

A topology-based algorithm using corrective distance matrix for air pollution tracing

Ziao Liu¹, Lin Gao¹, Yuyang Hong¹ and Haibo Hu^{*}

School of Software Engineering, Chongqing University, Chongqing, 410331, China.

ARTICLE INFO

Keywords:
Air pollution propagation
topology
graph visualization

ABSTRACT

The regional transmission of air pollutants is an important factor affecting climate and environment. Many existing studies are committed to mining potential propagation models. However, most models are highly dependent on data. Meanwhile, there is a lack of research on accurate pollutant tracing within and between regions. This paper presents a topology-based algorithm using correction distance matrix for air pollution tracing. The geographical distance is corrected according to meteorological factors. The distance relationship between monitoring stations is expressed as a directed graph. The distance threshold is used to filter the directed graph to obtain relevant subgraphs of the target station. By constructing pollutant concentration transfer equations, the reachability matrix of pollution propagation is calculated to get the pollution contribution of each node. The single point tracing result of the target station is obtained from the maximum contribution. Meanwhile, we design the air pollution tracing visualization system. The multi-dimensional views in the system are convenient for users to perceive propagation laws and explore results. Based on above work, we conduct an intraregional case study of Beijing in 2013 and an interregional case study of central and eastern China in 2014. The results of case studies confirm the effectiveness of the algorithm and visualization system in the field of air pollution propagation. Moreover, we evaluate the performance of the algorithm and get good feedback according to user's experience with MeteoInfo software and our system. Therefore, our work can accurately trace the source of air pollutants and provide technical support for meteorological experts.

1. Introduction

As one of the environmental problems, air pollution has attracted continuous attention all over the world. Relative research found that people's burden of diseases such as respiratory infection, heart disease, lung cancer may be led by air contamination(Jin et al., 2021). The Health Effects Institute report that the death rate of air pollution in 2019 has exceeded that of traffic accidents. Moreover, Parajuli et al. (2021) found that a single pollutant has more frequent and significant effects on human body. Vo et al. (2022) found that 80 % of PM2.5 will be deposited in the head airway, while 75 % of PM0.1 will be deposited in the alveolar area. In addition, environment pollution has a negative impact on the economic situation, resulting in the total loss rate of India reaching 1.36% of GDP in 2019 (Pandey et al., 2021). There has always been an interdependent relationship between urban sustainable development and air pollution(Ahmad et al., 2021). According to demographic monitoring, excessive SO₂ in the air leads to the reduction of the overall growth and development of plants(Cotrozzi, 2020). In order to effectively reduce the harm of air pollution to these aspects,

accurately determining the source of pollutants has become an essential part.

There has been a lot of research on the propagation of air pollution. The main research direction is to build an ideal propagation model and mine the symbiotic relationship under association rules. HYSPLIT4 (Stein et al., 2015) is a system that can not only calculate simple single trajectory problems but also solve complex distribution and precipitation problems through the smoke and particle method. Although HYSPLIT4 is most widely used in air pollution source identification, it ignores the problem of small-scale air transmission. Considering the influence of wind speed and direction on pollution propagation, some articles use Lagrange dispersion model (LDM) (Carvalho and De Vilhena, 2005) to avoid the numerical pseudo diffusion. On the other hand, some researches focus on mining hidden symbiotic relationships through association rules to infer the relationship between pollution transmission. For example, the visual exploration framework Co-Occurrence Pattern Exploration (COPE) (Li et al., 2018) looks for co-occurrence patterns based on time-space series to evaluate the intensity of occurrence patterns and make predictions. Fan et al. (2020) made a comprehensive analysis of the spatio-temporal variation of urban air pollution in China during 2014–2018. Guo et al. (2019) proposed Time Correlation Partition (TCP) tree, which provides situational awareness analysis between correlation data and variables. Huang et al. (2018) adopted CAMx adjoint propagation model to realize dynamic inversion and grid quantitative propagation of emission sources. In this way, the spatio-temporal characteristics of pollution sources and their contribution rates can

*This work was supported in part by the National Natural Science Foundation of China under Grant U1836114.

^{*}Corresponding author. Haibo Hu PhD, Associate Professor, School of Software Engineering, Chongqing University, Room 503, Software Building, No.55 South Road of University Town, Shapingba District, Chongqing, 410331 CN.

 haibo.hu@cqu.edu.cn (H. Hu)

ORCID(s): 0000-0003-3237-0690 (Z. Liu); 0000-0001-9115-730X (L. Gao); 0000-0003-3901-7769 (Y. Hong); 0000-0001-8442-5222 (H. Hu)

¹Authors contributed equally.

be evaluated quantitatively. The point-to-point propagation of pollutants can be mapped to the directed graph structure. The knowledge of graph theory and topological geometry can be feasible to apply to the pollution propagation model. Reachability matrix is often used to describe the degree that can be reached between nodes of a directed graph after a certain length of path. Topology-based reachability analysis plays an important role in exploring the relationship between two points of a directed graph. Xu and Hong (2012) proposed a finite solution method based on reachability matrix. The sufficient and necessary conditions for the reachability of deterministic and nondeterministic automata are obtained by using the semi-tensor product of matrices. Scarpato et al. (2019) proposed the Reachability Matrix Ontology (RMO) to describe the network structure and network security. RMO proposes an innovative method to investigate whether one node can reach another node through the ISO/OSI layer protocol. Kim et al. (2021) proposed a new incremental algorithm framework AORM, which is a reachability matrix framework of any order. AORM can handle both directed and disconnected networks, such as reference networks as well as various types of real data sets.

According to the feedback of meteorological experts, we find that there is a lack of effective visualization system to display data and analyze tracking results. Visual presentation relies on a lot of data analysis technology. Regression analysis (Demuzere et al., 2008), statistical analysis (Sharovsky et al., 2004) and correlation analysis (Qu et al., 2007) are often used as models for analyzing air quality data. The establishment of visualization system is helpful to discover the potential information and explore the value behind the data. There are a lot of researches focused on geographic information due to the close relationship between air pollution data and geographic information. By introducing the concept of air wrapping to draw the pollutant propagation map, Deng et al. (2019) developed the AirVis system, which is based on Word2Vec model to effectively extract transmission models. Ren et al. (2020) abstracted the complex inter regional propagation relationship into a dynamic network. In their research, visual analysis is introduced to explore the spatio-temporal multiple propagation model. At the same time, a model based on particle tracking is proposed to construct the pollution propagation network under multi-source factors. Deng et al. (2021) proposed VisCas system, which combines the extraction of pollution patterns with interactive visualization technology. Based on the triple optimization strategy, a novel influence view is designed to infer and explain the potential cascading patterns in the spatio-temporal context. Tian et al. (2021) developed an inline visualization system from the perspective of fine-grained air quality. Liu et al. (2021) proposed an integrated AQEyes system. By integrating multiple pipelines such as data preprocessor and characterizer, the unsupervised anomaly detection method is adopted to effectively check the anomalies. Qu et al. (2020) constructed the AirExplorer system. Combined with hierarchical piecewise linear representation and dynamic time warping, the time series patterns of interest are accurately

queried. Zhou et al. (2017) used Multidimensional Scaling (MDS) to convert data into 2D maps. Combining hierarchical clustering, Voronoi diagram and interactive technology help experts explore and extract the propagation modes.

