bin Chi^c

^a Dept. of Real Estate and Construction, Faculty of Architecture, The University of Hong Kong, Pokfulam, Hong Kong

^b Faculty of Architecture, The University of Hong Kong, Pokfulam, Hong Kong

^c Faculty of Built Environment, University of New South Wales, Sydney, Australia

Emails: wilsonlu@hku.hk (WL), waseljf@connect.hku.hk (JL), cwebster@hku.hk (CW), xuef@hku.hk (FX), u3004700@hku.hk (ZB), simon.binchi@gmail.com (BC).

This is the peer-reviewed post-print version of the paper:

Lu, W., Lou, J., Webster, C., Xue, F., Bao, Z., & Chi, B. (2020). Estimating construction waste generation in the Greater Bay Area, China using machine learning. *Waste Management*, 134, 78-88. Doi: [10.1016/j.wasman.2021.08.012](https://doi.org/10.1016/j.wasman.2021.08.012).

The final version of this paper is available at: <https://doi.org/10.1016/j.wasman.2021.08.012>.

The use of this file must follow the [Creative Commons Attribution Non-Commercial No Derivatives License](https://creativecommons.org/licenses/by-nc/4.0/), as required by [Elsevier's policy](https://www.elsevier.com/locate/elsevierpolicy).

Abstract

Reliable construction waste generation data is a prerequisite for any evidence-based waste management effort, but such data remains scarce in many developing economies owing to their rudimentary recording systems. By referring to several models proposed for estimating waste generation, this study aims to develop a reliable and accessible method for estimating construction waste generation based on limited publicly available data. The study has two objectives. Firstly, it aims to estimate construction waste generation by focusing on the Greater Bay Area (GBA) in China, one of the world's most thriving regions in terms of construction activities. Secondly, it aims to compare the strengths and weaknesses of various waste quantification models. 43 sets of annual socio-economic, construction-related and C&D waste generation data ranging from 2005 to 2019 were collected from the local government authorities. By analyzing the data using four types of machine learning models, namely multiple linear regression, decision tree, grey models, and artificial neural network, it is found that all calibrated models, with their respective strengths and weaknesses, can produce acceptable results with the testing R^2 ranging from 0.756 to 0.977. This study also reveals that the 11 cities in the GBA produced a total of about 364 million m³ of construction waste in 2018. The result can be used for monitoring the urban metabolism, quantifying carbon emission, developing a circular economy, valorizing recycled materials, and strategic planning of waste management facilities in the GBA. The research findings also contribute to the methodologies for estimating waste generation using limited data.

Keywords: Construction waste; waste quantification; Greater Bay Area, China; machine learning.

¹ Corresponding author details: Email: waseljf@connect.hku.hk, Tel: +852 94231962, Address: 7/F, Knowles Building, The University of Hong Kong, Pokfulam Road, Hong Kong

List of abbreviations:

ANN = artificial neural network

BIM = building information modeling

CART = classification and regression tree

30 C&D = construction and demolition

CO = total construction output

CP = construction productivity

CSA = classification system accumulation

CWDCS = construction waste disposal charging scheme

35 CWM = construction waste management

DT = decision tree

DS = development stages of an economy

FC = floor space completed

FCO = floor space under construction

40 FD = floor space of demolition

FM = factor modeling

FS = floor space of newly started buildings

GBA = greater Bay Area

GC = GDP per capita

45 GDP = gross domestic product

GM = grey models

GMC = grey model with convolution integral

GRA = grey relational analysis

GRC = generation rate calculation

50 KMO = Kaiser-Meyer-Olkin

LA = lifetime analysis

ML = machine learning

MLR = multiple linear regression

MSW = municipal solid waste

55 PCA = principal component analysis

PO = population

R^2 = coefficient of determination

SD = standard deviation

SV = site visit

60 VIF = variance inflation factor

63 **1. Introduction**

64 Construction waste, a term often used interchangeably with construction and demolition
65 (C&D) waste, is the solid waste generated by construction, renovation, or demolition
66 activities (HKEPD, 2015; USEPA, 2016). It comprises inert and non-inert materials including
67 concrete, steel, slurry, wood, glass, etc. C&D waste contributes significantly to
68 environmental degradation (Coelho & De Brito, 2012; Wang et al., 2010), consumes valuable
69 landfill space (Poon et al., 2004), causes geologic hazards and other undesirable
70 consequences (Lu, 2019; Perlez, 2016). Therefore, it needs to be carefully managed.

71
72 Information about waste generation is a prerequisite for many waste management strategies,
73 including planning landfill space, determining levies for polluters or subsidies for recyclers,
74 and scheduling companies' waste management policies. Since what cannot be measured
75 cannot be improved, estimation of waste generation at both regional and project levels has
76 begun to receive worldwide research attention. Wu et al. (2014), for instance, have reviewed
77 57 studies in C&D waste quantification. Examples of regional studies include Cochran et al.
78 (2007) exploring the accounting, generation, and composition of building-related C&D waste
79 in Florida, and Lu et al. (2017) who estimate that approximately 1.13 billion tons of C&D
80 materials were generated in China during 2014. This paper also has a regional focus.

81
82 For regions where waste generation data is regularly collected and released by official
83 recording systems, estimation is unnecessary (Lu et al., 2017). However, emerging regions
84 often do not have such systems in place. There are a plethora of studies on solid waste
85 quantification in the absence of direct data, where other statistics or signs such as population,
86 economic growth, construction expenditure, urban decay, and waste recycling levels, are
87 analyzed to inform waste generation. Such studies often adopt complicated algorithms to
88 estimate waste generation, but may exhibit overfitting where models report closely or exactly
89 fitting results in training datasets but poor results in testing datasets. In addition, few studies
90 have gone beyond the factors that can predict C&D waste generation to understand how
91 much each factor contributes to the prediction.

92
93 This study has two purposes. Firstly, it is to estimate construction waste generation in the
94 Greater Bay Area (GBA) of South China. The GBA is chosen for several reasons. It is among
95 China's most economically active areas, and one in which intense construction activity exists
96 in conflict with the severe environmental degradation it causes. The GBA comprises 11
97 regions including Hong Kong and Macau (both under the "one country, two systems"
98 constitutional framework), Shenzhen, Guangzhou, and others. Among these 11 regions,
99 economic development is imbalanced and recording systems vary in reliability. The second
100 purpose of our study is to compare the strengths and weaknesses of waste estimation
101 algorithms in terms of accuracy, scalability, and explanatory clarity, and also consider
102 overfitting issues.

103

104 2. Estimating solid waste generation

105 The amount of construction waste generated can be affected by an ocean of factors. Table 1
 106 summarizes the factors that have been used to predict C&D waste generation at a regional
 107 level. These factors are of two types: socio-economic or construction-related. Socio-
 108 economic factors include gross domestic product (GDP), GDP per capita, population, and
 109 others acting as indicators of socio-economic development and providing the context for
 110 construction industry development. It has been proven that C&D waste generation ascends in
 111 parallel with population expansion, urbanization, and economic development (Kofoworola &
 112 Gheewala, 2009; Zhao et al., 2011). Construction-related factors include total construction
 113 output, floor space of newly started buildings, floor space completed, and so on (see Table 1).
 114 Although it is impossible to obtain the direct amount of C&D waste generation, these factors
 115 with reasonable data availability and proper analytics can yield a satisfactory estimate of
 116 C&D waste generation.

117

118

Table 1 Factors impacting construction waste generation

Reference	Level	Socio-economic factors	Construction-related factors
Hsiao et al. (2002)	City	-	FS, FD
Kofoworola and Gheewala (2009)	Country	PO	FS, FCO
Zhao et al. (2011)	City	PO, GDP	FCO, FD
Song et al. (2015)	City	-	FC
Tam and Lu (2016)	Country	GDP, GC, DS	CO, CP, FS
Lu et al. (2017)	Country	GDP	CO, FCO, FC
Song et al. (2017)	Country	-	FC

119 Note:

- 120 1. Socio-economic factors: DS – development stages of an economy; GC – GDP per capita; PO – population
- 121 2. Construction-related factors: CO – total construction output; CP – construction productivity; FC – floor
- 122 space completed; FCO – floor space under construction; FD – floor space of demolition; FS – floor space
- 123 of newly started buildings

124

125 With the potential factors known, numerous methods have been proposed for estimating
 126 construction waste generation. Wu et al. (2014) categorize these methods into six types: site
 127 visit (SV), generation rate calculation (GRC), lifetime analysis (LA), classification system
 128 accumulation (CSA), factor modeling (FM), and others (e.g., BIM-based automated
 129 estimation). The SV method requires the investigator to conduct surveys on site, including
 130 direct measurement by surveying the weight or volume (Hoang et al., 2020; Lau et al., 2008)
 131 and indirect measurement by adopting other easily accessible indicators, such as hauling
 132 tickets (Bakchan & Faust, 2019). For GRC method, the total waste volume can be calculated
 133 through multiplying the quantity of a specific unit by its corresponding generation rate, e.g.,
 134 area-based calculation (Domingo & Batty, 2021; Hoang et al., 2021). The LA method
 135 assumes that all buildings must be dismantled after a certain period of lifetime and the C&D
 136 waste can be deduced from calculating the sum of the mass to be removed at expiration

137 (Huang et al., 2013). The CSA method combines the GRC method with a waste classification
138 to quantify each specific material (Hu et al., 2021). The FM method designs prediction
139 models based on accessible variables, such as linear regression models (Kern et al., 2018).
140 Others are infrequent methods that do not fit into any of the above categories, such as BIM-
141 based automated estimation (Guerra et al., 2019). To summarize, the methods to estimate
142 C&D waste generation have their own strengths and weaknesses.

