

1 **Clearer skies in Beijing – revealing the impacts of**
2 **traffic on the modeling of air quality**

3 Yanyan Xu*
4 Department of Civil & Environmental Engineering, MIT
5 77 Mass. Ave., Cambridge, MA 02139
6 Tel: 857-285-0982; Email: yanyanxu@mit.edu

7 Ruiqi Li
8 School of Systems Science, Beijing Normal University
9 No. 19 Xijiekouwai St., Haidian District, Beijing 100875, China
10 Department of Civil & Environmental Engineering, MIT
11 77 Mass. Ave., Cambridge, MA 02139
12 Tel: 857-207-5631; Email: liruiqi@mit.edu

13 Shan Jiang
14 Department of Civil & Environmental Engineering, MIT
15 77 Mass. Ave., Cambridge, MA 02139
16 Tel: 857-654-5066; Email: shanjiang@mit.edu

17 Jiang Zhang
18 School of Systems Science, Beijing Normal University
19 No. 19 Xijiekouwai St., Haidian District, Beijing 100875, China
20 Tel: 011-86-10-58802732; Email: zhangjiang@bnu.edu.cn

21 Marta C. González
22 Department of Civil & Environmental Engineering, MIT
23 Center for Advanced Urbanism, MIT
24 77 Mass. Ave., Cambridge, MA 02139
25 Tel: 617-715-4140; Email: martag@mit.edu

26 * Corresponding author

27 5339 words + 8 figures × 250 words + 0 tables = 7339 words
28 November 16, 2016

1 **ABSTRACT**

2 Urban air pollution imposes major environmental and health risks worldwide, and is expected to
3 become worse in the coming decades as cities expand. Detailed monitoring of urban air quality at
4 high spatial and temporal resolution can help to assess the negative impacts as a first step towards
5 mitigation. Improvement of air quality needs a variety of measures working together, including
6 controlling industrial pollution and mitigating automobile emissions. In contrast to the measurable
7 industrial pollution, in many of the developing countries, the impact and control of automobile
8 emissions on air quality is neither well understood nor well established. Moreover, the automobile
9 emission data sets are difficult to collect. In this paper, we present a data analysis framework
10 to uncover the impact of urban traffic on estimating air quality in different locations within a
11 metropolitan area. To that end, we estimate the traffic surrounding 24 air quality (AQ) monitoring
12 stations in Beijing, combining mobile phone data and road networks with a traffic assignment
13 model. We investigate how the amount of traffic surrounding each station can impact the modeling
14 of air quality index (AQI) observed by the stations. We separately estimate the contribution of
15 traffic information to the modeling of AQI with regression models in the summer and winter.
16 Further, we group the AQ monitoring stations into four classes, and show that in the summer, air
17 pollution in the inner city is generally more severe than that in the suburbs due to urban traffic;
18 while in the winter, air pollution in the south of Beijing surpasses that in the inner city, most likely
19 due to heating using coal.

1 INTRODUCTION

2 With the rapid urbanization and the acceleration of industrialization, today's air pollution has be-
3 come a global threat of human health, especially for the large scale and densely populated cities
4 in developing countries (1, 2). As pointed by the World Healthy Organization (WHO), in 2012
5 around 3.6 million people died – 16% of total global deaths – as a result of ambient air pollution
6 exposure, which makes it the largest environmental risk to the health of human beings. Moreover,
7 exposure to air pollutants is largely beyond the control of individuals and requires action by public
8 authorities at the national, regional and even international levels.

9 It is important to detect pollutants in the air, to explore their sources, and to model their
10 temporal and spatial patterns, in order to make policy recommendations to mitigate their negative
11 impacts. To better predict air quality (AQ), the relationship between the sources and AQ needs to
12 be examined and clarified. The sources of air pollution are usually divided in 4 categories: sta-
13 tionary, such as industries; mobile, such as transportation sources; area, such as agricultural areas,
14 cities, and wood burning fireplaces; and natural, such as dust and wildfires. The first two of them
15 are human related factors and represent research priorities in the literature. Mobile sources include
16 motor vehicles, marine vessels, and air-crafts. Among them, the exhaust emission of motor ve-
17 hicles is one of the primary factors that influence AQ in urban areas (3, 4). Consequently, clear
18 impacts between traffic and AQ may inform environmental policies. To examine the impact of
19 traffic on air pollution, McHugh *et al.* updated an atmospheric dispersion modeling system with
20 a traffic emissions database (5). Several studies measured the impacts of traffic and meteorology
21 on air pollution measuring data near roads (6, 7). While these studies are detailed on the chemical
22 processes, they do not cover the entire city. Using a data analysis perspective, Zheng *et al.* studied
23 the variations of air quality in space and time in the entire Beijing region via machine learning
24 techniques, combining multiple data sources including taxi data, number of facilities, and the road
25 network data (8). In a related work (9), they predicted air quality in each station, informed by his-
26 torical AQ and meteorological data, and weather forecasts without considering traffic conditions.

27 We focus our study in Beijing, which is one of the most congested and polluted cities in
28 China. Improving AQ in Beijing, is a top-priority locally, that has attracted the world's attention
29 in the past few years. These efforts are compromised by the rapid growth of motorization and
30 urbanization (4). Fig. 1 shows the noise-removed values of air quality index (AQI), wind speed,
31 and humidity from April 1, 2014 to May 1, 2015 in Beijing. The figure represents the variation
32 of AQI at 24 AQ monitoring stations within the Sixth Ring Road of Beijing. Higher AQI values
33 indicate worse air quality. Specifically, AQI values in the range of 0-50 are established as good
34 air, 51-100 moderate, 101-150 unhealthy for sensitive groups, 151-200 unhealthy, 201-300 very
35 unhealthy, and 301-500 hazardous. We see that in general the AQIs in Beijing in the observation
36 period are from moderate to unhealthy. However, they are more stable and lower in the summer
37 (May to October) than in the winter (November to March). Also the wind speed and humidity show
38 different patterns in the two seasons. In the present work, we seek to uncover the contributions of
39 traffic to the air pollution modeling in the summer and winter separately.

40 Our work contributes to the types of studies presented in Refs. (8, 9) in three major aspects.
41 First, we establish separate models—one for the winter and one for the summer—to gain better
42 understanding of seasonal effects of AQI. Second, to investigate the different spatial impacts of
43 urban traffic on AQI, we separately model their relationship by station, taking into account a set of
44 online publicly available daily traffic congestion index (TCI) reported by the local transportation
45 committee to reflect realistic daily traffic conditions. Finally, we enrich the on-line TCI with a