Inspired by above work, we propose a topology-based algorithm using correction distance for air pollution tracing. In order to maximize the data value, we preprocess the high-dimensional time series air pollution data. The reachability matrix between pollution monitoring stations is calculated by using classical topological geometry. According to the set thresholds, the correlation between monitoring stations is analyzed. Based on tracing results, we build the air pollution tracing visualization system. Our focus is on air quality monitoring stations. The goal of the system is to facilitate experts to explore the propagation law clearly and put forward urban air quality management plans.

In the remainder of this article, we describe the preprocessing work of the data set in Section 2. We present the details of the tracing algorithm in Section 3. In addition, we describe the visualization work in Section 4. We conduct two case studies with and between regions in Section 5. We evaluate the performance of algorithm and conduct user study in Section 6. We finally present the summary and outlook of this work in Section 7.

2. Data preprocessing

2.1. Data set

The data set used in this paper is high-resolution air pollution reanalysis data set of China released by the Institute of Atmospheric Physics, Chinese Academy of Sciences and other units. The data set includes the grid data of six conventional pollutants in China's ambient air quality standards. It is also the reanalysis data of the national ambient air quality monitoring network and nested grid air quality prediction model (naqpms) of China National Environmental Monitoring Station (CNEMC)(Kong et al., 2021). <https://doi.org/10.11922/scienceb.00053>

The data set involves 2751 monitoring stations, covering 375 cities in China, with a time range from 2013 to 2018. Stations have collected daily mean data and hourly data in January from 2013 to 2018. Each data information includes the concentration of six air pollutants: PM_{2.5}, PM₁₀, NO₂, SO₂, O₃ and CO₂. In addition, each station also recorded the hourly or daily wind direction, wind speed, temperature, humidity and air pressure information, as well as the longitude and latitude data of the station. The basic information is in Table 1:

2.2. Inverse geocoding

The position variables in the data set are expressed in longitude and latitude. However, We hope to obtain provinces, cities, districts and counties. This needs to obtain the geographic information of each city and carry out inverse geographic coding according to the longitude and latitude. In this paper, our local geographic information database is used to realize the inverse geocoding. We code the longitude and

Table 1
Data set variables

Data type	Name	Unit	Explain
Pollutant	PM _{2.5}	μg/m ³	Diameter ≤ 2.5 μm
	PM ₁₀	μg/m ³	Diameter ≤ 10 μm
	NO ₂	μg/m ³	Nitrogen dioxide
	SO ₂	μg/m ³	Sulphur dioxide
	O ₃	μg/m ³	Ozone
	CO	μg/m ³	Carbon monoxide
	U	m/s	Latitudinal wind speed
	V	m/s	Radial wind speed
	RH		Relative humidity
Meteorology	TEMP	K	
	PSFC	Pa	Ground pressure

Algorithm 1 The topology-based algorithm for air pollution tracing

```

Input: data: air pollution data set; Dth: distance threshold;
       Pth: pollution threshold; i: target station number; t:
       initial date; r: sliding time window unit; times: sliding
       times
Output: k: source station number
1: for station i do:
2:   move the sliding window unit r;
3:   while times > 0 do:
4:     load data;
5:     compute humidity diffusion coefficient H;
6:     compute wind direction and wind speed diffu-
       sion coefficient W;
7:     compute real spherical distance distance;
8:     use comprehensive evaluation mode to get the
       modified distance matrix D;
9:     if Dij > Dth then:
10:      compute reachability matrix A
11:      use logarithm and diagonal differentiation to
       get A''';
12:      compute the maximum value A''ik;
13:      if A''ik > Pth and i ≠ k then:
14:        station k is the source of pollution at time
       t;
15:      end if
16:    end if
17:  end while
18: end for

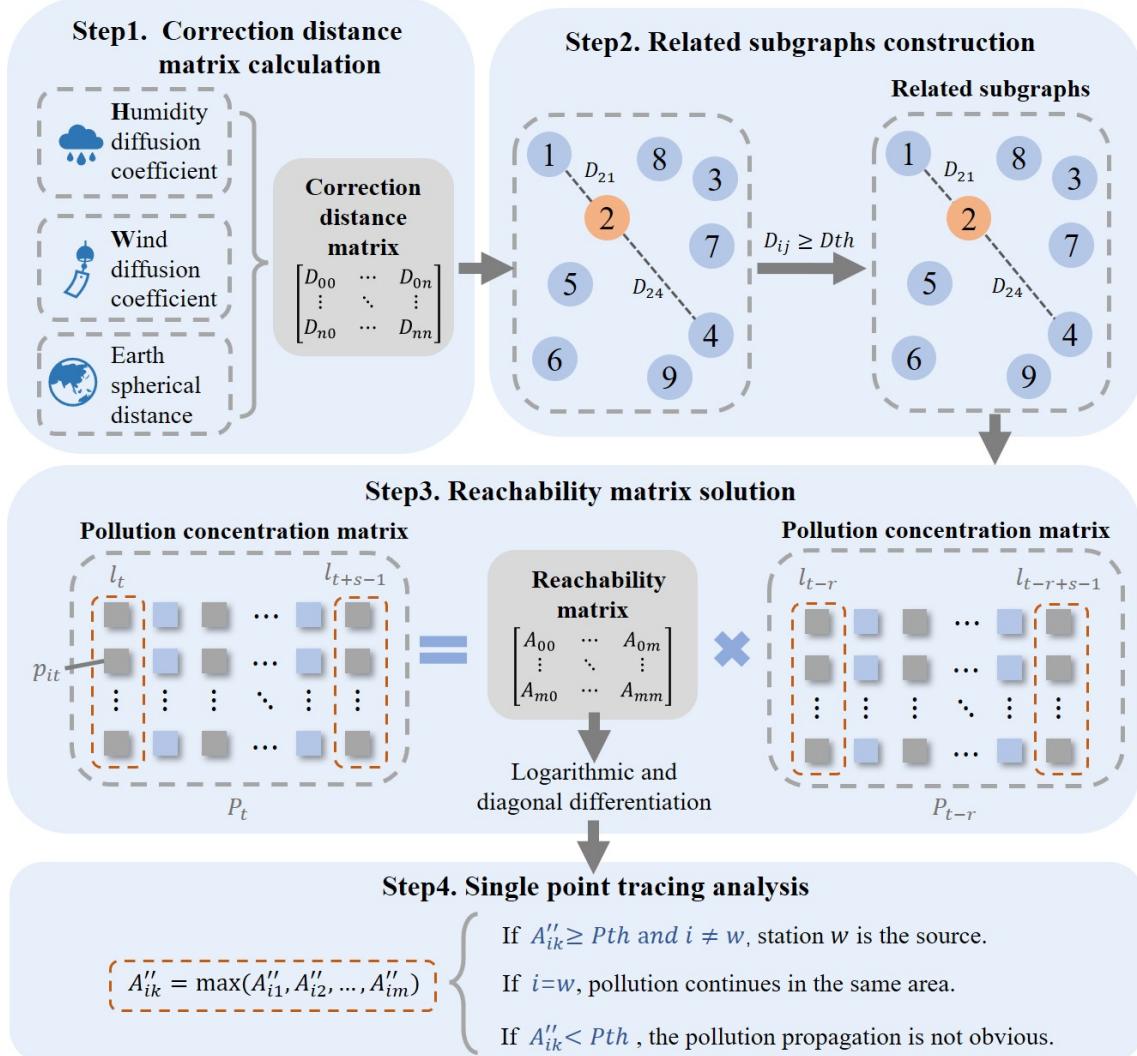
```