143

144 In this study, estimating waste generation in the GBA is a problem at a regional level.
145 Previous studies have engaged FM methods, particularly in MSW or WEEE. This study
146 therefore will deploy FM methods as well while keeping an eye on other methods as
147 reviewed above. Table 2 provides a detailed analysis of previous studies examining locations,
148 levels, methods, models, data, and estimate performance. It can be seen that the machine
149 learning (ML) models, such as multiple linear regression (MLR), grey model (GM), artificial
150 neural network (ANN), and decision tree (DT), are among the most frequently adopted FM
151 methods. ML is a computer program that “optimize a performance criterion using example
152 data or past experience” (Alpaydin, 2020). The ability to automatically learn from data and
153 adapt to changes is the key to ML applications. ML algorithms may not learn everything from
154 data, but can still identify some patterns or regularities which are presented to humans, either
155 explicitly or implicitly. The ML models or algorithms as shown in Table 2, including the type
156 and size of data and model performance, provide very useful references for this study.

157

Table 2 Machine learning models used to estimate the solid waste generation at a regional level

Method	Model	Waste Type	Region	Region level	Type of data	No. of data	Level of data collection	Model input	Model perform. (R^2)	Reference
Linear regression	MLR	MSW	Nigeria	Country	Panel data (Monthly)	166	Household	Household income, household size, educational background, social status, occupation, and season of the year	0.88	Afon and Okewole (2007)
	MLR	MSW	Mexico	City	Panel data	181	Household	Education, income per household, and number of residents	0.51	Benítez et al. (2008)
	MLR	MSW	Iraq	City	Cross-sectional data	150	Household	Hotel size, expenditure, area and number of staffs	0.80	Abdulredha et al. (2018)
	MLR	MSW	China	City	Panel data (Yearly)	10	City	Population, total consumption expenditure	0.94	Yuan et al. (2012)
	MLR	MSW	Vietnam	City	Cross-sectional data	100	City	Household size and household income	0.36	Thanh et al. (2010)
	MLR	MSW	India, Nepal, Pakistan, Bangladesh, Sri Lanka	Country	Panel data (Yearly)	35	Country	Population, GDP per capita, illiteracy rate	-	Khajuria et al. (2010)
	MLR	MSW	Europe	Country level, City level	Panel data (Yearly)	86	Country level, City level	GDP, population, infant mortality rate, household size, life expectancy at birth	0.65	Beigl et al. (2004)
Grey model	GM(1,1)	WEEE	China	Country	Panel data (Yearly)	13	Country	Quantity of household appliances	-	Zhao et al. (2016)
	GMC(1,n), NBGMC(1,n)	WEEE	USA	City	Panel data (Yearly)	13	City	Population Density, Household income	0.99	Duman et al. (2019)
	GM (1, 1), GM (1, 1)- α , GM (1, n) and GMC (1, n)	MSW	Thailand	Country	Panel data (Yearly)	13	Country	Household expenditure, household size, employment, population density, and urbanization	-	Intharathirat et al. (2015)
	GM(1,1), GIM(1), GPPM(1) and GLPM(1)	MSW	China	City	Panel data (Yearly)	14	City	GDP, population, size, total retail sales, consumption of gas, water and electricity, personal salary	0.98	Liu and Yu (2007)
	GM(1,n)	MSW	China	City	Panel data (Yearly)	10	City	GDP, population, household expenditure, total sales of consumer goods	0.68	Wang et al. (2012)
	GM(1,1), GM(1,n)	MSW	China	City	Panel data (Yearly)	9	City	GDP, population, retail sales, consumer spending	0.95	Zhang (2013)
	Neural network	ANN, ANFIS, DWT-ANN, DWT-ANFIS, GA-ANN, GA-ANFIS	MSW	India	City	Panel data (Yearly)	19	City	Previous waste generation	0.87
ANN, GM(1,1), MLR		MSW	China	Country	Panel data (Yearly)	16	Country	GDP, population, urbanization, energy consumption	0.93	Chhay et al. (2018)
ANN, SVM, ANFIS, kNN		MSW	Australia	City	Panel data (Monthly)	216	City	Previous waste generation	0.98	Abbasi and El Hanandeh (2016)
ANN, ANFIS, SVM, LSSVM, FSVM, MLR		MSW	Iran	City	Cross-sectional data	105	Household	Hospital's wards, staff, ownership type, inpatients	0.92	Golbaz et al. (2019)

	ANN	WEEE	USA	Country	Panel data (Yearly)	9	Country	Previous waste generation	-	Milojkovic and Litovski (2008)
	ANN, PCA-MLR	MSW	Iran	City	Panel data (Weekly)	158	City	Previous waste generation, number of waste truck	-	Noori et al. (2009)
	ANN	MSW	Serbia	Country	Cross-sectional data	54	Country	Income, employment, age, education, housing condition	0.96	Batinić et al. (2011)
	ANN	MSW	India	City	Cross-sectional data	98	City	Population, previous waste generation, longitude, latitude, tax	-	Patel and Meka (2013)
	ANN	MSW	Iran	City	Panel data (Weekly)	144	City	Previous waste generation, number of waste truck	0.75	Jalali and Nouri (2008)
Decision tree	DT, ANN	MSW	Canada	City	Cross-sectional data	1553	City	Population, income, employment, education, housing condition	0.72	Kannangara et al. (2018)
	DT, SVM, RNN	MSW	Colombia	City	Panel data (Monthly)	60	City	Population, socio-economic stratification, latitude and altitude	-	Meza et al. (2019)
Other machine learning methods	GM-SVR	C&DW	China	Country	Panel data (Yearly)	30	Country	Total floor areas completed	0.99	Song et al. (2017)
	GBRT	MSW	USA	City	Panel data (Weekly)	41412	Building	Building attributes, socio-economic and demographic feature, weather	0.87	Kontokosta et al. (2018)
System dynamics	System dynamics model	MSW	USA	City	Panel data	3	City	Population, income, household size, and employment	0.99	Dyson and Chang (2005)
Time series	ARIMA	C&DW	China	City	Panel data (Yearly)	9	City	Total floor areas completed	-	Song et al. (2015)

Note:

1. MSW – municipal solid waste; WEEE – waste electrical and electronic equipment; C&DW – C&D waste
2. Panel data – data refers to multi-dimensional data frequently involving measurements over time; Cross-sectional data – data collected by observing many subjects at the same point of time or without regard to differences in time
3. MLR – Multiple linear regression; GM – Grey model; GMC – Grey model with convolution integral; NBGMC – Nonlinear grey Bernoulli model with convolution integral; GIM – Grey index model; GPPM – Grey parabola power model; GLPM – Grey logarithm power model; ANN – Artificial neural network; ANFIS – Adaptive neuro-fuzzy inference system; DWT – Discrete wavelet theory; GA – Genetic algorithm; SVM – Support vector machine; kNN – k-nearest neighbors; SVR – Support Vector Regression; LSSVM – Least squares support vector regression; FSVM – Fuzzy logic support vector regression; PCA – Principal component analysis; DT – Decision tree; RNN – Recurrent neural network; GBRT – Gradient boosting regression tree; ARIMA – Autoregressive integrated moving average
4. R^2 is the best testing performance of all models if there is more than one model in the literature. If R^2 wasn't given, it is calculated by the authors according to the prediction results in the literature.

3. The Greater Bay Area

The GBA comprises the two SARs (special administration regions) Hong Kong and Macau and nine municipalities in China's Guangdong province. In 2019, the GBA occupied a total area of about 56,000 km² and had a GDP of USD 1,679.5 billion (CMAB, 2020). Although it occupies less than 1% of China's land area, the GBA's contribution to national GDP is up to 12% (Cheung, 2019), making it one of the most economically vibrant regions in China. To accommodate its economic activities, massive construction activities have been undertaken or are underway to materialize the supporting infrastructure and building in the GBA. Meanwhile, huge amounts of C&D waste have been produced. For example, Hong Kong generated about 18.12 million tons of construction waste in 2018 (HKEPD, 2019). If not properly managed, such vast quantities of waste are bound to hinder the sustainable development of the GBA and cause harm to the inhabitants. In an extreme case, a construction waste landslide in 2015 in Shenzhen resulted in 73 deaths and ruined over 30 buildings (Perlez, 2016).

Having recognized the importance of proper construction waste management (CWM), some regions in the GBA have deployed response strategies. In 2006, Hong Kong launched its construction waste disposal charging scheme (CWDCS) under which contractors are charged HK\$71 to \$200 per ton for waste mandatorily disposed of at designated facilities (Bao et al., 2020). Facing increased pressure after the tragedy in 2015 to better manage its construction waste, Shenzhen has closed all landfills so that contractors are forced to reduce, reuse and recycle construction waste (Bao & Lu, 2020). In some exemplar sites, zero waste is pursued (Lu, Bao, et al., 2021). In recognition of the imbalance in demand and supply among GBA regions, construction waste material sharing has been actively explored. In fact, since 2006 Hong Kong has been sending its construction waste materials to Jiangmen through an official channel for land reclamation (Lu et al., 2020). However, such efforts are still too piecemeal and discrete. Integrated policies and measures are being sought, but reliable data is a prerequisite for their formulation. Hong Kong and Macau have a long-established recording system with detailed waste generation data, which enables better CWM practice. For example, Ahmed and Zhang (2021) developed a multi-stage network-based model to reduce the logistics cost for inert waste management and validated it using sufficient data from Hong Kong. However, other GBA regions may only possess broad socio-economic background data and lack accurate waste quantity and distribution (Ma et al., 2020). Nonetheless, based on previous studies (Li et al., 2020), it is possible to estimate construction waste generation by extrapolating from data-rich to data-scarce regions.

4. Research methods

210 This study adopted a four-step research method, including (i) collecting relevant data, (ii) selecting alternative models, (iii) developing the selected models, and (iv) cross-validation of models.