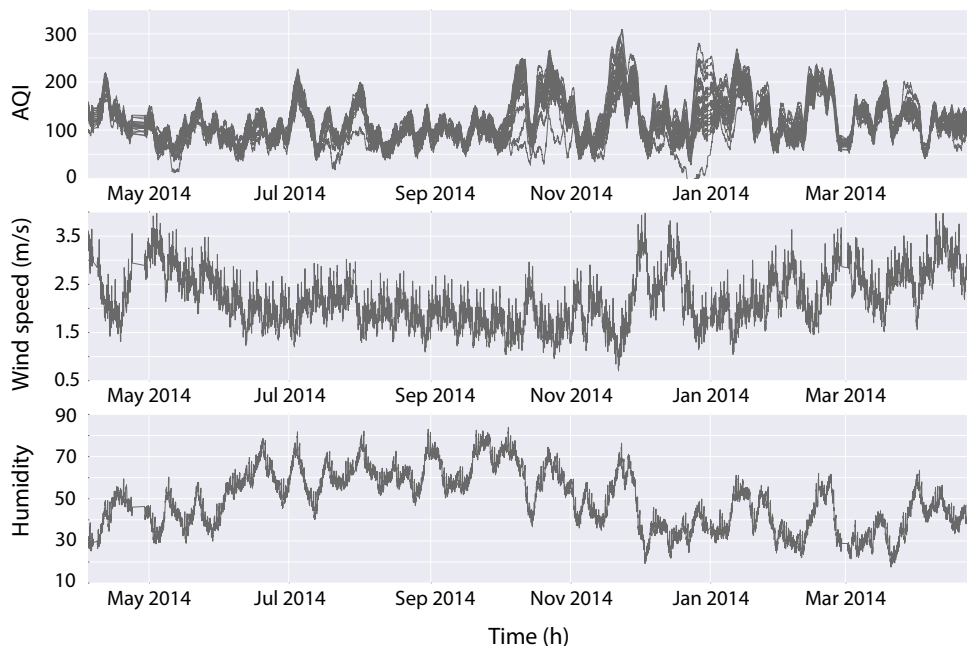


FIGURE 1 Variations of AQI, wind speed, and humidity from April 1, 2014 to May 1, 2015 in each of the 24 air quality monitoring stations within the Sixth Ring Road of Beijing.

1 travel demand model. We calculate the collective travel time (CTT) of all vehicles surrounding
 2 the AQ monitoring station, which is estimated from a mobile phone data based travel demand
 3 model and traffic assignment model integrated with the TCI. Relating traffic with the actual number
 4 of drivers and their origins and destinations is crucial to mitigate congestion in the urban road
 5 network, which can take into account AQ impact.

6 In the next sections, first, we discuss the mobile phone meta data and results from the
 7 call detail records (CDR) to inform a travel demand model. Second, we analyze the importance of
 8 traffic information to the prediction of AQI and the diversity of AQ in space per season. Concluding
 9 remarks and directions for future work are given in the last section.

10 DATA AND METHODOLOGY

11 Travel demand estimation from mobile phone data

12 We estimate the travel demand for the 19.4 million residents living in the urban area of Beijing.
 13 This is commonly referred to the region within the Sixth Ring Road, shown in Fig. 2a, which
 14 has 5.6 million privately-owned vehicles registered in 2013 (10). To our knowledge, our work
 15 constitutes the first traffic estimates of the region based on mobile phone data for Beijing.

16 Alexander *et al.* and Colak *et al.* outlined a general framework to obtain Origin-Destination
 17 (OD) matrices from massive mobile phone data (11, 12). We apply the same methods to extract
 18 trips of users, and estimate the person and vehicle travel demand by combining them with census
 19 data within the Sixth Ring Road of Beijing. Fig. 2a shows the map of Beijing with the AQ stations
 20 marked by blue circles. We focus our study in the inner area marked in darker green.

21 First, we extract stay locations of massive anonymous users from raw mobile phone data,
 22 and labeling activities with *home*, *work* and *other*. Second, we infer number of trips among the

1 stay locations of users by different time of the day and by purpose. Combining with census data,
 2 we expand mobile phone users to total population, and estimate an OD matrix for an average day.
 3 Next (an innovative step proposed by this study), we generate a series of day-specific OD matrices
 4 by using local reported daily traffic congestion index for the city, which allows us to fluctuate the
 5 average daily OD to reflect the realistic daily traffic conditions. We then assign the daily vehicle
 6 demand to the road network.

7 *Mobile phone data*

8 The mobile phone dataset contains 100,000 users with their call detailed records (CDR) and data
 9 detailed records (DDR) for December 2013. Each record of the CDR and DDR data has a hashed
 10 ID, time-stamp, longitude, and latitude of the cell tower when the phone communicated with it.
 11 According to Voronoi tessellation, the average distance between towers is 332 meters (with a me-
 12 dian of 254 meters), representing the spatial resolution in the study. Fig. 2b shows the flow between
 13 tracts for the morning peak (6am-10am), obtained using the mobile phone data as proxy for sur-
 14 veys, with the methods detailed below. Fig. 3a shows the average number of phone usage records
 15 per day that a user has during the whole month. As we see the majority of users are active with an
 16 average of 15 records per day.

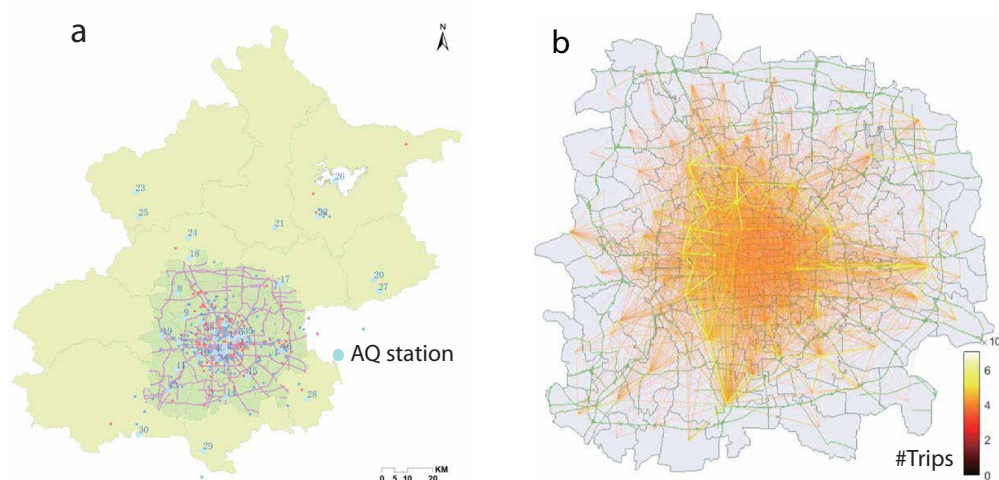


FIGURE 2 Study Area. (a) The boundary of the whole Beijing, it is about $16,410\text{km}^2$. The city area is the greener area, within the Sixth Ring Road, marked by the outer purple line. The blue circles are 35 air quality (AQ) station. (b) Trips between origin destination (OD) pairs in the Morning peak (7am-10am) in urban Beijing.

17 Mobile phone carriers use methods to execute tower-to-tower call balancing to improve
 18 their service. This generates signal jumps that introduce noises, appeared as fast and long move-
 19 ments beyond a travel speed limit. To eliminate this artifact, various methods have been reviewed
 20 in (13). One of the simplest yet effective methods is to remove the next record if the the inferred
 21 speed between two records is beyond reasonable speed limit. However, it heavily relies on the cor-
 22 rectness of the first record. To improve its accuracy, we check if the first record is a noise —if the
 23 speed between the first and the second record is beyond a predefined speed limit, we then remove
 24 the first record. We repeat this process until there is no artificial jumps between two records. Next,
 25 we distinguish stay-point and pass-by from the remaining records.