According to the research of Li et al. (2017) and Li et al. (2018), wind direction and wind speed diffusion coefficient W is calculated in Eq. 3.1. Where d_A and d_B represent the azimuth represented by the wind direction at grid cells A and B respectively. d_M represents the azimuth of AB direction. F is the influence of wind direction on diffusion. The diffusion distance of downwind direction is small and that of upwind direction is large. The weights w_A and w_B are the influence of wind speed on diffusion. For the downwind direction, the wind speed and diffusion distance are negatively correlated. For the upwind direction, the wind speed and diffusion distance are positively correlated. L_{AB} represents the type of adjacency relationship between A and B . If it is edge adjacency, $L_{AB} = 1$. If it is point adjacency, $L_{AB} = 2$. v_A and v_B represent the wind speed. $sgn(x)$ represents that the higher the wind speed is, the greater the influence on the diffusion distance becomes, and vice versa.

(2) \rightarrow

The real spherical distance $dist$ between monitoring stations can be calculated by Gaussian coordinate conversion of longitude and latitude in the data set. The humidity diffusion coefficient H is multiplied by the real spherical distance to obtain the corrected humidity distance. Multiply the coefficient W by the real spherical distance to obtain

$$H = \frac{0.0101T^{1.75} \sqrt{\frac{1}{M_A} + \frac{1}{M_B}}}{P[(\sum vol_A)^{1/3} + (\sum vol_B)^{1/3}]} \quad (1)$$


Fig. 1: Flow chart of algorithm steps

the corrected wind distance. The comprehensive evaluation model is used to give weight to the above three distances and the modified distance matrix D is obtained in Eq. 3.

$$D = w_1 \times dist + w_2 \times W \times dist + w_3 \times H \times dist \quad (3)$$

3.2. Related subgraphs construction

The distance matrix D between monitoring stations is regarded as a directed graph, where D_{ij} represents the influence distance of station j on station i . We suppose the ranges of i and j are positive integers from 0 to n . Determine the target station and filter the D_{ij} using the distance threshold D_{th} . Finally, the relevant subgraph of the target station are obtained. We assume that the relevant subgraph contains m stations at all.

3.3. Reachability matrix solution

The reachability matrix describes the degree that can be achieved after a certain length of path between the nodes of the directed connection graph. In this model, it is used to represent the pollution contribution between monitoring

stations in the relevant subgraph. A_{ij} represents the pollution contribution of station j to station i in the relevant subgraph.

$$A = \begin{bmatrix} A_{00} & A_{01} & \cdots & A_{0m} \\ A_{10} & A_{11} & \cdots & A_{1m} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m0} & A_{m1} & \cdots & A_{mm} \end{bmatrix} \quad (4)$$

The solution of the reachability matrix depends on the pollutant concentration transfer equation and the solution process is as follows. In Eq. 5, P_t is the pollution concentration matrix of the monitoring station in the relevant subgraph. r is the sliding time window unit.

$$P_t = A \times P_{t-r} \quad (5)$$

According to the relevant subgraph, the pollutant concentration matrix of the monitoring station in the time interval is constructed in Eq. 6. s represents the time interval,

206 l_t is the concentration column vector of each station in the
 207 correlation subgraph at time t . p_{it} represents the pollution
 208 concentration of station i at time t .

$$\begin{aligned} P_t &= [l_t, l_{t+1}, \dots, l_{t+s-1}] \\ P_{t-r} &= [l_{t-r}, l_{t-r+1}, \dots, l_{t-r+s-1}] \\ l_t &= [p_{1t}, p_{2t}, \dots, p_{nt}]^T \end{aligned} \quad (6)$$

209 We use logarithm and diagonal differentiation to prepro-
 210 cess the reachability matrix for subsequent tracing analysis:
 211 1) Log the reachability matrix A to the matrix A' .

$$A'_{ij} = \begin{cases} \ln(A_{ij} + 1) & A_{ij} \geq 0 \\ -\ln(|A_{ij}| + 1) & A_{ij} < 0 \end{cases} \quad (7)$$

212 2) diagonally differentiate matrix A' to the matrix A'' .

$$A''_{ij} = \begin{cases} A'_{ij} & i = j \\ A'_{ij} - A'_{ji} & i \neq j \end{cases} \quad (8)$$

213 3.4. Single point tracing analysis

Eq. 9 shows the process of extracting the station with the highest absolute value of contribution to the target station as the tracing result to achieve single point tracing. Take the corresponding row or column of the target station i in the reachability matrix as a set and calculate the maximum value A''_{ik} .

$$A''_{ik} = \max(A''_{i1}, A''_{i2}, \dots, A''_{im}) \quad (9)$$

214 Set the pollution threshold Pth . If $A''_{ik} \geq Pth$ and $i \neq k$,
 215 then k is the tracing result of station i and A''_{ik} is the impact
 216 intensity of pollution source. If $A''_{ik} < Pth$, it is considered
 217 that the pollution transmission is not obvious in this time
 218 period. That is, the tracing result of station i is not obtained.
 219 If $i = k$, the pollution in the area where station i is located
 220 continues.