4.1 Data collection

215 Factors that impact C&D waste generation were carefully selected from the literature, based on our own knowledge, and also depending on data availability. In the end, selected factors were (1) population (PO), (2) GDP per capita (GC), (3) total construction output (CO), (4) floor space of newly started buildings (FS), and (5) floor space of buildings completed (FC).

220 Data relating to these five factors (i.e., the model input) was collected from statistical yearbooks of the National Bureau of Statistics of China. C&D waste generation data (i.e., model output) is quite limited compared to MSW data. However, with increasing C&D waste and city management system maturity, local governments have begun to pay more attention to effective C&D waste management. For example, the Guangzhou Bureau of Ecology and
225 Environment started incorporating C&D waste generation data into its annual solid waste management report in 2016, while the Shenzhen Housing and Construction Bureau began to provide C&D waste generation data from 2014. Hong Kong and Macau's C&D waste generation data has been available from the Hong Kong Environment Protection Department since 2005 and the Macau Environmental Protection Bureau since 2010, respectively.
230 Shanghai has relatively good waste generation data and although not within the GBA is comparable with Shenzhen in terms of construction and waste management. Therefore, statistics from Shanghai were also collected for this study.

In total, 43 sets of data were collected based on the availability of annual C&D waste
235 generation data, as shown in the Supplementary Materials. Shanghai and Hong Kong measure C&D waste generation by weight (ton), while other cities measure by volume (cubic meter). For a better comparison, the weight unit was aligned to the volume unit by the bulk density of C&D waste. Lu, Yuan, et al. (2021) calculated the bulk density by analyzing 4.9 million truckloads of C&D waste. The results show that the 5% to 95% percentile interval of bulk
240 density ($\rho_{[5\%,95\%]}$) is [0.266 tons/m³, 0.971 tons/m³], with the mean value ($\bar{\rho}$) of 0.528 tons/m³ and the median value ($\rho_{0.5}$) of 0.476 tons/m³. In this study, the density of C&D waste for conversion took a value of 0.5 tons/m³.

The descriptive statistics of the collected data and the Pearson correlation coefficient between
245 the model inputs and outputs can be found in the Supplementary Materials. The correlation analysis can serve as a preliminary screening of factors for further modeling (Kannangara et al., 2018). The correlations can be considered negligible when the absolute value of the

Pearson correlation coefficient is less than 0.3 (Pallant, 2011). In this study, all such values are above 0.6, demonstrating that these factors can be adopted for modeling. The strength of the correlation is in order of PO, CO, FS, FC, and GC.

4.2 Model selection

In view of the strong ability to extract experience from the previous data, four ML models, namely MLR, DT, GM, and ANN were adopted in this study. MLR and DT have strong interpretability, so were selected to provide explanatory results. GM and ANN are good at fitting and GM in particular has been used extensively in small-sample prediction. These two models were chosen to provide accurate predictive results.

4.3 Model development

4.3.1 Multiple linear regression (MLR)

The MLR model adopts a linear equation shown as Eq. (1):

$$Y = B_0 + \sum_{i=1}^m B_i x_i \quad (1)$$

where Y is the output, i.e., the total C&D waste volume in this study; x_i is the input factors; B_0 and B_i are the regression coefficient; m is the number of data points. The MLR model was built by fitting this linear equation to input data. The logarithmic transformation operation was applied to (i) improve data normality; and (ii) reduce the difference between data in the same dimension. There are of course significant differences in data of different-sized cities. For example, in 2019, C&D waste of Guangzhou is around 50 times that of Macau. When fitting an MLR model, a slight change in the regression coefficient (slope) can lead to large variances in the prediction values for smaller cities. Furthermore, negative predictive values might occur, which is obviously not reasonable. Therefore, this study takes $[\lg(x_i)]$ as the inputs and $[\lg(Y)]$ as the output. In this way, the range of data in a specific dimension is reduced, and there are no negative prediction values.

Before MLR modeling, a principal component analysis (PCA) was conducted to (i) identify the principal components for MLR; and (ii) eliminate the multicollinearity between the inputs, which might affect the outcome of MLR (Pallant, 2011). The variance inflation factor (VIF), shown as Eq. (2), can be used to evaluate multicollinearity.

$$VIF_x = \frac{1}{1 - R^2} \quad (2)$$

where R^2 is the coefficient of determination for the input X . A VIF value of less than 10 is acceptable (Abdulredha et al., 2018). In this study, after PCA, the VIF values of all the extracted components were reduced to 1.0, i.e., the perfect VIF value, indicating that the multicollinearity was eliminated and the data was suitable for MLR. Five MLR models (identified by PCA-MLR- t , $t = 1, 2, \dots, 5$) were built based on the first t component(s). The PCA results can be found in the Supplementary Materials.

4.3.2 Decision tree (DT)

A decision tree has one root node, internal nodes, and leaf nodes (Tayefi et al., 2017). At the root or internal nodes, the division happens where the information gain reaches its maximum, and the purity of data contained in the sub-nodes increases. The data is divided into smaller groups recursively until certain criteria are met. The classification and regression tree (CART) algorithm was used in this study. The stop criterion for division was set by limiting the number of data points at each leaf node. CART is a binary tree built by greedy algorithm. This means the binary division only reaches the local optimum, without considering the best partition for all the nodes (Kannangara et al., 2018). In this case, the results might be local minima. To solve this problem, the CART algorithm was trained repeatedly with different initial training data. Its performance and errors were analyzed comprehensively. Moreover, since CART is prone to overfitting, some trivial branches were removed by post-pruning, increasing the generality of the CART model (Bramer, 2007). In this study, the cases where the minimum number of data points at leaf nodes (minimum leaf size) is 1, 2, 4, 6, 8, 10, respectively, were investigated. These models are identified by DT- t ($t = 1, 2, 4, 6, 8, 10$).

4.3.3 Grey model (GM)

Construction waste volume is interpreted by a great number of factors, requiring the multivariate grey model, GM(m,n). In GM(m,n), m denotes the order of differential equations, while n denotes the number of variables, including input variables and output variables (Duman et al., 2019). The first-order GM, GM($1,n$), has been widely used in terms of prediction and proven to bear high accuracy. The grey model with convolution integral, GMC($1,n$), one of the variants of GM, can achieve higher accuracy than GM($1,n$) (Intharathirat et al., 2015). In this study, GM($1,n$) and GMC($1,n$) were used to predict C&D waste generation and are presented as differential equations, Eq. (3) and Eq. (4), respectively:

$$x_1^{(0)}(k) + az_1^{(1)}(k) = \sum_{i=2}^n b_i x_i^{(1)}(k) \quad (3)$$

$$x_1^{(0)}(k) + az_1^{(1)}(k) = \sum_{i=2}^n b_i x_i^{(1)}(k) + u \quad (4)$$

where a , b_i and u are model parameters; $x_i^{(0)}(k)$, denotes k th element of the sequence of i th factors; $x_i^{(1)}(k)$ denotes the k th accumulated generating operation (AGO) values of the sequence of i th factor; $z_1^{(1)}(k) = 0.5x_1^{(1)}(k) + 0.5x_1^{(1)}(k-1)$; for output factor, $i = 1$, for input factors, $i = 2, 3, \dots, n$.

Before constructing the GM, grey relational analysis (GRA) was carried out to rank the factors according to their grey relational grades. Each input factor was compared with the model output in regard to variation tendency in order to determine the grey relational grade (Hsu & Wang, 2009). According to the ranked factors for GRA, 10 GMs were trained,

including [GM(1,2) – GM(1,6)] and [GMC(1,2) – GMC(1,6)]. Each GM(1,*n*) or GMC(1,*n*) considered one model output factor and the first (*n*-1) input factor(s).

325

4.3.4 Artificial neural network (ANN)

An ANN consists of three kinds of layers: input, hidden, and output. The neurons of the input layer and output layer are equal to the number of inputs and outputs, respectively. There can be a single hidden layer or multiple hidden layers, and each hidden layer can have multiple neurons, leading to the diversiform ANN architectures.

330

This study adopts a feed-forward neural network with a single hidden layer. This ANN architecture has been employed frequently and has performed well (Ojha et al., 2017). Given enough neurons in the single hidden layer, this ANN model can handle arbitrarily complex problems (Duka, 2014). In this study, the number of neurons in the hidden layers was taken as 3, 5, 10, 15, 30, and 50, respectively (identified by ANN-*t*, *t* = 3, 5, 10, 15, 30, 50). The sigmoid transfer function served as the activation function within ANN (Kannangara et al., 2018). The network was trained with the Levenberg-Marquardt backpropagation algorithm (Yu & Wilamowski, 2011). Similar to MLR, to avoid negative prediction values, the logarithm of collected data was taken as the inputs and outputs of the ANN model.

335

340

4.4 Cross validation

To validate these models, the data set was randomly divided into a training set and a testing set by the ratio of 80:20 (Azadi & Karimi-Jashni, 2016). Different partitions of a dataset result in different models, and some are quite sensitive to training data (Cunningham et al., 2000). These models might fall into the local minima, meaning the model is not the optimal solution, especially for DT and ANN. Therefore, 50 iterations of random partitions were performed for each model, and each model was trained 50 times. The 50 iterations were averaged for performance evaluation. The coefficients of determination (R^2) were used to evaluate the training and testing performance. They can be calculated by Eq. (5):

345

350

$$R^2 = 1 - \frac{\sum_{k=1}^m (\hat{Y}_i - Y_i)^2}{\sum_{k=1}^m (Y_i - \bar{Y})^2} \quad (5)$$

where m is the number of data points; \hat{Y}_i is the forecast value of the total C&D waste volume; Y_i is the actual value of the total C&D waste volume; \bar{Y} is the average value of Y_i .