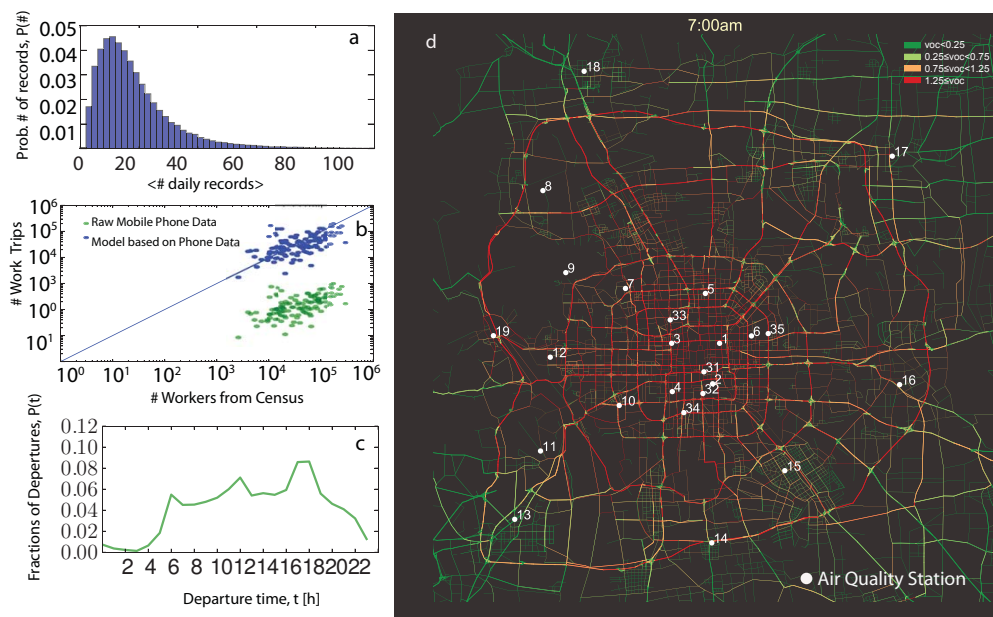


FIGURE 3 Traffic model of Beijing. (a) The distribution average daily records among the 100,000 mobile phone users (b) Validation of the estimated number of work trips vs. employment information from Census data (c) Estimated fraction of trip departures per hour of the day. (d) Estimated volumes of cars in the streets per time of the day.

1 We improve upon the stay-point algorithm presented in (13, 14) as follows. (i) we apply
 2 a temporal agglomeration algorithm. The temporally consecutive records within a certain radius
 3 (e.g., 500 meter) are bundled together with a updated stay duration from the start time of first record
 4 to the end time of last one. (ii) We then label the records as pass-by points and stays, according
 5 to the stay duration threshold (e.g., 10 minutes) based on the local context in Beijing. In analysis
 6 hereafter, we only focus on the stays. We then combine all the spatially adjacent stay points for a
 7 user (within a threshold) as his or her stay regions, which will be later labeled as *home*, *work*, and
 8 *other*. For this spatial agglomeration, we use R-tree to accelerate the computation (15). R-tree is
 9 a type of spatial B-tree, a spatial search balancing tree that checks the boundaries of elements to
 10 make the search faster (see details in Fig. 4). We then get a mapping relation between stay points
 11 and stay regions.

12 *Stay detection and activity labeling*

13 We then estimate the type of each stay location for every user, classified as *home*, *work* or *other*.
 14 The most visited location during weekday nights and weekends are labeled as *home*, and the most
 15 visited one during weekday working hours (at least 500 meters away from home) is labeled as
 16 *work*, and the rest are labeled as *other*. We assume that within 500 meters, it is not necessary to
 17 travel by car.

18 *Vehicle demand estimation*

19 After labeling the activity type, we estimate residential and working population within each zone
 20 (i.e, a Voroni polygon generated from towers), and calculate an expansion factor by dividing the
 21 number of phone users by total population for each zone. We aggregate the population data at

Algorithm 1: Spatial Agglomeration by R-tree (Python)

```

1 import index from rtree;
2 tempStay2Stay = dict();
3 idx = index.Index();
4 for node in tempStays do
5   | idx.inserts(tempStay);
6   | #degenerate the rectangular to a point when inserting the tempStays;
7 for node in tempStays do
8   | VectorState[node] = idx.intersection(node's square buffer);
9   | #search the buffer region of each node to see how many nodes are in its neighbor in the
   | | constructed rtree (i.e., idx above);
10 while the sum of StateVector is not 0 do
11   | choose the node_i with maximum value in StateVector;
12   | intersection = idx.intersectionnode i's square buffer if the len(intersection)==StateVector[i]
   | | then
13   | | #cluster the nodes within node i's buffer, get the mapping relationship to the most central
   | | | one for nodes j within the square buffer;
14   | | for node_j in intersection do
15   | | | tempStay2Stay[node_j] = node_i;
16   | | | StateVector[node_j]=0;
17   | | | idx.delete(node_j);
18   | | else
19   | | | StateVector[i] = len(intersection);
20   | | | continue;
21 return tempStay2Stay;

```

FIGURE 4 Algorithm used for detection of stay regions from mobile phone traces

1 the 100 by 100 meter grid level obtained from WorldPop¹ to the *Jiedao* level (census zones
2 comparable to towns in U.S.). We compared the total population obtained from WorldPop with
3 the Beijing Census data (2010) at the *Jiedao* resolution, and they are in good agreement. We
4 compare the home-work trips generated by our model with the census employment statistics at the
5 *Jiedao* level, only taking into account the phone users with labeled work location. We find that
6 our employment estimation is in reasonable agreement with the Beijing 2nd Economic Census (see
7 Fig. 3b).

8 Trips are then assigned a trip purpose: home-based-work (commuting), home-based-other,
9 and non-home-based, according to the inferred locations of two consecutive stays. We then get an
10 overall average departure time distribution from all the trips normalized by the number of active
11 days, and an expansion factor for each user. Although a travel survey from Beijing is not available
12 to us at the moment, this method has been approved in other cities with their travel surveys (11,
13 12, 16). In Fig. 3c, we show the estimated fraction of trips per hour in an average day.

14 We obtained OD matrices by different time periods of an average weekday according to
15 the departure time at both the Voronoi polygon and census tract level, where the number of trips
16 are expanded by the expansion factors. To consider trips made by motorized vehicles, we weigh
17 obtained person trips by vehicle ownership rates at the district level which is larger than *Jiedao*

¹<http://www.worldpop.org.uk/data/methods/>

1 (e.g, with 18 districts in Beijing). According to the 2013 Beijing Year Book (10), due to local
 2 traffic regulation policy, around 20% of cars are restricted not to travel on the road according to
 3 their car license numbers. We multiply 0.8 by all trips, as each day two license ending-numbers
 4 are restricted by the city. The other factor is the vehicle usage rate— many people who own cars
 5 tend to use subways rather than driving to avoid traffic congestion in peak hours. Consequently,
 6 we assume a factor of 80% for all tracts, and this step is yet to be improved with more accurate
 7 car usage rate data, which is not available at high resolution. Finally, with a traffic assignment
 8 model (17), we assign the vehicle ODs to the road network resulting estimates of travel time and
 9 car volumes for each segment of the road network.