221 Use time units r as a sliding window. Move the sliding
 222 window and repeat the above steps to draw the pollutant
 223 propagation path. While the tracing results are obtained, the
 224 influence degree of each transmission process can also be
 225 obtained.

226 4. Visualization of our work

227 Although the algorithm has been able to give the tracing
 228 results in the third part, it can't intuitively tell users more
 229 detailed information about the source place and the target
 230 place. Therefore, we design a visual analysis system, so that
 231 users can more intuitively see the tracing results and detailed
 232 information of tracing locations. The exploration pipeline of
 233 air pollution tracing visualization system is shown in Fig. 2.

235 4.1. Option window

236 The leftmost part of the system is the option window,
 237 which is composed of views A, B and C. The user selects
 238 the date through view A (the time range is from January
 239 1, 2013 to December 31, 2018). View B uses the method
 240 of transforming word cloud to enable users to select the
 241 research area. In this view, the user can determine whether
 242 the area is selected according to the bold of the font. At the
 243 same time, users can select multiple areas. View C inputs
 244 control variables for tracing results. Users can generate trac-
 245 ing results by setting distance threshold, pollution threshold,
 246 sliding window size, sliding step and sliding times.

247 4.2. Basic information window

248 The basic information window are determined by views
 249 A and B. View D shows the AQI thermodynamic map of
 250 the selected area over a specific time range. The six colors in
 251 View D (varying from pale white to dark blue) correspond to
 252 each of the six levels of pollution. View E shows the change
 253 of climatic conditions during this period, specifically wind
 254 direction and wind speed. The color from white to converge
 255 with blue shows wind speed decrease. The color from white
 256 to converge with red represents wind speed increase.

257 4.3. Tracing result window

258 View F is the pollution propagation path map. After the
 259 user determines the parameters such as time, tracing range
 260 and pollution threshold, the tracing results of the algorithm
 261 can be visually displayed by the figure. The line between the
 262 two places indicates the transmission route of pollutants and
 263 the direction of the arrow indicates the transmission direc-
 264 tion of pollutants. We also draw pie charts at the locations of
 265 each monitoring station. And the three components of each
 266 pie chart are W_1 , W_2 , and W_3 respectively, which form a one-
 267 to-one correspondence with the distance weight distribution
 268 chart (view I).

269 4.4. Multi-view analysis window

270 4.4.1. Trend comparison chart

271 View G is called the trend comparison chart. This view
 272 is designed to show the AQI trend comparisons between the
 273 source city and the target city. Among them, the red column
 274 is the AQI value of the source city and the blue column is the
 275 AQI value of the target city. Through the AQI value display
 276 for four consecutive days, we can see the trend change of air
 277 pollutants in the source and target points. On the right axis
 278 is the range of the AQI value. Through a lot of tests, we find
 279 that setting the threshold to 0-1000 can cover all situations.
 280 The lowest coordinate axis is the number of days displayed.

281 4.4.2. Influence degree sankey chart

282 View H is called the influence degree sankey chart. The
 283 tracing results of polluted areas reflects the relationship be-
 284 between data flow. This view is convenient for meteorological
 285 experts to observe the relationship between areas, to clarify
 286 the trend of pollutants between areas. The left vertexes are
 287 the source areas. The right vertexes are the target areas. The
 288 edges of the view constitute the pollution relationship. The

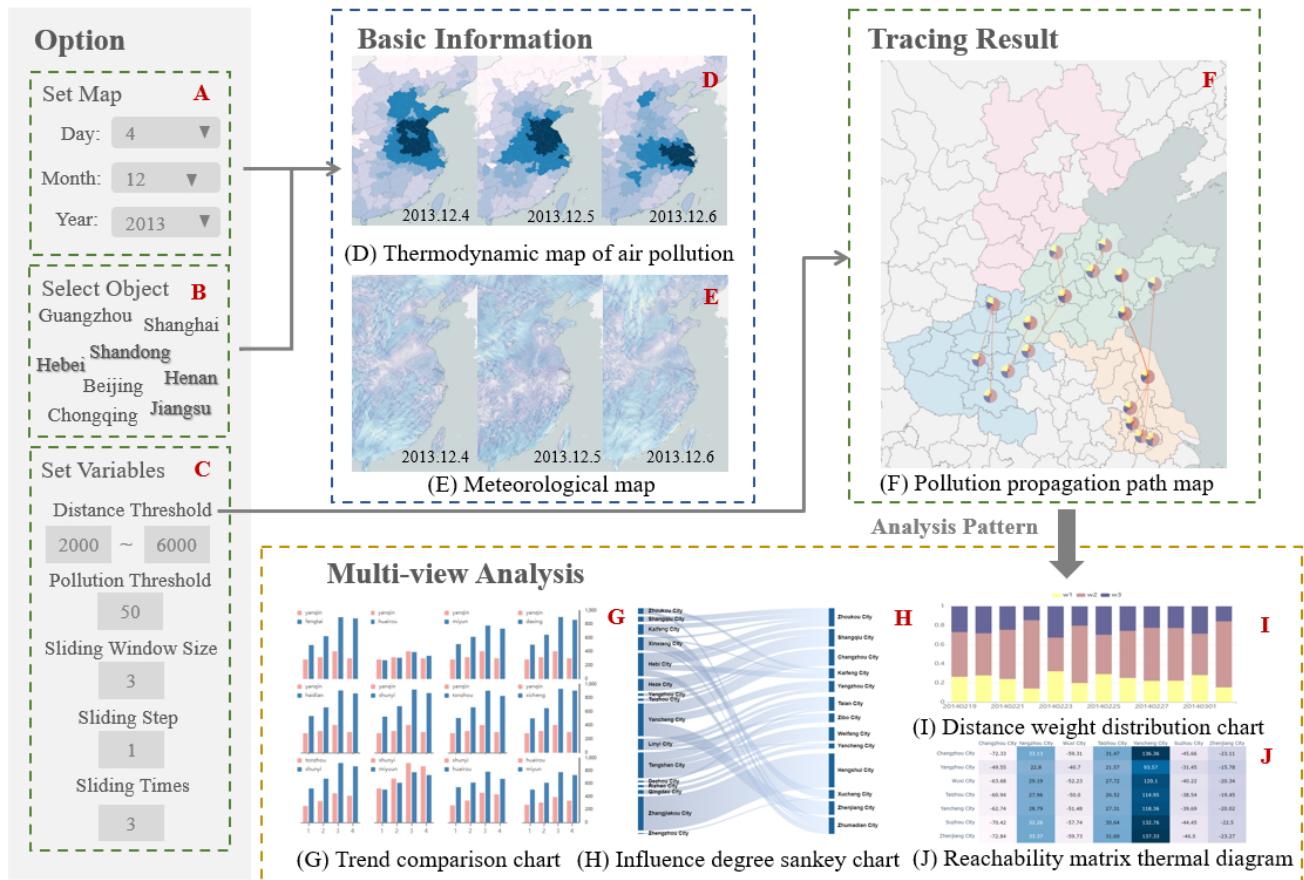


Fig. 2: The exploration pipeline guides users to employ the air pollution tracing visualization system. After setting variables in option window, users can explore the tracing results from algorithm with the help of basic information window. The multi-view analysis window helps users make some analysis work better.