355

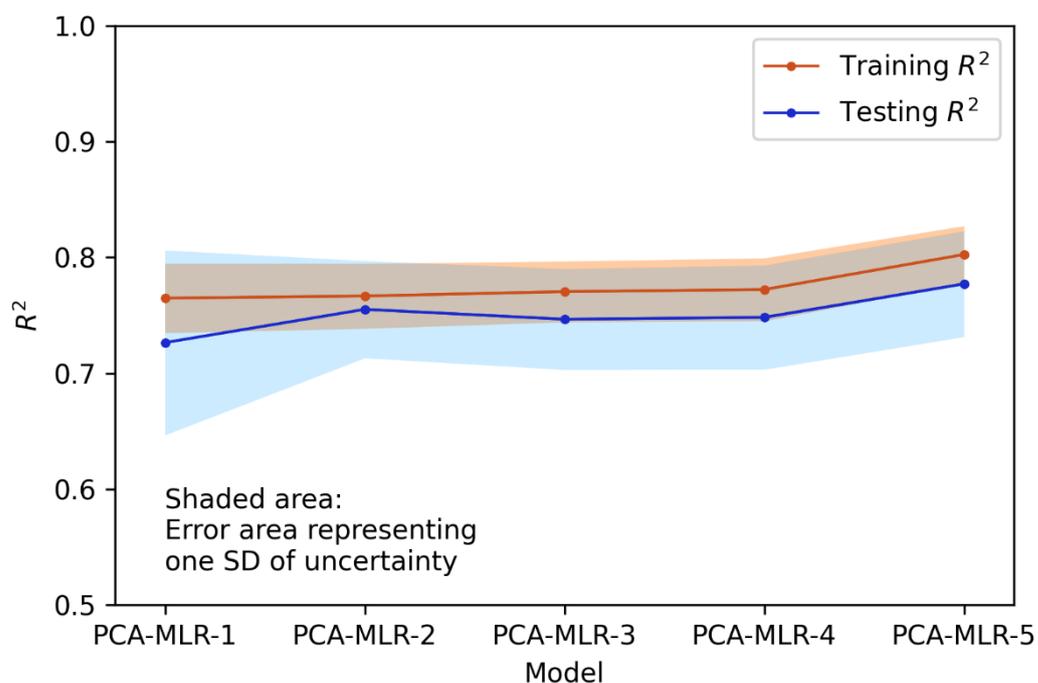
5. Analyses, results and findings

5.1 Multiple linear regression (MLR)

The first t identified components were the inputs of PCA-MLR- t . Accordingly, five models were trained, and their average performance results are shown in Fig. 1. The PCA-MLR-1 only employed Component 1 and obtained the training R^2 of 0.765, meaning that 76.5% of

360 the variance can be explained by Component 1. As more components were added to the
 model, the training performance improved slightly. The testing R^2 shows a similar trend to
 the training R^2 . The five components identified by PCA contain 100% information of the data
 set. When modeling, ignoring any one of them can result in information loss. That is the
 reason for the performance improvement as the number of components increases. The best-
 365 performing model, PCA-MLR-5, has a training R^2 of 0.803 and a testing R^2 of 0.777.

The error area, i.e., the shaded area in Fig. 1, illustrates the stability of the model
 performance. The width of the error area is equal to twice the standard deviation (SD), with
 its center at the average value. Models with narrower error areas can deliver more reliable
 370 results. In Fig. 1, the widths of the testing error area are near twice those of the training error
 area. The widths of error areas change not too much for different MLR models, which
 demonstrates these MLR models have almost the same stability.



375 **Fig. 1** Average training and testing performance of MLR models

The PCA-MLR-5 model was trained 50 times based on different data partitioning. Among
 these models, the one with the closest training and testing R^2 to the average R^2 occurred in the
 17th trial. This model has a training R^2 of 0.791 and a testing R^2 of 0.787, shown as Eq. (6):

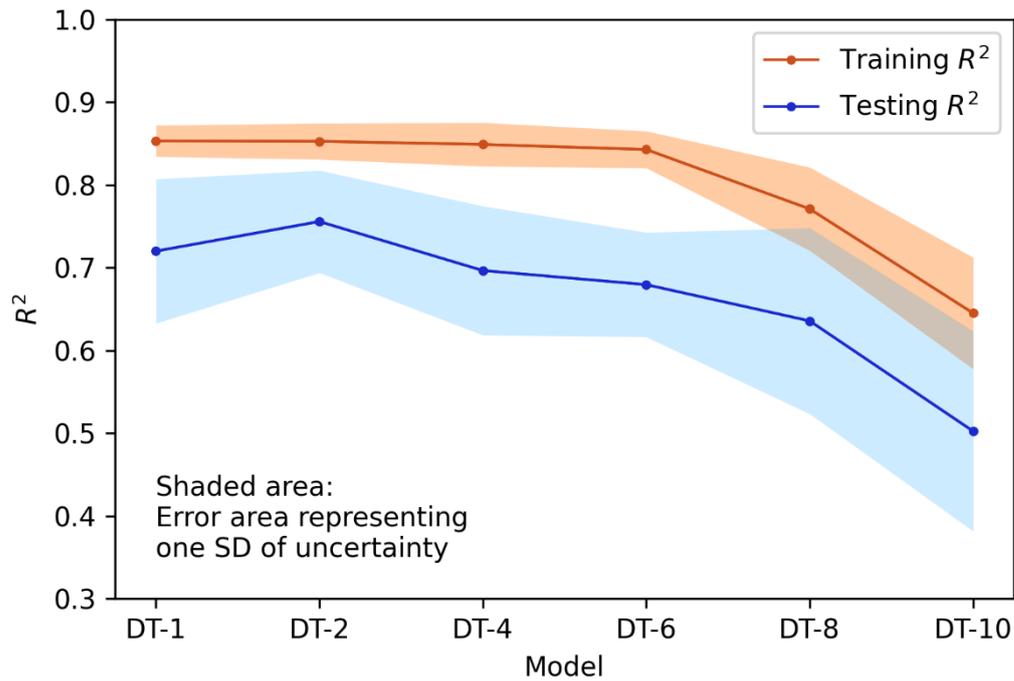
380

$$Y = 3.510 + 0.028 * PO - 0.176 * GC + 0.434 * CO + 0.079 * FS - 0.050 * FC \quad (6)$$

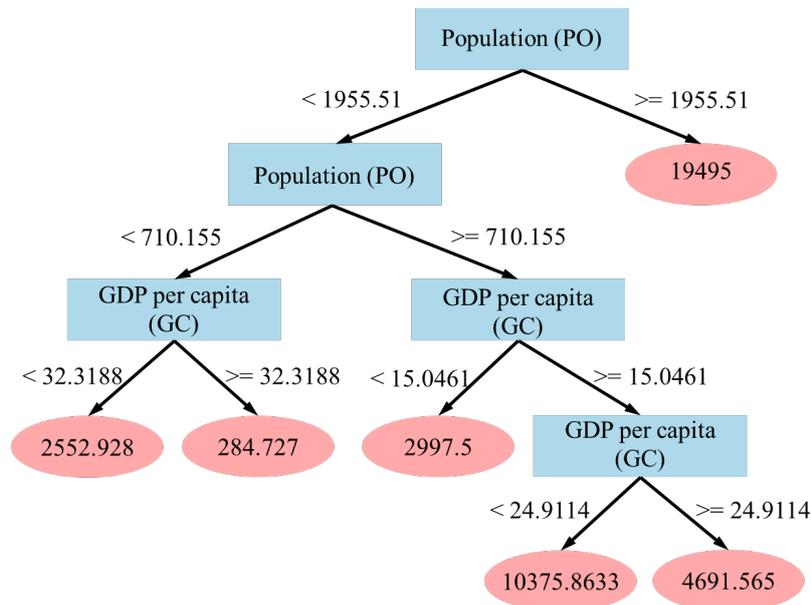
It is worth noting that the model inputs have been normalized to have the same average and variance values, so that the regression coefficients in this model are comparable. It is found that CO has the most significant positive influence, which may be due to the fact that CO directly reflects the level of construction activity. GC has a negative impact on C&D waste generation. In general, high GC means good living conditions for people, sound infrastructural facilities, and high levels of government management. Therefore, GC may contain information about the management level of C&D waste. A higher management level of C&D waste may result in less waste generation. PO, FS, and FC have relatively small coefficients. The regression coefficient of FC is negative, which is probably because the annual FC data fluctuate greatly for some cities.

5.2 Decision Tree (DT)

The complexity of DT was controlled by the minimum leaf size. In general, the bigger the leaf size, the simpler the model. The results under different minimum leaf sizes are presented in Fig. 2(a). The training R^2 decreases as the minimum leaf size increases, while the testing R^2 rises slightly first and then goes down significantly. The decrease in training R^2 is because the DT model becomes so simple that it is not able to accurately define the rules existing in the data. The variations of testing R^2 are closely related to the overfitting and underfitting problems. When the minimum leaf size is small, the model is too complicated to generalize the trained model to the testing data, and the problem is overfitting. When the minimum leaf size is too large, the model is simple and not fully developed, and the problem is underfitting. The optimal model is DT-2, i.e., with a minimum leaf size of 2. Its training and testing R^2 are 0.853 and 0.756, respectively. When the model is simple, the SD is at a high level because the simple model cannot handle these data. The error area is narrowed with more complicated models. Largely, the SDs of DT models are similar to those of MLR models.