10 *Day-specific travel demand estimation*

11 We extend the average 24-hour demand calculated from mobile phone data to day-specific ODs
 12 using data reported on traffic congestion index (TCI). TCI is published by Beijing Transportation
 13 Research Center (BJTRC) (18) and ranges from 0 to 10. As explained by BJTRC, 0 indicates all
 14 vehicles in the road network traveling in free flow speed; 10 indicates the travelers on average take
 15 double free-flow travel time on the road segments. TCI reflects the degree of congestion, other
 16 than the fraction of travel demand. We use them, however, as a source of information to generate
 17 variations in demand, with the following equation:

$$f_d = \frac{TCI_{max} + TCI_d}{TCI_{max} + TCI_{mean}} \quad (1)$$

18 where f_d is the demand factor on the d th weekday in our data set; TCI_d is the value of TCI on the
 19 d th weekday; TCI_{max} and TCI_{mean} are the maximum and mean TCI of all weekdays, respectively.
 20 As a result, the zone-to-zone OD matrix is scaled with the demand factor f_d for each weekday. In
 21 our experiments, f_d ranges from 0.65 to 1.31. This means that during weekdays from April 2014
 22 to May 2015, we allow for fluctuations in traffic congestion, introducing a degree of uncertainty
 23 in the proposed travel demand estimates, enriched by the variations reported by in the TCI on the
 24 same days over which we will model the AQIs.

25 *Traffic Assignment*

26 To estimate the traffic state and travel time of drivers, we assign the vehicle demand to the road
 27 network using a user equilibrium (UE) model. A UE model assumes that all of the travelers in
 28 the road network try to find their routes with respect to the shortest travel time (19, 20). The
 29 road network of Beijing within the Sixth Ring Road is extracted from OpenStreetMap (21). We
 30 extracted or estimated requisite attributes of road segments, including free flow speed, capacity,
 31 length, and number of lanes, from OpenStreetMap.

32 The road network is represented as a directed acyclic graph (DAG), $\mathcal{G}(\mathcal{N}, \mathcal{E})$, where \mathcal{N} is
 33 the set of all nodes, \mathcal{E} is the set of all edges. In our implementation of the UE model, the anticipated
 34 travel time on each edge e is calculated by the Bureau of Public Roads (BPR) function:

$$t_e = \left(1 + \alpha \left(\frac{v_e}{C_e} \right)^\beta \right) \times t_e^f \quad (2)$$

35 where v_e is the number of vehicles attempting to use edge e per hour; C_e is the capacity of the edge;
 36 t_e^f is the free flow travel time on edge e and is estimated using the limit speed of the edge; α and β
 37 are two coefficients and we are using $\alpha = 0.18$ and $\beta = 4$ in our experiments.

1 To solve the UE model, we minimize the distance between the optimal solution and the
 2 current solution in an iterative process (22, 23). In our work, the distance is measured using the
 3 following equation:

$$r_g = 1 - \frac{\sum_{o \in \mathcal{O}, d \in \mathcal{D}} t'_{od} f_{od}}{\sum_{e \in \mathcal{E}} t_e v_e} \quad (3)$$

4 where \mathcal{O} and \mathcal{D} are the set of origin and destination nodes in the road network; f_{od} is the demand
 5 of flow from o to d ; t'_{od} is the shortest travel time of trip (o, d) in the current iteration. Further
 6 details of the implementation of assignment can be found in (17).

7 Fig. 3d shows the assignment results during morning peak hour. The color of each road
 8 segment reflects the volume-to-capacity (VoC). A larger VoC indicates that the road is used by a
 9 larger number of vehicles compared with its capacity. As seen from the figure, a large proportion
 10 of the urban roads are in congestion during the morning peak.

11 To verify that the assignment results are reliable and robust, we compare the travel time of
 12 5,000 OD pairs with top number of commuters during the morning peak hour with the travel time
 13 provided by Gaode (24), which is a leading traffic navigation company in China. Fig. 5a shows
 14 the comparison of travel times, suggesting that our estimated travel times and Gaode's are quite
 15 close for most of the trips. Fig. 5b presents the distribution of commuting time of the top 5000 OD
 16 pairs. The distribution indicates that our assignment model provides reliable estimates of travel
 17 time delay in the peak hour.

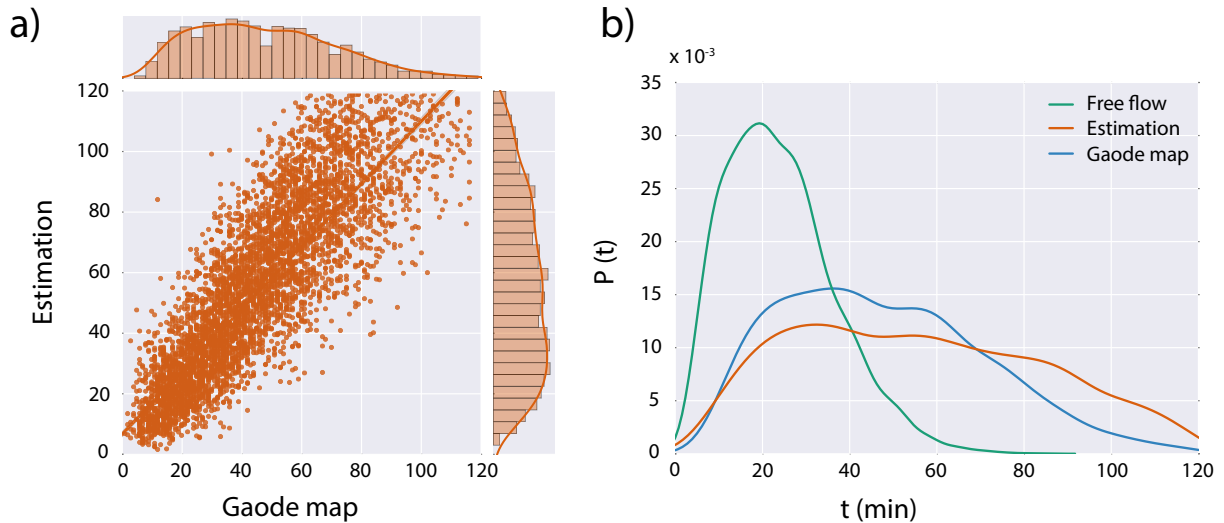


FIGURE 5 Commuting time validation with Gaode travel time. (a) The scatter plot of 5000 trips with top commuters. (b) The distribution of travel time with three modes: free flow, our estimation, and Gaode map.