flow is the degree of influence. The greater the degree of influence, the greater the flow. We can usually find that multiple areas will cause varying degrees of pollution to a area. At the same time, one area will also cause varying degrees of pollution to different areas. Therefore, the area with the largest flow value will become the focus of researchers.

4.4.3. Distance weight distribution chart

View I is called distance weight distribution chart. From the introduction of the algorithm, we can know that in order to consider meteorological factors, we modify the original distance. A comprehensive evaluation model is adopted to assign weights to each correction distance. We all visualize these three weights. w_1 represents the weight of the true distance of the sphere. w_2 is the weight of humidity correction distance. w_3 is the weight of wind direction and wind speed correction distance. Given the fact that the sum of the three weights always adds up to 1, we use stacked histograms. The abscissa is the time point and the ordinate is the weight value. By comparing the distribution of weights at different time points, it is found that the degree of meteorological impact changes. At the same time, the maximum weight represents the dominant factor.

4.4.4. Reachability matrix thermal diagram

View J is called the reachability matrix thermal diagram. From the introduction of the algorithm, we can know that the reachability matrix represents the pollution contribution between monitoring stations in the relevant subgraph. We use visual methods to facilitate meteorologists to find more relevant information about urban pollution. The matrix items are represented by squares. The color of the box indicates the size of the contribution. The greater the contribution is, the darker the color is. The contribution is positive or negative. Users can query the specific value of the contribution according to the tag. According to the algorithm, the station with the greatest contribution is the tracing result.

5. Case studies

The effectiveness of the algorithm is verified in this part. In order to explore the path of pollution propagation in the region, we take Beijing, China as the research object. The algorithm is also suitable for the research between regions. Therefore, we conducted pollution propagation restoration for the severe haze event in central and eastern China.

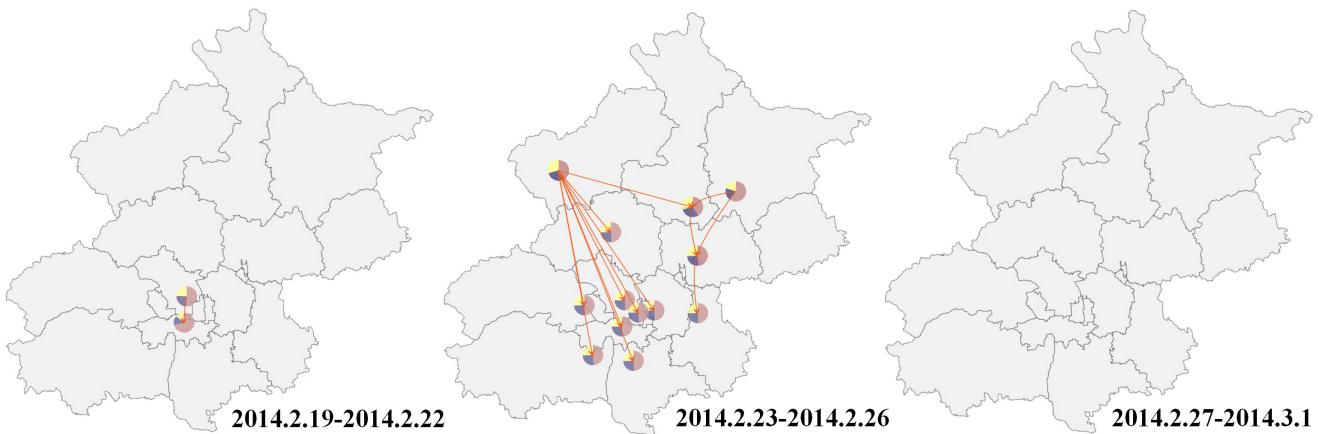


Fig. 3: Pollution propagation path map of Beijing

331 Combined with the follow case studies, the algorithm and
332 system are further improved.

5.1. Beijing pollution event

333 As the capital of China, Beijing's industrial economy is
334 developing rapidly. However, the excessively rapid industrial
335 development has led to a sharp deterioration of air conditions
336 in Beijing. Many smog incidents have seriously endangered
337 the lives and health of local people. From 2013 to 2017, in
338 order to prevent the spread of air pollutants in the region or
339 even outside the region, Beijing stepped up the rectification
340 of Beijing's industry and formulated a number of measures
341 against air pollution. As a typical city of air pollution trans-
342 mission, we choose Beijing as the first case to analyze. Based
343 on the data set and related work, we find that Beijing had
344 a very obvious phenomenon of pollution transmission from
345 February 23, 2014 to February 26, 2014. In order to study
346 the process of air pollution event in Beijing, we set the length
347 of the algorithm's sliding window as 3 days/cycle and the
348 sliding step as 1 day/time. Meanwhile, in order to achieve
349 the optimal effect, we adjust the distance threshold to 300
350 and the pollution threshold to 200. The result of tracing path
351 is shown in Fig. 3. It shows the source of pollutants and
352 propagation paths in districts of Beijing.
353

5.1.1. Early stage of pollution event (2014.2.21-2014.2.23)

354 The data set shows that the air pollution index was good
355 before February 23, with the AQI remaining between 51 and
356 100. The tracing results only show that the pollution spread
357 from Fengtai District of Beijing to Haidian District, with a
358 low degree of influence.

5.1.2. Middle stage of pollution event (2014.2.24-2014.2.26)

359 At this stage, the overall AQI of Beijing increased to 300-
360 350, reaching the level of severe pollution. Compared with
361 the early stage, the tracing results become more complex
362 and diverse. We can know the pollution sources of each area

367 through Fig. 3. Meanwhile, The degrees of pollution in the
368 polluted area are shown in Fig. 4. We find that Tongzhou
369 District was seriously affected since the pillar pointing to it
370 is the thickest.