(a) Average training and testing performance of DT models



(b) Regression decision tree

Fig. 2 The performance of DT models and the selected decision tree

415

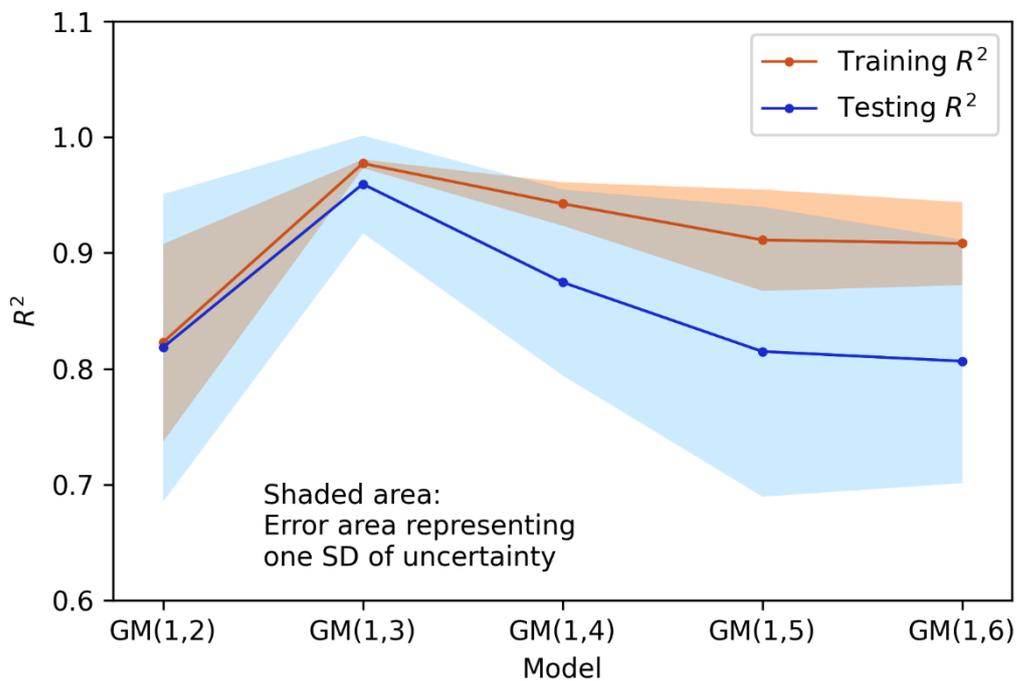
For DT-2, the DT model with the closest training and testing R^2 to the average training and testing R^2 is selected among 50 trials, in order to explore the tree structure. The selected decision tree appeared on the 19th trial, shown in Fig. 2(b). It has a training R^2 of 0.854 and a testing R^2 of 0.764. Other DT-2 models that have similar R^2 present the same tree structure

420 with only the regression value at the leaf nodes different. Thus, it is reasonable to perceive
 425 the regression process by this model.

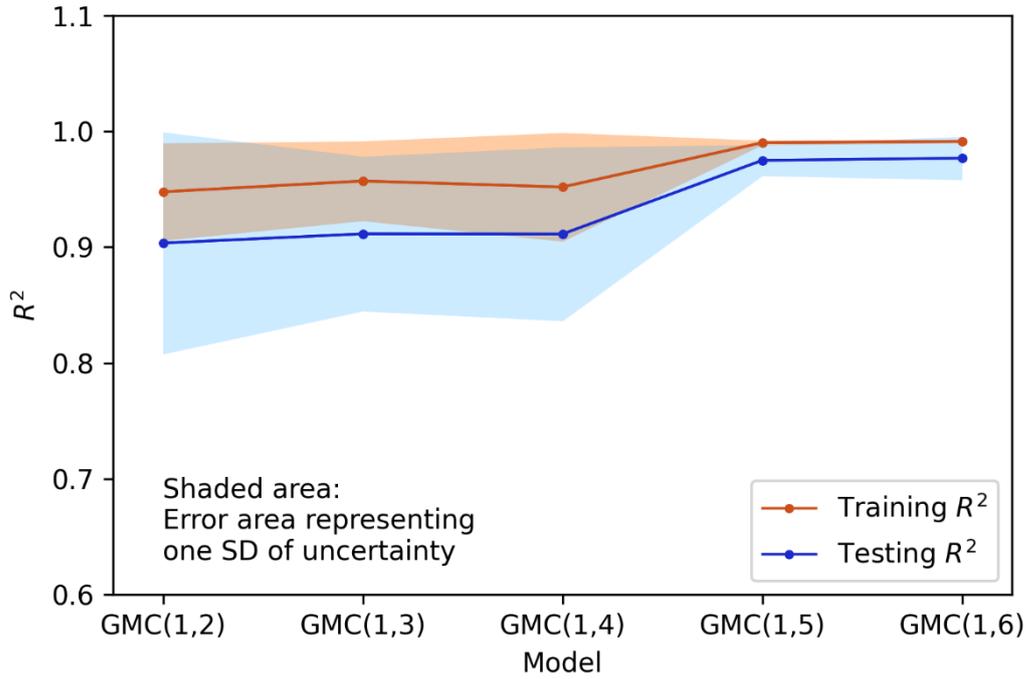
5.3 Grey model (GM)

The GRA results show the FS (0.898), CO (0.865), and FC (0.863) have almost the same
 425 grey relational grades. PO (0.767) is the fourth input factor, followed by GC (0.590).

With the ranked input factors from GRA, 10 grey models were built based on different input
 factors. The training and testing results are shown in Fig. 3. For GM(1,*n*), the best fitting
 model is GM(1,3), with a training R^2 of 0.977 and a testing R^2 of 0.959, involving FS and
 430 CO. With the number of factors increasing, the model performance is not stable. For
 GMC(1,*n*), the grey models perform more steadily as more input factors are fed into the
 models. The best GMC model is GMC(1,6), with a training R^2 of 0.991 and a testing R^2 of
 0.977. Moreover, the error area also contains important information. Most of the GM(1,*n*)
 435 models have high SDs, indicating these models are unstable. However, the case of GMC(1,*n*)
 is different. For GMC(1,5) and GMC(1,6), the widths of the training error area almost shorten
 to perfectly zero, and the testing SDs are also lowered to an acceptable level.



(a) GM(1,*n*)



(b) GMC(1,n)

Fig. 3 Average training and testing performance of grey models

Among 50 trials of GMC(1,6), the model with the closest training and testing R^2 to the
 445 aforementioned average R^2 was found in the 18th trial, with a training R^2 of 0.993 and a
 testing R^2 of 0.967. In Eq. (4), the model parameter is identified: $a = -0.0694$, $b =$
 $[-0.0461, -0.3489, 0.4774, -0.1028, 17.6166]$, $u = -2079.5973$. With these model
 parameters, the selected model is determined. This model has great fitting ability but poor
 interpretive ability.

5.4 Artificial neural network (ANN)

ANN modeling needs a dataset for validation, which can be regarded as a training process. In
 each training epoch, the validation data measures the model's generalization ability. The
 training process is terminated when the generalization no longer improves. The previous 80%
 455 of data was divided into 70% for training and 10% for validation. Six ANN models were
 trained, and the results are shown in Fig. 4. The training R^2 goes up with the number of
 neurons in the hidden layer growing. As shown in Fig. 4, the testing R^2 first increases and
 then decreases. The increase in testing R^2 is due to the enhanced fitting ability as the under-
 fitting model becomes more complicated. The decrease in the testing R^2 means poor
 460 generalization to test data and unreliable predictive performance.

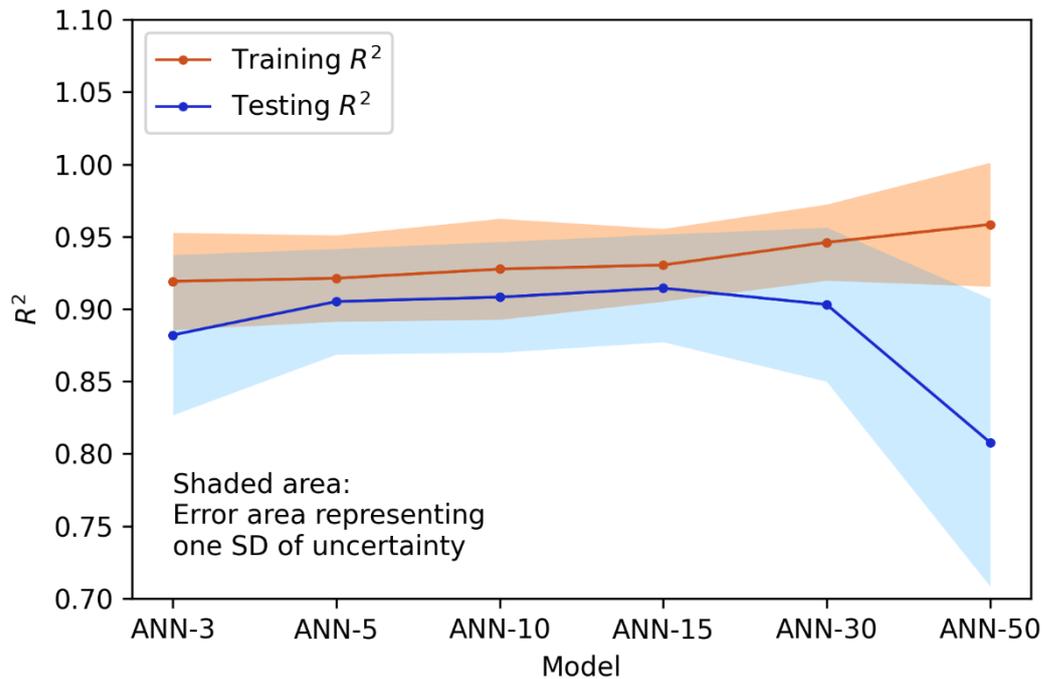


Fig. 4 Average training and testing performance of ANN models

465

In this study, the best-performing model is ANN-15, with 15 neurons in the hidden layer. It has a training R^2 of 0.930 and a testing R^2 of 0.914. Both underfitting and overfitting models have larger SDs. The best-performing model is relatively stable in terms of prediction performance. Among 50 trials of ANN-15, the model with the closest training and testing R^2 to the aforementioned average R^2 occurred in the 35th trial. This model has a training R^2 of 0.925 and a testing R^2 of 0.918. However, this model has many nested sub-structures, making it difficult to trace the influence of each input factor. The meaning of each parameter in neurons is elusive due to the high complexity of the model architecture.