18 Measuring traffic feature by AQ monitoring station

19 The coverage radius for an AQ monitoring stations in the city ranges from 500 meters to 4 kilo-
 20 meters. We define a $2km \times 2km$ square-buffer surrounding each station to examine the relation
 21 between traffic around the station and its AQ. This enables us to identify stations that are more
 22 sensitive to local traffic. By assigning the day-specific vehicle ODs (extended by the TCIs) to the

1 road network, we estimate vehicle numbers in the streets by hour for different days. We then es-
 2 timate the volume of vehicles and travel times for each road segment for each of the days. Since
 3 the traffic-related air pollution is not only related to the vehicle volumes but also with the time
 4 they spend (to approximate emission) in the road network, we calculate the collective travel time
 5 (CTT) within the buffer area of the AQ monitoring station as a traffic feature to model the AQI.
 6 The collective travel time is calculated as $t_c = \sum_{e \in \mathcal{B}} v_e t_e$, where \mathcal{B} is the set of roads in the buffer
 7 area. Besides, the total VoC is also calculated as the summation of VoC on all roads in the station
 8 buffer area.

9 Fig. 6 shows the CTT and total VoC per hour per station. The CTT and VoC in each station
 10 buffer are obtained by assigning the average demand to the road network. As shown in the figure,
 11 there are three peak hours on weekdays in Beijing. The CTT at four stations: 1, 3, 5, 31 are
 12 significantly higher than others, and three of them also have high VoC. Besides, most stations with
 13 heavy traffic are located within the Fourth Ring Road of Beijing. In the next section we discuss the
 14 use the CTT as a predictive feature for AQI.

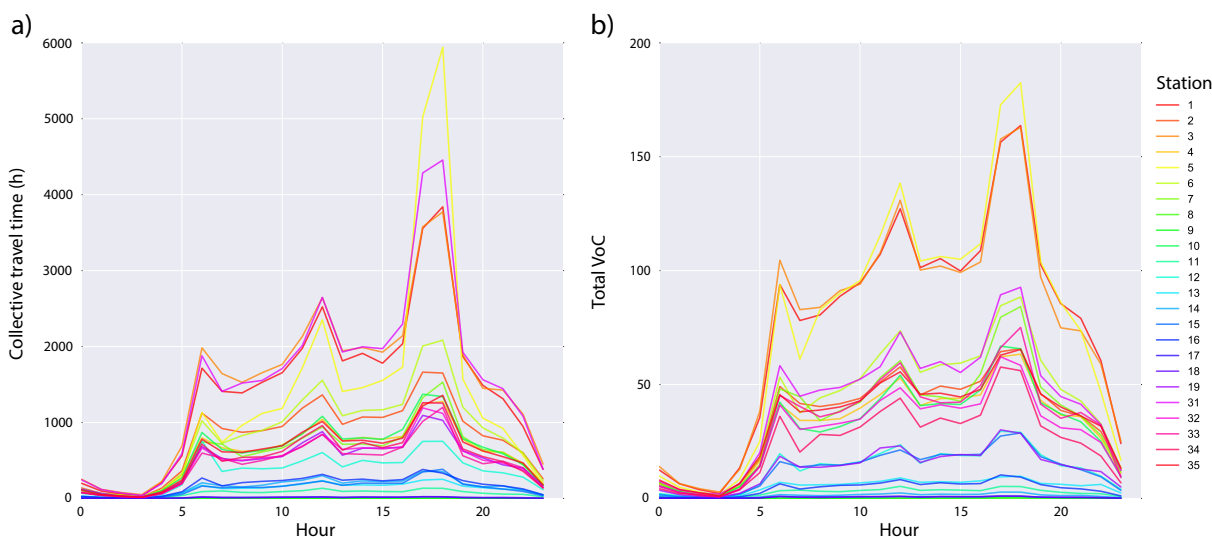


FIGURE 6 Travel demand information per hour in buffer areas of each station (a) The collective travel time. (b) Total car volume over street capacity (VoC).

15 Impact of traffic to modeling of air quality

16 Although traffic is regarded as one of the most critical influences on air pollution in urban areas,
 17 the impact of traffic is still not well measured and understood. Zheng *et al.* predicted the AQI in
 18 Beijing with features related to meteorology, number of taxi trips, road properties, point of interests
 19 (POIs), and traffic related features (e.g., speeds from taxi data) (8). They built a single prediction
 20 model for the entire city. That is, the model was trained using data from all AQ stations in Beijing,
 21 disregarding the spatial variations of AQ. In a later work of the same team, they predicted future
 22 AQ in each station, but without considering the traffic factor in the station (9). We argue that
 23 a city-wide model cannot identify the spatial variations reflecting the importance of local traffic
 24 feature for the AQI by station, which is important in relating AQ with transportation policy. In
 25 this work, we investigate this aspect, modeling the AQI in each of the 24 monitoring integrating a

1 travel demand model. The location and ID of the stations are shown in Fig. 3d. We can see that
 2 some stations are located in zones with heavier traffic than others.

3 To evaluate the impact of traffic on air pollution, we model the AQI using the meteorology
 4 and traffic information in the same hour. The meteorological features include wind speed, wind
 5 direction, humidity, temperature, and pressure. The traffic features include the TCI and proposed
 6 CTT. We divide the data set into two parts: summer (from May 1, 2014 to September 30, 2014) and
 7 winter (from December 1, 2014 to March 31, 2015). For each part, we train a estimation model for
 8 each station under the three aforementioned scenarios, and use the raw AQIs as response. More-
 9 over, as people are more concerned with air quality during daytime, we select the samples from
 10 6:00am to 8:00pm everyday. After eliminating the missing data, the summer data set contains
 11 about 430 sample-hours per station; the winter data set contains about 530 sample-hours per sta-
 12 tion, corresponding to only 31 and 38 days with complete data, respectively. To avoid the overlap
 13 between training and testing sets, the first 70% sample-hours are used to train the models, and the
 14 last 30% hours are used to test. Subsequently, we estimate the AQI with two distinct models, lin-
 15 ear regression and non-linear random forest model (25, 26). To obtain stable estimation, we repeat
 16 training the model 20 times at each station. At each time, we randomly select 90% data from the
 17 training datasets to train the model. The average value of the 20 estimations is regarded as the final
 18 estimation of AQI at the station.

19 RESULTS

20 Analysis of AQI estimation

21 To assess the impact of traffic features on air quality, we first calculate the relative feature impor-
 22 tance of three feature sets, meteorology, TCI, and CTT in two regression models. A linear regres-
 23 sion model, and a random forest. For the linear regression, we use the Lindeman, Merenda and
 24 Gold (LMG) method to quantify the contribution of individual feature sets to modeling AQI (27).
 25 For random forest, the importance of a feature set is calculated through the difference of training
 26 accuracy with and without the feature set. The estimation accuracy of AQI is calculated by:

$$p = \left(1 - \sum_{i=1}^N \frac{|\hat{AQI}_i - AQI_i|}{AQI_i} \right) \times 100\% \quad (4)$$

27 where \hat{AQI}_i is the estimated value of the i th sample; N is the number of samples in the testing set.

28 Fig. 7a and Fig. 7c illustrate the relative feature importance of meteorology, TCI, and CTT
 29 in summer, with linear regression and random forest, respectively. As can be seen, meteorology
 30 is the leading factor at most stations. However, linear regression suggests TCI is less important
 31 than CTT, while random forest suggests TCI and CTT have equal level of contribution to AQI. The
 32 importance of features to AQI estimation in winter are divergent for the two regression models,
 33 as shown in Fig. 7e and Fig. 7g. Such diversity between two models reflects the AQI in winter is
 34 more difficult to model than summer.