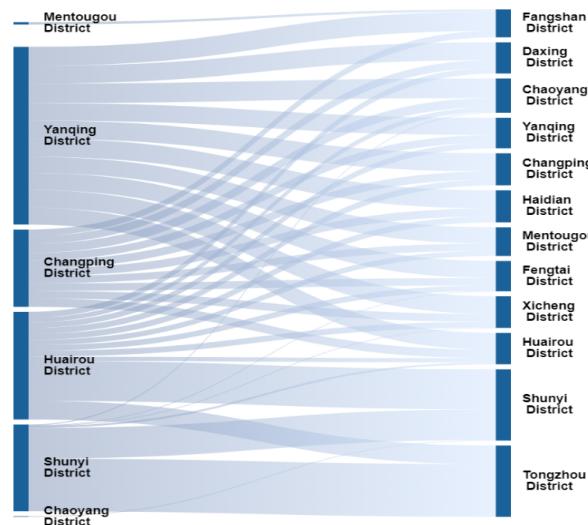


Fig. 4: Influence degree sankey chart of Beijing

371 To analyze the pollution factors in Tongzhou District, we
372 draw Fig.5. Fig. 5 shows that the main factor of transmis-
373 sion is wind. The northwest wind gradually strengthened in
374 Tongzhou area from 23rd to 26th, which makes Shunyi in
375 the northwest spread the pollutants to Tongzhou. The con-
376 centration of pollutants in Shunyi is higher, leading to more
377 pollutants spreading to Tongzhou. Combined with the Fig.
378 6, we can see that the pollutants in Tongzhou rose rapidly in
379 these four days, reaching a maximum of 460.38. All these
380 information above shows that Tongzhou was infected by
381 Shunyi and Huairou because of the wind.

382 Simultaneously, Fig. 7 shows that Yanqing District had
383 a wide spread of influence and had a certain degree of pollu-
384 tion impact on 9 surrounding areas, including Changping,

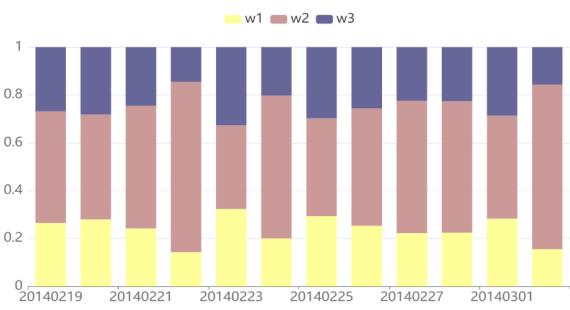


Fig. 5: Distance weight distribution chart of Beijing

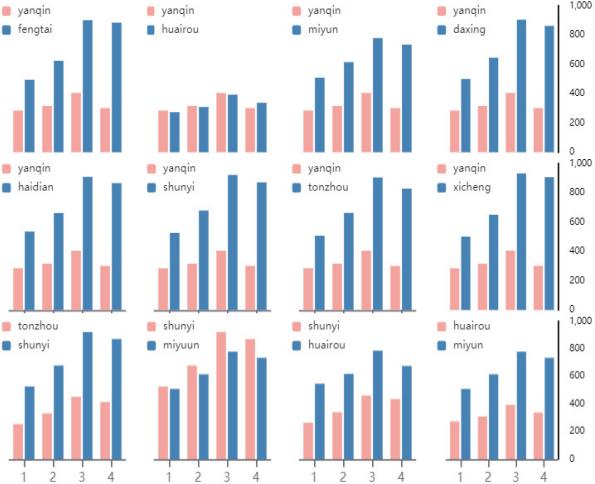


Fig. 6: Trend comparison chart of Beijing

Fengtai District	Daxing District	Yanqing District	Huairou District	Changping District	Chaoyang District	Haidian District	Xicheng District	Tongzhou District	Mentougou District	Shunyi District
-108.66	-171.07	304.3	75.04	152.53	-41.68	-10.43	-64.45	-109.61	-59.34	34.7
District	-94.65	-149.01	265.07	65.36	132.87	-36.29	-9.08	-56.14	-95.47	-51.69
District	-106.05	-166.95	296.98	73.23	148.66	-40.67	-10.18	-62.79	-106.97	-57.92
District	-103.9	-163.58	290.98	71.75	145.86	-39.85	-9.97	-61.63	-104.8	-56.75
District	-105.83	-166.61	296.37	73.08	148.56	-40.59	-10.16	-62.77	-106.75	-57.79
District	-108.73	-171.18	304.5	75.09	152.63	-41.7	-10.44	-64.49	-109.68	-59.38
District	-105.88	-166.7	296.53	73.13	148.63	-40.62	-10.17	-62.8	-106.81	-57.81
District	-98.53	-155.13	275.95	68.05	136.32	-37.79	-9.46	-58.44	-99.39	-53.81
District	-100.92	-158.9	282.67	69.72	141.68	-38.72	-9.69	-59.87	-101.81	-55.1
District	-103.1	-162.31	288.73	71.19	144.73	-39.53	-9.89	-61.15	-103.99	-56.32
District	-99.59	-156.78	278.99	68.77	139.8	-38.19	-9.56	-59.07	-104.45	-54.39

Fig. 7: Reachability matrix thermal diagram of Beijing

385 Huairou and Haidian. Although Yanqing District spread
386 pollution to surrounding areas, the overall trend of trans-
387 mission was from the northwest to the south of Beijing.
388 Combined with the dominant factor shown in the trend
389 comparison chart, we infer that the wind in Beijing was
390 mainly from the west, and the wind tended to be excessive
391 to the northwest. To verify the visual results, we investigate
392 the meteorological conditions in Beijing at this stage and get
393 the same results. This verification confirm the accuracy of
394 the model inference and confirm the validity of the model.

5.1.3. Late stage of pollution event (2014.2.27-2014.3.1)

The air pollution level in Beijing improved significantly on 27th. The AQI index was below 50, reaching an excellent level. The tracing results only show the pollution influence of Daxing District on Xicheng District, and the pollution influence degree was extremely low which is almost negligible. The humidity changed greatly from 68.096 on the 26th to 25.298 in this stage. The temperature dropped more and the possibility of rainfall was great. Weather conditions greatly promoted the deposition of pollutants and severely reduced the serious pollution from 23rd to 26th. The actual situation shows that under the action of strong cold air and north wind, pollutants were effectively removed. The monitoring data shows that the concentration of fine particles was between $3\text{--}21 \mu\text{g}/\text{m}^3$, which means the air quality reached the first-class excellent level.