470

475 **5.5 Summary of models**

The best-performing models for each modeling method are PCA-MLR-5, DT-2, GMC(1,6), and ANN-15. The predicted results are shown in Fig. 5, with all the points evenly distributed on both sides of the 45-degree line. The GMC(1,6) has the best performance with the highest training and testing R^2 , followed by ANN-15. These two models are well known for their strong fitting ability. The DT-2 and PCA-MLR-5 models rank third and fourth, respectively. Both of these two models have strong interpretability. The MLR model can tell the major predictors. However, an MLR model can only depict the linear part of a system, which is why it cannot achieve high accuracy. The DT model presents clear logical rules in a tree-based manner, understandable to a person without knowledge of mathematics and statistics. To some extent, the MLR and DT models sacrifice their ability to fit but gain stronger interpretability as compensation. The four best-performing models were also adopted to

480

485

forecast C&D waste generation in GBA cities in 2018. The forecast results show that 11 cities in the GBA produced about 364 million m³ of C&D waste in 2018. The details about the predicted and forecast results of the four best-performing models can be found in the Supplementary Materials.

490

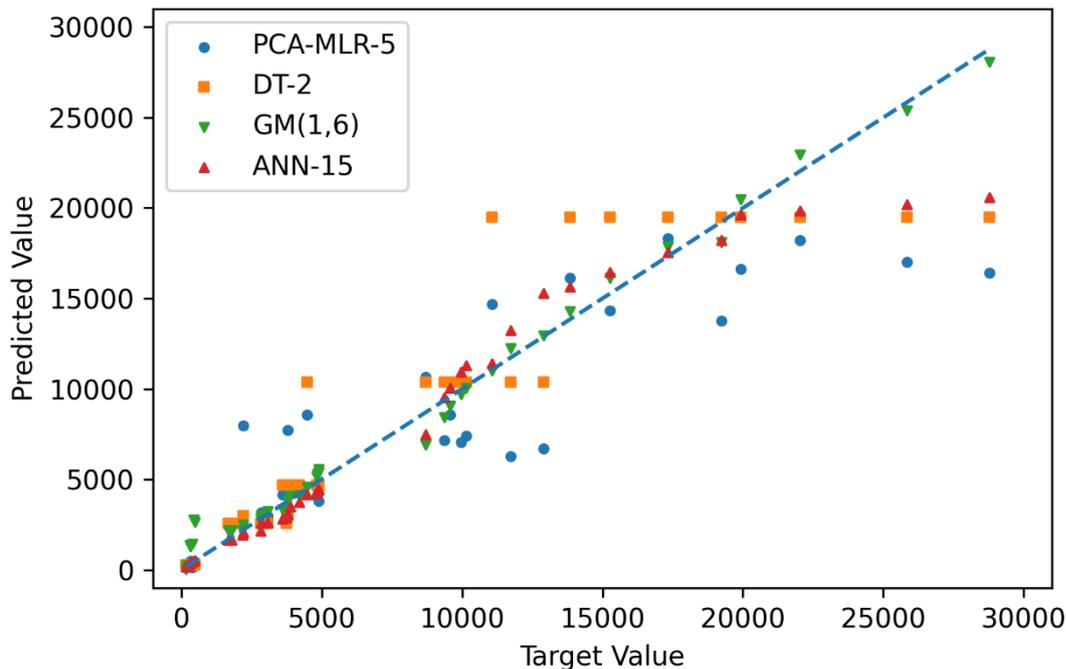


Fig. 5 The predicted results of the four best-performing models

495 **6. Discussion**

6.1 Methodological contributions

Numerous algorithms and models are emerging for estimating waste generation. In this study, four popular and representative ML models, namely MLR, GM, ANN, and DT, were selected and built to estimate C&D waste generation in the GBA, China. Compared with the relevant studies summarized in Table 2, all models in this study are able to deliver acceptable results, implying that the selected factors can explain most variations in C&D waste generation in a region. GM has the highest R^2 . GM is a series of differential equations by nature. The meanings of parameters in such equations are sometimes overly abstruse. Existing studies tend to report its applications in forecasting solid waste and emphasize its high accuracy. Little literature has explained the meaning of the parameters. Likewise, the results recorded in this study do not provide much explanative information other than accurate predictive value.

500

505

ANN, as a powerful ML technique, has also been widely adopted in estimating solid waste generation. ANN is capable of modeling arbitrarily complex nonlinear relations between

510

inputs and outputs, as long as there are enough neurons in the hidden layer. More neurons in the hidden layer make the model entail more model parameters and lead to a huge and complex architecture. These complicated models can also fall into the overfitting trap. In this study, the best-performing ANN models had almost equal performance. Therefore, 515 overfitting is less an issue in this study. Nevertheless, there is no guarantee that this model will perform well on other data sets. Moreover, it is hard to decipher the mechanism behind the models to transfer inputs to outputs. Like GM, the results from ANN are just the predictive value of C&D waste generation without explaining why these values are obtained.

520 Compared with nonlinear models as mentioned above, an apparent benefit of linear models is that the model parameters have their practical meanings. In this study, the MLR model illustrates that factors having a great impact include CO and GC, which may demonstrate the level of construction activity and construction waste management, respectively.

525 DT provides R^2 results of more than 0.75 for both training and testing, which is within the acceptable range. Two factors, namely PO and GC, were used in building the decision tree. Although there is no denying that these two factors do matter, some of the information (e.g., different levels of construction activity) may be neglected. The DT model only produces six predictive values, which may have discrepancies with reality. Nevertheless, it still can give a 530 rough but reliable estimation for reference.

Some of the ML models can achieve high accuracy by developing very intricate models with strong fitting ability. Such models do have significance in research, but may experience the problems of overfitting in practice. One solution is to train the model with more data. It is 535 obviously not feasible in forecasting C&D waste generation in regions with poor statistics in presence. However, encouragingly, some well-managed cities have started collecting data about C&D waste generation in recent years. Another solution is to avoid complex models. When it is impossible to incorporate all influencing factors to produce a deterministic model, there is a wisdom to “think less”, especially in the case of insufficient data. A simple model 540 may be more robust, reliable and interpretative. Therefore, this study calls for paying the same attention to simple and indicative models as complex and accurate models.

To summarize, each ML model has its own strength and weakness (see Table 3). Among those considered, GM and ANN results are more accurate, while MLR and DT contain more 545 understandable information. The better solution is to look at them more comprehensively. It is vital to not only try those models with high accuracy, but also employ interpretative models when estimating C&D waste generation.

Table 3 The strength and weakness of prediction models in this study

Prediction model	Strength	Weakness
Multiple linear regression (MLR)	Strong interpretability; Simple implementation with lower time complexity.	Inability to describe the nonlinear part of datasets.
Decision tree (DT)	Strong interpretability; Ability to model the nonlinear relationship.	Prone to fall into the local optima; Prone to overfitting.
Grey model (GM)	Possible to produce accurate prediction results; Only require small datasets.	Poor interpretability.
Artificial neural network (ANN)	Strong prediction performance; Ability to modeling various complex nonlinear relationships.	Poor interpretability; Prone to fall into the local optima; Prone to overfitting.

6.2 Practical implications

The findings of this study mainly have several practical implications for researchers, policymakers, or environmental protection groups. Firstly, the information, e.g., the estimated C&D waste generation in a region, can be used for examining urban metabolism with a view to developing a circular economy. Urban metabolism is widely applied to describe how material, food, energy, and water consumed by urban as an eco-system to support its growth and reproduce, and consequently generate products and by-products (e.g., GHG, pollutants, and waste) (Kennedy et al., 2007; Wolman, 1965). The amount of C&D waste generation is an indispensable parameter to understand the urban, in particular industrial eco-system metabolism (Zhang et al., 2018). It is also a useful indicator to understand the efficiency of a circular economy system (MacArthur, 2013), which aims to turn some of the waste materials into more circular uses.

Secondly, the estimated C&D waste generation amount can be used for a series of evidence-based policy-making. For example, it can be used for planning the waste management capacity in a region, e.g., the landfill space, the existing and expected 3R capacities. Planners often face the problem of a lack of data when performing this practice. Based on the magnitude of the problem and waste management capacity, policymakers can further make proper arrangements on incentives for recyclers and penalties for polluters. The incentives, including subsidies, tax deduction, and low-cost land usage, have been adopted previously to help recyclers bolster profitability (Bao et al., 2019). Penalties, such as CWDCS in Hong Kong, could impel polluters to minimize C&D waste generation (Lu & Tam, 2013). The information can also be used for inter-regional coordination. For example, the boundary of an urban metabolism system is extended to several regions owing to the globalization of construction resources. Under this circumstance, policymakers are exploring extended producer responsibility (Xu et al., 2021) or cross-jurisdiction waste material sharing (Lu et al., 2020). The reliable estimation of waste generation in this study will provide one of the most important information pieces for these policy-making efforts.

Last but not least, the amount of waste generation can be used for a series of public engagement activities. For example, by presenting the capacity of recycling and landfill and the generation of C&D waste, the urgency of the problem can be better sensed by the general public. As a result, it may better urge stakeholders to avoid the Not-In-My-Back-Yard mindset (Bao et al., 2021), and to more consciously engage in pursuing a circular economy (Ruiz et al., 2020). Performed periodically, this estimate will provide a longitudinal data set, which shows the trend of the CWM performance, hopefully, will allow people to achieve a virtual circle between built environment development and natural environment protection.