35 Fig. 7b and Fig. 7d present the distribution of the estimated accuracy in all stations in the
 36 summer, with the linear regression and random forest, respectively. The red distribution is obtained
 37 with model trained with all features, while the blue one is obtained with model trained without
 38 traffic features (TCI and CTT). Integrating the traffic features with the meteorological features, the
 39 accuracy decreases in some stations. This indicates that the traffic information in a given hour has
 40 not direct impact in the AQI in the same hour. The impact of traffic to air quality may be delayed

- 1 for more than one hour. Similar results were obtained in winter, as shown in Fig. 7f and Fig. 7h.
- 2 From these results, we notice that although the traffic information has significant importance in the
- 3 training phase of regression models, it can not promote the estimation of AQI in the testing phase.

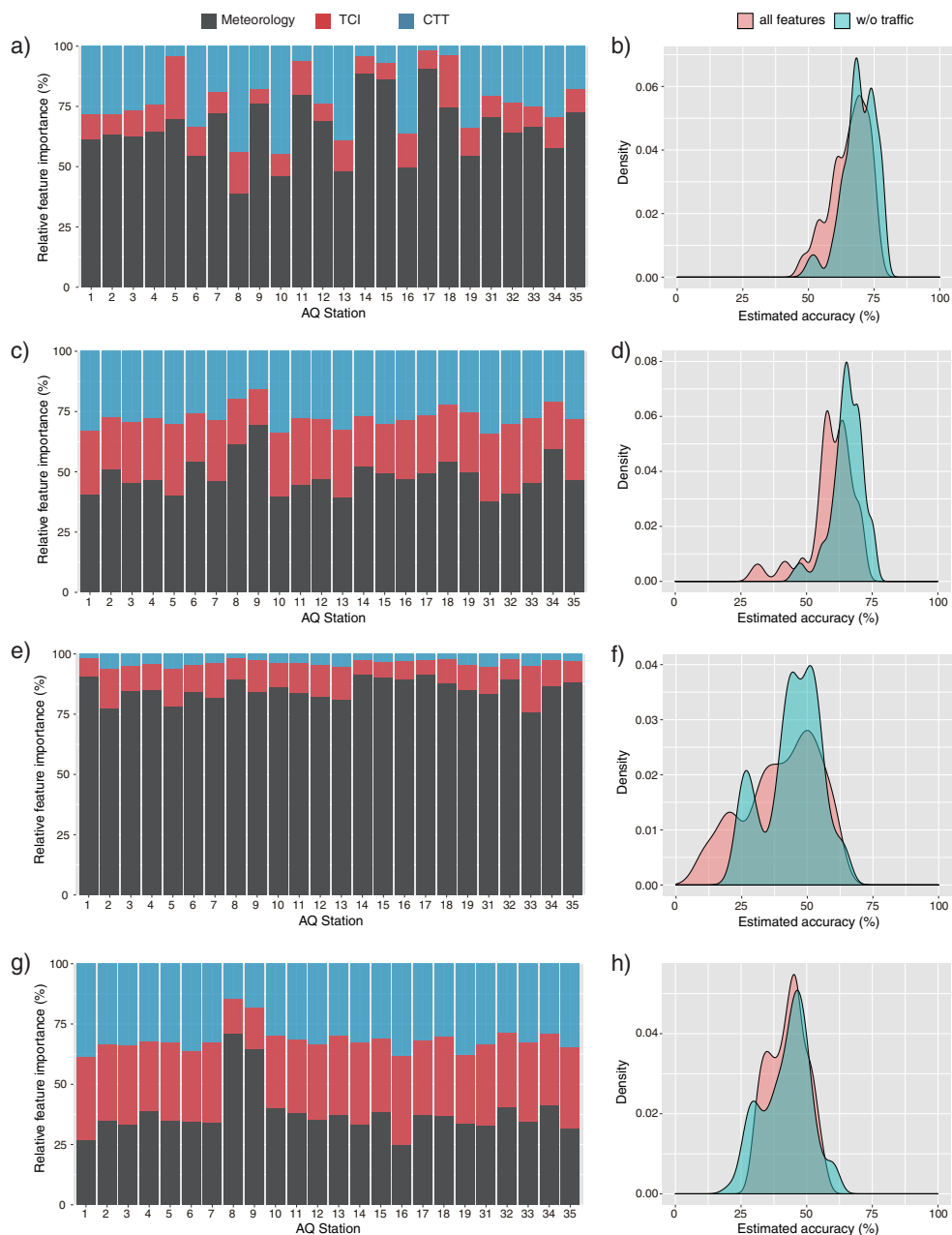


FIGURE 7 Feature importance and AQI modeling accuracy. (a-b) Relative feature importance per AQ station and the distribution of estimated accuracy of all AQ stations with linear regression in the summer. (c-d) Relative feature importance and the distribution of estimated accuracy with random forest in the summer. (e-h) Results in the winter, e and f are results of linear regression, g and h are results of random forest.

1 **Spatial diversity of AQ monitoring stations**

2 We further analyze the different relationship between AQI and traffic demand information among
3 the 24 stations. In Fig. 8a and 8c, we plot the median value of AQI and CTT at each station in
4 summer and winter, respectively. As shown in the figures, the CTT at the 24 stations are distinctly
5 separated in two groups: heavy and light traffic stations. Heavy traffic stations are located in the
6 inner urban area, while lighter traffic stations are located in the suburbs. To divide the AQIs, we
7 use 100 as a threshold—according to the U.S. Environmental Protection Agency, a AQI higher
8 than 100 is regarded as unhealthy.

9 Finally, we partition the 24 stations into four groups: healthy with light traffic, healthy with
10 heavy traffic, unhealthy with light traffic, and unhealthy with heavy traffic, shown in Fig. 8a-d with
11 different colors. As seen from results in Fig. 8a and 8b, the median AQIs of all light traffic AQ
12 stations (green) are under 100, which indicates that around these stations, the air quality on most
13 days in the summer are healthy. For the stations with heavy traffic, only two of them (station 12
14 and 34) are unhealthy. Station 12 is located at the West Fifth Ring Road; station 34 is located at
15 the South Third Ring Road; and both of them suffer with busy traffic. Fig. 8c and 8d show the
16 results in the winter. In general, the air pollution in winter is much severer than that in the summer.
17 Consequently, AQ at some stations (e.g., station 10, 15, 31, 32 and 35) change from healthy in
18 the summer to unhealthy in the winter, while only stations 12 and 5 improve their AQIs in the
19 winter. Interestingly, these stations are all located in the southern area of Beijing. Meanwhile,
20 from the map of major coal power plants in and around Beijing in Fig. 8e, we observe there are
21 some large-capacity power plants at the south-eastern area of Beijing, e.g. Hebei province and
22 Tianjin. This argument has been demonstrated in literature (28): in the winter, the air pollution in
23 the north China is more critical than the south because of the burning of coal for heating. On the
24 other side, the traffic is heavier in the inner core for both winter and summer. Therefore, we argue
25 that the degraded air quality in the southern area of Beijing reaching the unhealthy limits, is likely
26 not related to traffic but due to heating by coal sources.