The study of Beijing pollution event is a tracing study in a polluted region. According to above analysis, the dominant factor of the event is wind direction and wind speed. Our algorithm and visualization system make an accurate and comprehensive analysis of the event process. The transmission path of pollutants is shown for relevant experts to make follow-up treatment research. Relevant experts show that the tracing results are effective and beneficial to pollution control and policy-making.

5.2. Severe haze event in central and eastern China

The central and eastern regions of China have less wind in winter, and the meteorological conditions are not conducive to the diffusion of air pollutants from central and southern North China to the Yangtze River Delta region. Internal settlement and diffusion of air pollutants in the central and eastern regions often occur.

In December 2013, the most serious long-term and large-scale spread process of high-concentration particulate matter occurred in the middle and eastern regions of China. This event is called the severe haze event in central and eastern China in December 2013(Li et al., 2015). The pollution spread of this incident involved Tianjin, Hebei, Shandong, Jiangsu, Shanghai and other provinces. The daily average concentration of $\text{PM}_{2.5}$ exceeded $150 \mu\text{g}/\text{m}^3$, and even reached $300\text{--}600 \mu\text{g}/\text{m}^3$ in some areas, with the air quality index reaching level 6 (Fig. 8).

Therefore, we carry out a series of studies on the tracing process of air pollution in some central and eastern provinces (Hebei, Shandong, Jiangsu and Henan). We hope that the algorithm could identify effective pollution sources and provide suggestions for further treatment and policy intervention. This paper studies the track of PM_{2.5} in four central and eastern provinces from December 1st, 2013 to December 10th, 2013. To verify whether the tracing results are in line with the actual situation, this paper uniformly sets the sliding window value to 3 days/cycle, the sliding step size to 1 day/time, the distance threshold to 2000 km, and the pollution threshold to 50. The second distance threshold is set as 6000 km to effectively study the cases of medium

and long-distance tracing process. The final tracing results are shown in Fig. 9.

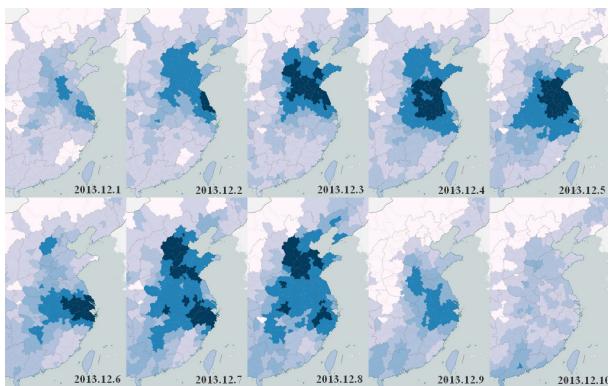


Fig. 8: Thermodynamic map of haze event

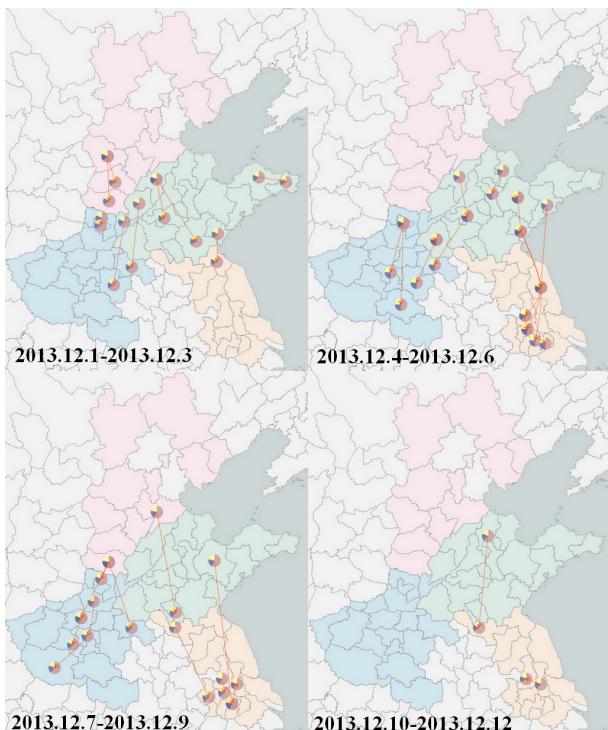


Fig. 9: Pollution propagation path map of haze event

5.2.1. Early stage of haze event (2013.12.1-2013.12.3)

Since December 1, PM_{2.5} had appeared in Shandong and Henan. The tracing results show that the phenomenon of fine particulate pollutants transmitted to the south in Hebei, Liaocheng, Shandong and Puyang, Henan also gradually transmitted fine particulate pollutants to the southwest, resulting in the increase of PM_{2.5} in Shangqiu and Zhoukou. The internal pollutants in Shandong not only affected the southwest cities but also transported fine particles of pollutants to the southeast of Shandong and the northeast of Jiangsu. For example, Rizhao City affected the aggravation of pollution in Lianyungang City. According to the diffusion

and aggravation of pollution, the scope of pollution was expanded to most areas of north and east China on December 2. On December 3, the pollution situation continued in most parts of the East. AQI indexes of several cities were higher than 200 and up to more than 400.

5.2.2. Middle stage of haze event

(2013.12.4-2013.12.6, 2013.12.7-2013.12.9)

The middle stage of the pollution transmission event is the period when the pollution was the most serious and the pollution transmission was more obvious and complicated. Therefore, we focus on using views to analyze the middle stage of pollution transmission event. Fig. 8 shows that Jiangsu was seriously polluted and Fig. 9 shows that it had a wide range of pollution transmission. Firstly, we analyze that the wind is the dominant factor in the stage. According to the meteorological map (Fig. 10), we find that the wind direction gradually changed to the northeast in December, which caused pollutants spreading from Zhangjiakou City to Hengshui City. Meanwhile, pollution spreaded more seriously in Jiangsu, starting from Yancheng to Changzhou City, Yangzhou City and other cities in the southeast of Jiangsu.

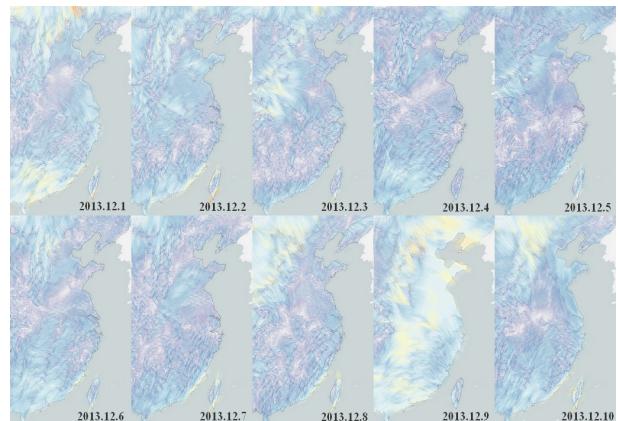


Fig. 10: Meteorological map of haze event

Although the atmospheric diffusion conditions improved during December 5 and 6, the pollution influence in Jiangsu was still increasing. In this period, the AQI in southeast Jiangsu increased sharply, exceeding 600 and reaching the serious pollution level.