7. Conclusion

Data on waste generation at a regional level is of paramount importance to devising proper waste management strategies, but many regions, in particular emerging ones, lack reliable data of this kind. Focusing on the Greater Bay Area (GBA) in China, one of the most economically dynamic areas in the world, this study estimates construction waste generation using limited, publicly available data and proper data analytics. The five factors of population, GDP per capita, total construction output, floor space of newly started buildings, and floor space of buildings completed were adopted. The data analysis results show that these factors can explain most of the variations of C&D waste generation and the coefficients of determination (R^2) reach the level of 0.75 or above. Construction waste generation in individual regions of the GBA can be estimated. These are useful data for developing waste management strategies, for example, monitoring the urban metabolism of input (e.g., materials, energy) and output (e.g., waste), quantifying carbon emission and impacts on climate changes, planning waste management facilities (e.g., recycling plants or landfills), promoting cross-jurisdictional waste material sharing, and so on. This method of estimating construction waste is a useful reference for other regions considering their own dilemma over development and environment.

This study also contributes to the methodology for estimating waste generation. Four types of popular and powerful ML models, namely multiple linear regression (MLR), decision tree (DT), Grey models (GM), and artificial neural network (ANN) were selected and compared by their strengths and weaknesses. The four models all achieved high accurate predictions of waste generation, as evidenced in the high R^2 . Amongst them, GM and ANN have higher prediction accuracy but are more like “black boxes”, not readily accessible to readers. One should also avoid overfitting issues when using the models. In contrast, MLR and DT have slightly lower prediction accuracy but allow more information understandable to readers about the predicting mechanism.

This study has its share of limitations. Firstly, it is based on limited data points, regardless of the best efforts paid to data collection. More model calibration works are expected in the

future using other methodological approaches (e.g., Geographic Information System) when more data is available. Secondly, although it is legitimate to use data-rich regions to extrapolate data-scarce ones, individual features of each region (e.g., Hong Kong's long leading role in waste management; Shenzhen's ambitious zero waste initiative) are yet to be considered in the estimation. Thirdly, estimating future waste generation based on present data is inherently inaccurate. Hence, researchers should adopt a dynamic perspective, monitor the modeling effects, and adjust if necessary. Finally, the biggest motivation of such estimation works is to apply the results in real life. Future studies are encouraged to implement this study in the GBA and receive further verification.

Acknowledgment

This research is supported by the Strategic Public Policy Research Funding Schemes (Project Number: S2018.A8.010.18S) from the Hong Kong SAR Government.

References

- Abbasi, M., & El Hanandeh, A. (2016). Forecasting municipal solid waste generation using artificial intelligence modelling approaches. *Waste management*, 56, 13-22.
- Abdulredha, M., Al Khaddar, R., Jordan, D., Kot, P., Abdulridha, A., & Hashim, K. (2018). Estimating solid waste generation by hospitality industry during major festivals: A quantification model based on multiple regression. *Waste management*, 77, 388-400.
- Afon, A. O., & Okewole, A. (2007). Estimating the quantity of solid waste generation in Oyo, Nigeria. *Waste management & research*, 25(4), 371-379.
- Ahmed, R. R., & Zhang, X. (2021). Multi-stage network-based two-type cost minimization for the reverse logistics management of inert construction waste. *Waste management*, 120, 805-819.
- Alpaydin, E. (2020). *Introduction to machine learning*. MIT press.
- Azadi, S., & Karimi-Jashni, A. (2016). Verifying the performance of artificial neural network and multiple linear regression in predicting the mean seasonal municipal solid waste generation rate: A case study of Fars province, Iran. *Waste management*, 48, 14-23.
- Bakchan, A., & Faust, K. M. (2019). Construction waste generation estimates of institutional building projects: Leveraging waste hauling tickets. *Waste management*, 87, 301-312.
- Bao, Z., Lee, W. M., & Lu, W. (2020). Implementing on-site construction waste recycling in Hong Kong: Barriers and facilitators. *Science of The Total Environment*, 747, 141091.
- Bao, Z., & Lu, W. (2020). Developing efficient circularity for construction and demolition waste management in fast emerging economies: Lessons learned from Shenzhen, China. *Science of The Total Environment*, 138264.
- Bao, Z., Lu, W., Chi, B., Yuan, H., & Hao, J. (2019). Procurement innovation for a circular economy of construction and demolition waste: Lessons learnt from Suzhou, China. *Waste Management*, 99, 12-21.
- Bao, Z., Lu, W., & Hao, J. (2021). Tackling the "last mile" problem in renovation waste management: A case study in China. *Science of The Total Environment*, 790, 148261.
- Batinić, B., Vukmirović, S., Vujić, G., Stanisavljević, N., Ubavin, D., & Vukmirović, G. (2011). Using ANN model to determine future waste characteristics in order to achieve specific waste management targets-case study of Serbia. *Journal of Scientific & Industrial Research*, 70, 513-518.

- Beigl, P., Wassermann, G., Schneider, F., & Salhofer, S. (2004, June 14-17). *Forecasting municipal solid waste generation in major European cities* 2nd International Congress on Environmental Modelling and Software, Osnabrück, Germany.
- Benítez, S. O., Lozano-Olvera, G., Morelos, R. A., & de Vega, C. A. (2008). Mathematical modeling to predict residential solid waste generation. *Waste management*, 28, S7-S13.
- 670 Bramer, M. (2007). Avoiding overfitting of decision trees. *Principles of data mining*, 119-134.
- Cheung, E. (2019). *Greater Bay Area: 10 facts to put it in perspective*. South China Morning Post. Retrieved 16/11/2020 from <https://www.scmp.com/native/economy/china-economy/topics/great-powerhouse/article/3002844/greater-bay-area-10-facts-put>
- 675 Chhay, L., Reyad, M. A. H., Suy, R., Islam, M. R., & Mian, M. M. (2018). Municipal solid waste generation in China: influencing factor analysis and multi-model forecasting. *Journal of Material Cycles and Waste Management*, 20(3), 1761-1770.
- CMAB. (2020). *Overview of Greater Bay Area*. CMAB. Retrieved 16/11/2020 from <https://www.bayarea.gov.hk/en/about/overview.html>
- 680 Cochran, K., Townsend, T., Reinhart, D., & Heck, H. (2007). Estimation of regional building-related C&D debris generation and composition: Case study for Florida, US. *Waste management*, 27(7), 921-931.
- Coelho, A., & De Brito, J. (2012). Influence of construction and demolition waste management on the environmental impact of buildings. *Waste management*, 32(3), 532-541.
- 685 Cunningham, P., Carney, J., & Jacob, S. (2000). Stability problems with artificial neural networks and the ensemble solution. *Artificial Intelligence in medicine*, 20(3), 217-225.
- Domingo, N., & Batty, T. (2021). Construction waste modelling for residential construction projects in New Zealand to enhance design outcomes. *Waste management*, 120, 484-493.
- 690 Duka, A. V. (2014). Neural network based inverse kinematics solution for trajectory tracking of a robotic arm. *Procedia Technology*, 12(1), 20-27.
- Duman, G. M., Kongar, E., & Gupta, S. M. (2019). Estimation of electronic waste using optimized multivariate grey models. *Waste management*, 95, 241-249.
- 695 Dyson, B., & Chang, N.-B. (2005). Forecasting municipal solid waste generation in a fast-growing urban region with system dynamics modeling. *Waste management*, 25(7), 669-679.
- Golbaz, S., Nabizadeh, R., & Sajadi, H. S. (2019). Comparative study of predicting hospital solid waste generation using multiple linear regression and artificial intelligence. *Journal of Environmental Health Science and Engineering*, 17(1), 41-51.
- 700 Guerra, B. C., Bakchan, A., Leite, F., & Faust, K. M. (2019). BIM-based automated construction waste estimation algorithms: The case of concrete and drywall waste streams. *Waste management*, 87, 825-832.
- HKEPD. (2015). *What is construction waste?* HKEPD. Retrieved 16/11/2020 from <http://www.epd.gov.hk/epd/misc/cdm/introduction.htm>
- 705 HKEPD. (2019). *Monitoring of solid waste in Hong Kong*. HKEPD. Retrieved 16/11/2020 from <https://www.wastereduction.gov.hk/sites/default/files/msw2018.pdf>
- Hoang, H. N., Ishigaki, T., Kubota, R., Tong, K. T., Nguyen, T. T., Nguyen, G. H., Yamada, M., & Kawamoto, K. (2021). Financial and economic evaluation of construction and demolition waste recycling in Hanoi, Vietnam. *Waste management*, 131, 294-304.
- 710 <https://doi.org/10.1016/j.wasman.2021.06.014>
- Hoang, N. H., Ishigaki, T., Kubota, R., Tong, T. K., Nguyen, T. T., Nguyen, H. G., Yamada, M., & Kawamoto, K. (2020). Waste generation, composition, and handling in building-related construction and demolition in Hanoi, Vietnam. *Waste management*, 117, 32-41.