27 **CONCLUSION**

28 In this paper, we studied the contribution of traffic related features to the air quality index in the
29 same hour in 24 monitoring stations in Beijing. We integrated mobile phone meta data and publicly
30 available daily traffic congestion index (TCI) to define the traffic features. First, we estimate zone-
31 to-zone vehicle travel using mobile phone data, census data, vehicle usage rate, and road network
32 information. Second, we generate day-specific hourly ODs using TCIs. The day-specific ODs
33 are then assigned to the road network, and the maximum collective travel time (CTT) surrounding
34 each AQ station area is estimated per day in the studied period. Based on the meteorological
35 data, the TCI, and our estimates of CTT, we built two regression models for each station in the
36 summer and the winter. The results show that the traffic information has significant importance in
37 the training phase of the regression model. However, it cannot promote the estimation accuracy in
38 the testing phase. The main reasons may be: (i) the air pollution generated by automobile can not
39 be reflected by AQI immediately; (ii) the regression models do not capture the relations between
40 traffic features and AQI effectively due to the limited period of observation and sample size of to
41 generate the travel demand model.

42 Moreover, to relate the impact of traffic on air quality in space, we categorize the 24 stations
43 within Sixth Ring Road of Beijing into four groups. We find that the stations with heavy traffic
44 are in the inner core of the city both in winter and summer. The stations with unhealthy levels of

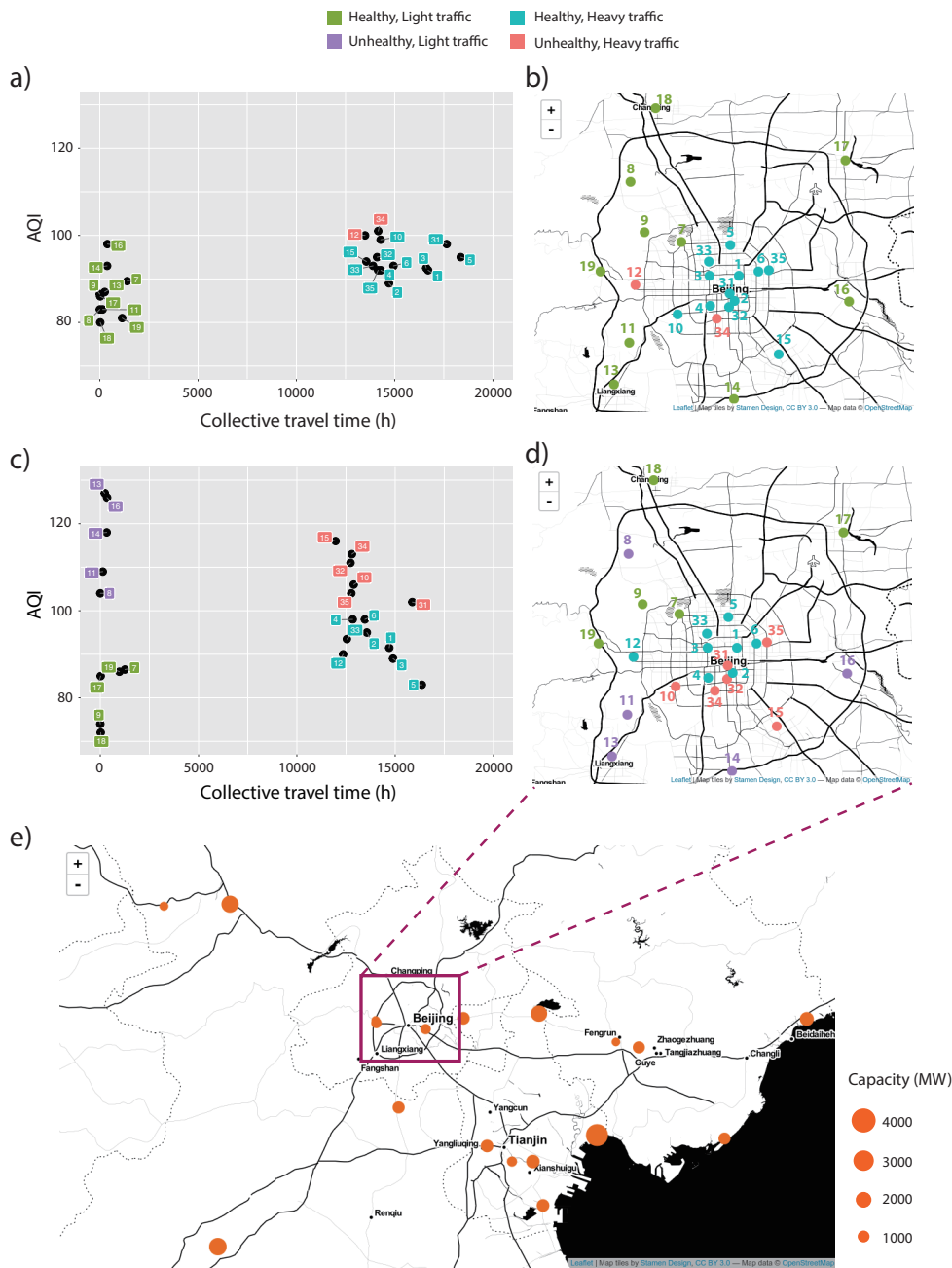


FIGURE 8 The spatial separation of stations according to AQI and CTT. (a, b) Stations separation results in summer. (c, d) Stations separation results in winter. (e) The location of major coal power plants in and around Beijing.

1 air pollution appear in the winter and are concentrated in the southern area of Beijing. Based on
 2 these observations, it suggests that the coal heating rather than traffic contributes significantly to
 3 the degraded air quality in south Beijing in the winter.

4 The presented framework is portable, as the data sets employed here can be easily obtained
 5 for other cities. The traffic estimation model is of low cost in computation and data require-

1 ments. This work also provides a data pipeline to categorize AQ monitoring stations more affected
2 by traffic congestion, and to estimate AQIs based on meteorology data, traffic congestion index,
3 and travel demand estimates from mobile phone meta data. There are important avenues for fu-
4 ture work, these include: (i) to further investigate the variation of specific pollutants such as NO_2 ,
5 $PM_{2.5}$ and PM_{10} in space; (ii) to employ disaggregated vehicle models to detect the bottlenecks
6 of congestion in the road network, with sensitivity analyses for the effects of unknown parameters
7 (such as presences of buses and trucks, which are important sources of vehicle emissions); (iii)
8 to validate the potential sources of pollution, integrating aerial images (from providers of remote
9 sensing data such as Planet Labs) with longer and more detailed observations of pollutant sources
10 and presence of vehicles.

11 ACKNOWLEDGMENTS

12 We acknowledge Jinhua Zhao for enlightening discussions. This work would not have been pos-
13 sible without the kind support, information and data provided by Zheng Chang. This work was
14 funded in part by the MIT-Environmental Solutions Initiative, the MIT Samuel Tak Lee (STL)
15 Real Estate Entrepreneurship Lab, the New England UTC 25, and the Center for Complex Engi-
16 neering Systems (CCES) at KACST.