- 715 Hsiao, T., Huang, Y., Yu, Y., & Wernick, I. (2002). Modeling materials flow of waste concrete from construction and demolition wastes in Taiwan. *Resources Policy*, 28(1-2), 39-47.
- Hsu, L. C., & Wang, C. H. (2009). Forecasting integrated circuit output using multivariate grey model and grey relational analysis. *Expert systems with applications*, 36(2), 1403-1409.
- 720 Hu, R., Chen, K., Chen, W., Wang, Q., & Luo, H. (2021). Estimation of construction waste generation based on an improved on-site measurement and SVM-based prediction model: A case of commercial buildings in China. *Waste management*, 126, 791-799.
- Huang, T., Shi, F., Tanikawa, H., Fei, J., & Han, J. (2013). Materials demand and environmental impact of buildings construction and demolition in China based on dynamic material flow analysis. *Resources, Conservation and Recycling*, 72, 91-101.
- 725 Intharathirat, R., Salam, P. A., Kumar, S., & Untong, A. (2015). Forecasting of municipal solid waste quantity in a developing country using multivariate grey models. *Waste management*, 39, 3-14.
- 730 Jalali, G. Z. M., & Nouri, R. E. (2008). Prediction of municipal solid waste generation by use of artificial neural network: A case study of Mashhad. *International Journal of Environment Research*, 1(2), 13-22.
- Kannangara, M., Dua, R., Ahmadi, L., & Bensebaa, F. (2018). Modeling and prediction of regional municipal solid waste generation and diversion in Canada using machine learning approaches. *Waste management*, 74, 3-15.
- 735 Kennedy, C., Cuddihy, J., & Engel-Yan, J. (2007). The changing metabolism of cities. *Journal of industrial ecology*, 11(2), 43-59.
- Kern, A. P., Amor, L. V., Angulo, S. C., & Montelongo, A. (2018). Factors influencing temporary wood waste generation in high-rise building construction. *Waste management*, 78, 446-455.
- 740 Khajuria, A., Yamamoto, Y., & Morioka, T. (2010). Estimation of municipal solid waste generation and landfill area in Asian developing countries. *Journal of Environmental Biology*, 31(5), 649-654.
- Kofoworola, O. F., & Gheewala, S. H. (2009). Estimation of construction waste generation and management in Thailand. *Waste management*, 29(2), 731-738.
- 745 Kontokosta, C. E., Hong, B., Johnson, N. E., & Starobin, D. (2018). Using machine learning and small area estimation to predict building-level municipal solid waste generation in cities. *Computers, Environment and Urban Systems*, 70, 151-162.
- Lau, H., Whyte, A., & Law, P. (2008). Composition and Characteristics of Construction Waste Generated by Residential Housing Project. *Int. J. Environ. Res*, 2(3), 261-268.
- 750 Li, X., Chertow, M., Guo, S., Johnson, E., & Jiang, D. (2020). Estimating non-hazardous industrial waste generation by sector, location, and year in the United States: A methodological framework and case example of spent foundry sand. *Waste management*, 118, 563-572.
- 755 Liu, G., & Yu, J. (2007). Gray correlation analysis and prediction models of living refuse generation in Shanghai city. *Waste management*, 27(3), 345-351.
- Lu, W. (2019). Big data analytics to identify illegal construction waste dumping: A Hong Kong study. *Resources, Conservation and Recycling*, 141, 264-272.
- 760 Lu, W., Bao, Z., Lee, W. M., Chi, B., & Wang, J. (2021). An analytical framework of “zero waste construction site”: Two case studies of Shenzhen, China. *Waste management*, 121, 343-353.
- Lu, W., Lee, W. M., Bao, Z., Chi, B., & Webster, C. (2020). Cross-jurisdictional construction waste material trading: Learning from the smart grid. *Journal of Cleaner Production*, 277, 123352.

- 765 Lu, W., & Tam, V. W. (2013). Construction waste management policies and their effectiveness in Hong Kong: A longitudinal review. *Renewable and sustainable energy reviews*, 23, 214-223.
- Lu, W., Webster, C., Peng, Y., Chen, X., & Zhang, X. (2017). Estimating and calibrating the amount of building-related construction and demolition waste in urban China. *International Journal of Construction Management*, 17(1), 13-24.
- 770 Lu, W., Yuan, L., & Xue, F. (2021). Investigating the bulk density of construction waste: A big data-driven approach. *Resources, Conservation and Recycling*, 169, 105480.
- Ma, M., Tam, V. W., Le, K. N., & Li, W. (2020). Challenges in current construction and demolition waste recycling: A China study. *Waste management*, 118, 610-625.
- 775 MacArthur, E. (2013). Towards the circular economy. *Journal of industrial ecology*, 2, 23-44.
- Meza, J. K. S., Yepes, D. O., Rodrigo-Illari, J., & Cassiraga, E. (2019). Predictive analysis of urban waste generation for the city of Bogotá, Colombia, through the implementation of decision trees-based machine learning, support vector machines and artificial neural networks. *Heliyon*, 5(11), e02810.
- 780 Milojkovic, J., & Litovski, V. (2008). Comparison of some ANN based forecasting methods implemented on short time series. 9th Symposium on Neural Network Applications in Electrical Engineering, Belgrade, Serbia.
- Noori, R., Abdoli, M., Ghazizade, M. J., & Samieifard, R. (2009). Comparison of neural network and principal component-regression analysis to predict the solid waste generation in Tehran. *Iranian Journal of Public Health*, 74-84.
- 785 Ojha, V. K., Abraham, A., & Snášel, V. (2017). Metaheuristic design of feedforward neural networks: A review of two decades of research. *Engineering Applications of Artificial Intelligence*, 60, 97-116.
- 790 Pallant, J. (2011). *SPSS Survival Manual: A step by step guide to data analysis using SPSS*. Open University Press.
- Patel, V., & Meka, S. (2013). Forecasting of municipal solid waste generation for medium scale towns located in the state of Gujarat, India. *International Journal of Innovative Research in Science, Engineering and Technology*, 2(9), 4707-4716.
- 795 Perlez, J. (2016). *China Cites Negligence as Cause of Landslide That Killed 73*. The New York Times. Retrieved 16/11/2020 from <https://www.nytimes.com/2016/07/17/world/asia/china-cites-negligence-as-cause-of-landslide-that-killed-73.html>
- Poon, C. S., Yu, A. T., & Jaillon, L. (2004). Reducing building waste at construction sites in Hong Kong. *Construction Management and Economics*, 22(5), 461-470.
- 800 Ruiz, L. A. L., Ramón, X. R., & Domingo, S. G. (2020). The circular economy in the construction and demolition waste sector—a review and an integrative model approach. *Journal of cleaner production*, 248, 119238.
- Song, Y., Wang, Y., Liu, F., & Zhang, Y. (2017). Development of a hybrid model to predict construction and demolition waste: China as a case study. *Waste management*, 59, 350-361.
- 805 Song, Z., Li, Y., & Huang, Z. (2015). Analysis and Forecast of Construction Waste Based on ARIMA Model. 5th International Conference on Civil Engineering and Transportation, Guangzhou, China.
- 810 Soni, U., Roy, A., Verma, A., & Jain, V. (2019). Forecasting municipal solid waste generation using artificial intelligence models—a case study in India. *SN Applied Sciences*, 1(2), 162.
- Tam, V. W. Y., & Lu, W. (2016). Construction waste management profiles, practices, and performance: a cross-jurisdictional analysis in four countries. *Sustainability*, 8(2), 190.

- 815 Tayefi, M., Esmaeili, H., Karimian, M. S., Zadeh, A. A., Ebrahimi, M., Safarian, M.,
Nematy, M., Parizadeh, S. M. R., Ferns, G. A., & Ghayour-Mobarhan, M. (2017). The
application of a decision tree to establish the parameters associated with hypertension.
Computer methods and programs in biomedicine, 139, 83-91.
- 820 Thanh, N. P., Matsui, Y., & Fujiwara, T. (2010). Household solid waste generation and
characteristic in a Mekong Delta city, Vietnam. *Journal of Environmental Management*,
91(11), 2307-2321.
- USEPA. (2016). *Advancing Sustainable Materials Management: 2014 Fact Sheet*. USEPA.
Retrieved 16/11/2020 from [https://www.epa.gov/sites/production/files/2016-
11/documents/2014_smmfactsheet_508.pdf](https://www.epa.gov/sites/production/files/2016-11/documents/2014_smmfactsheet_508.pdf)
- 825 Wang, C. Q., Wei, X. D., & Wang, X. L. (2012). Prediction of Municipal Solid Waste
Production in Changchun City Based on Gray Model GM (1, 5). *Applied Mechanics and
Materials*, 178, 799-803.
- 830 Wang, J., Yuan, H., Kang, X., & Lu, W. (2010). Critical success factors for on-site sorting of
construction waste: a China study. *Resources, Conservation and Recycling*, 54(11), 931-
936.
- Wolman, A. (1965). The metabolism of cities. *Scientific American*, 213(3), 178-193.
- Wu, Z., Ann, T., Shen, L., & Liu, G. (2014). Quantifying construction and demolition waste:
An analytical review. *Waste management*, 34(9), 1683-1692.
- 835 Xu, J., Ye, M., Lu, W., Bao, Z., & Webster, C. (2021). A four-quadrant conceptual
framework for analyzing extended producer responsibility in offshore prefabrication
construction. *Journal of Cleaner Production*, 282, 124540.
- Yu, H., & Wilamowski, B. M. (2011). Levenberg-marquardt training. *Industrial electronics
handbook*, 5(12), 1.
- 840 Yuan, A., Wu, C., & Huang, Z. W. (2012). The prediction of the output of municipal solid
waste (MSW) in Nanchong city. *Advanced Materials Research*, 518, 3552-3556.
- Zhang, Y. (2013). The prediction of the generation of municipal solid waste based on grey
combination model. *Advanced Materials Research*, 807, 1479-1482.
- 845 Zhang, Y., Lu, W., Tam, V. W.-Y., & Feng, Y. (2018). From urban metabolism to industrial
ecosystem metabolism: A study of construction in Shanghai from 2004 to 2014. *Journal
of Cleaner Production*, 202, 428-438.
- Zhao, M., Zhao, C., Yu, L., Li, G., Huang, J., Zhu, H., & He, W. (2016). Prediction and
analysis of WEEE in China based on the gray model. *Procedia Environmental Sciences*,
31, 925-934.
- 850 Zhao, W., Ren, H., & Rotter, V. (2011). A system dynamics model for evaluating the
alternative of type in construction and demolition waste recycling center—The case of
Chongqing, China. *Resources, Conservation and Recycling*, 55(11), 933-944.