17 REFERENCES

- 18 [1] Dockery, D. W., C. A. Pope, X. Xu, J. D. Spengler, J. H. Ware, M. E. Fay, B. G. Ferris Jr,
19 and F. E. Speizer, An association between air pollution and mortality in six US cities. *New*
20 *England journal of medicine*, Vol. 329, No. 24, 1993, pp. 1753–1759.
- 21 [2] Hoek, G., B. Brunekreef, S. Goldbohm, P. Fischer, and P. A. van den Brandt, Association
22 between mortality and indicators of traffic-related air pollution in the Netherlands: a cohort
23 study. *The lancet*, Vol. 360, No. 9341, 2002, pp. 1203–1209.
- 24 [3] Guo, S., M. Hu, M. L. Zamora, J. Peng, D. Shang, J. Zheng, Z. Du, Z. Wu, M. Shao,
25 L. Zeng, et al., Elucidating severe urban haze formation in China. *Proceedings of the Na-*
26 *tional Academy of Sciences*, Vol. 111, No. 49, 2014, pp. 17373–17378.
- 27 [4] Kelly, F. J. and T. Zhu, Transport solutions for cleaner air. *Science*, Vol. 352, No. 6288, 2016,
28 pp. 934–936.
- 29 [5] McHugh, C., D. Carruthers, and H. Edmunds, ADMS–Urban: an air quality management
30 system for traffic, domestic and industrial pollution. *International Journal of Environment*
31 *and Pollution*, Vol. 8, No. 3-6, 1997, pp. 666–674.
- 32 [6] Vardoulakis, S., B. E. Fisher, K. Pericleous, and N. Gonzalez-Flesca, Modelling air quality
33 in street canyons: a review. *Atmospheric environment*, Vol. 37, No. 2, 2003, pp. 155–182.
- 34 [7] Baldauf, R., E. Thoma, M. Hays, R. Shores, J. Kinsey, B. Gullett, S. Kimbrough, V. Isakov,
35 T. Long, R. Snow, et al., Traffic and meteorological impacts on near-road air quality: Sum-
36 mary of methods and trends from the Raleigh near-road study. *Journal of the Air & Waste*
37 *Management Association*, Vol. 58, No. 7, 2008, pp. 865–878.

- 1 [8] Zheng, Y., F. Liu, and H.-P. Hsieh, U-Air: when urban air quality inference meets big data.
2 In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery*
3 *and data mining*, ACM, 2013, pp. 1436–1444.
- 4 [9] Zheng, Y., X. Yi, M. Li, R. Li, Z. Shan, E. Chang, and T. Li, Forecasting fine-grained air
5 quality based on big data. In *Proceedings of the 21th ACM SIGKDD International Conference*
6 *on Knowledge Discovery and Data Mining*, ACM, 2015, pp. 2267–2276.
- 7 [10] *Beijing Regional Statistic Year Book*. [http://www.bjstats.gov.cn/nj/qxnj/2014/zk/](http://www.bjstats.gov.cn/nj/qxnj/2014/zk/indexch.htm)
8 [indexch.htm](http://www.bjstats.gov.cn/nj/qxnj/2014/zk/indexch.htm), 2016, [Online; accessed 12-January-2016].
- 9 [11] Alexander, L., S. Jiang, M. Murga, and M. C. González, Origin–destination trips by purpose
10 and time of day inferred from mobile phone data. *Transportation Research Part C: Emerging*
11 *Technologies*, 2015.
- 12 [12] Çolak, S., L. P. Alexander, B. G. Alvim, S. R. Mehndiretta, and M. C. González, Analyzing
13 cell phone location data for urban travel: current methods, limitations and opportunities. In
14 *Transportation Research Board 94th Annual Meeting*, 2015, 15-5279.
- 15 [13] Jiang, S., G. A. Fiore, Y. Yang, J. Ferreira Jr, E. Frazzoli, and M. C. González, A review of
16 urban computing for mobile phone traces: current methods, challenges and opportunities. In
17 *Proceedings of the 2nd ACM SIGKDD International Workshop on Urban Computing*, ACM,
18 2013, p. 2.
- 19 [14] Zheng, Y. and X. Xie, Learning travel recommendations from user-generated GPS traces.
20 *ACM Transactions on Intelligent Systems and Technology (TIST)*, Vol. 2, No. 1, 2011, p. 2.
- 21 [15] Guttman, A., *R-trees: a dynamic index structure for spatial searching*, Vol. 14. ACM, 1984.
- 22 [16] Toole, J. L., S. Colak, B. Sturt, L. P. Alexander, A. Evsukoff, and M. C. González, The path
23 most traveled: Travel demand estimation using big data resources. *Transportation Research*
24 *Part C: Emerging Technologies*, 2015.
- 25 [17] Çolak, S., A. Lima, and M. C. González, Understanding congested travel in urban areas.
26 *Nature communications*, Vol. 7, 2016.
- 27 [18] *Traffic congestion index in Beijing*. [http://www.bjtrc.org.cn/PageLayout/](http://www.bjtrc.org.cn/PageLayout/IndexReleased/Realtime.aspx)
28 [IndexReleased/Realtime.aspx](http://www.bjtrc.org.cn/PageLayout/IndexReleased/Realtime.aspx), 2016, [Online; accessed 30-March-2016].
- 29 [19] Wardrop, J. G., ROAD PAPER. SOME THEORETICAL ASPECTS OF ROAD TRAFFIC
30 RESEARCH. *Proceedings of the institution of civil engineers*, Vol. 1, No. 3, 1952, pp. 325–
31 362.
- 32 [20] Friesz, T. L. and D. Bernstein, *Foundations of Network Optimization and Games*. Springer,
33 2016.
- 34 [21] *OpenStreetMap*. <https://www.openstreetmap.org>, 2016, [Online; accessed 18-April-
35 2016].

- 1 [22] Dial, R. B., A path-based user-equilibrium traffic assignment algorithm that obviates path
2 storage and enumeration. *Transportation Research Part B: Methodological*, Vol. 40, No. 10,
3 2006, pp. 917–936.
- 4 [23] Nie, Y. M., A class of bush-based algorithms for the traffic assignment problem. *Transporta-*
5 *tion Research Part B: Methodological*, Vol. 44, No. 1, 2010, pp. 73–89.
- 6 [24] *AMP direction API*. <http://lbs.amap.com/api/webservice/reference/direction/>,
7 2016, [Online; accessed 16-May-2016].
- 8 [25] Breiman, L., Random forests. *Machine learning*, Vol. 45, No. 1, 2001, pp. 5–32.
- 9 [26] Liaw, A. and M. Wiener, Classification and Regression by randomForest. *R News*, Vol. 2,
10 No. 3, 2002, pp. 18–22.
- 11 [27] Lindeman, R. H., P. Merenda, and R. Gold, *Introduction to bivariate and multivariate analy-*
12 *sis*, 1980.
- 13 [28] Chen, Y., A. Ebenstein, M. Greenstone, and H. Li, Evidence on the impact of sustained
14 exposure to air pollution on life expectancy from China’s Huai River policy. *Proceedings of*
15 *the National Academy of Sciences*, Vol. 110, No. 32, 2013, pp. 12936–12941.