From December 7 to December 9, affected by the high-altitude cyclone, the relative humidity in Shandong increased and the wind speed decreased. The air condition basically stopped the spread of internal pollutants. However, there was a wide spread in Henan. From Fig. 9, since December 7, the northeast of Henan (such as Hebi City and Xinxiang City) had a pollution impact on the cities in central Henan. From December 8 to 9, the pollutants moved further to the southwest, affecting as far as Nanyang City. According to Fig. 11, the AQI values of all cities in Henan showed an upward trend during this stage. The AQI value of Xuchang City reached 600 and the main influent was PM_{2.5}. It reached the level of moderate pollution.

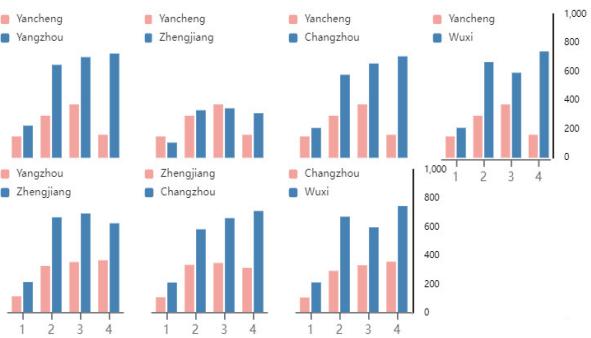


Fig. 11: Trend comparison chart of haze event

According to Fig.10, we can see that see an increase in wind speed on December 9, resulting in a decrease of pollutants. The same information can also be reflected in Fig 9. We can see that due to the increase of the wind speed on the 9th, the number of pollution tracks from December 10th to December 12th also decreased sharply.

5.2.3. Late stage of haze event (2013.12.10-2013.12.12)

During this period, the pollution transmission in various provinces and cities was greatly reduced. The pollution transmission events basically stop. According to the query of relevant meteorological factors, there was a large area of rainfall in the later stage of transmission and the humidity increased in many cities in the southeast. This leaded to a large settlement of pollutants. The inference obtained in this paper is basically consistent with the objective conditions.

5.3. Effectiveness analysis of algorithm

The severe haze event in central and eastern China is a case of pollution transmission involving many provinces. According to the above research, our algorithm can effectively restore the pollution propagation paths. Assisted by the visualization system, relevant experts can better explore the law of pollution propagation. Experts believe that inter regional collaborative governance is essential in air pollution control. The division of multi regions should consider the links between regions in economy, transportation, policy and so on. Therefore, the case further verifies the effectiveness of our work.

6. Discussion

6.1. Universality analysis of algorithm

This part mainly discusses the universality test of algorithm input. The algorithm is less dependent on the data set than the traditional forward propagation and backward trajectory algorithms. The algorithm can also guarantee the accuracy of backward trace tracing results when the requirement of data set is not high.

The common PSCF factor method (Begum et al., 2005) for identifying sources based on air flow track analysis requires obtaining the residence time of air mass in each

region. However, it is difficult to collect this data set. HYS-PLIT4 and LDM models are widely used to trace the particle sedimentation at different times, so they require higher height data. However, in the process of studying backward trajectory propagation, it is difficult to accurately obtain the data of every influence factor at every moment. Inaccurate impact factor data have a significant impact on both process variables and final results. Therefore, the above algorithms or models are highly dependent on the data set. The deviation of data sets often leads to error in tracing results.

Our algorithm requires less precise data input. In addition, user can effectively adjust some thresholds and parameters based on the system. It is of great help to remedy the defect of data accuracy and obtain satisfactory tracing results. Therefore, the algorithm is generally applicable to most data sets. It surpasses most existing algorithms and models in universality.

6.2. Sensitivity analysis of algorithm

This part is the sensitivity analysis of the algorithm. By analyzing the parameters of the algorithm, we find that the distance threshold and the pollution threshold are parameters with great uncertainty. So we focus on the sensitivity analysis of distance threshold and pollution threshold. We take Beijing from February 23, 2014 to February 26, 2014 as a case for analysis. The threshold value is taken as the variable, and the number of polluted routes output by the algorithm is taken as the dependent variable. We control the distance threshold value and pollution threshold value within the interval of [0,300] respectively, and collect the number of pollution routes output by the algorithm. Finally, we draw a line chart (Fig. 12). The abscissa is the value of the threshold, and the ordinate is the number of polluted routes. The blue line shows the relationship between the pollution threshold and the number of pollution routes. The orange line shows the relationship between the distance threshold and the number of polluted routes. According to the figure, when the distance threshold is constant, the number of pollution routes decreases with the increase of the pollution threshold. When the pollution threshold is constant, the number of pollution routes increases with the increase of distance threshold.

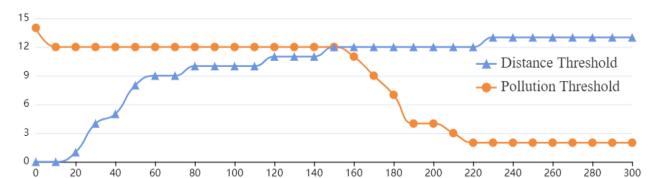


Fig. 12: Sensitivity analysis of algorithm

In the sensitivity analysis of this case, the number of pollution routes output by the algorithm changed within a reasonable range by adjusting the distance threshold and pollution threshold values respectively. This conforms to the principle that the output of the algorithm varies with the input parameters. At the same time, the range of variation

588 is within a reasonable range, which further proves that our
 589 algorithm has strong stability.

590 6.3. User study

591 The visualization system we designed is a tool to feed
 592 back the algorithm output to users in time. Users can further
 593 analyze the backward trajectory results based on the output
 594 of the algorithm, so as to explore the causes of pollution and
 595 explore the propagation law.

596 After investigation, we find that MeteoInfo (Wang,
 597 2014) is a software for effective interaction of backward
 598 trajectory algorithm, which is widely used in the case study
 599 of atmospheric propagation. The software well simulates
 600 the results of HYSPLIT4 with Concentration Weighted Tra-
 601 jectory (CWT) algorithm. However, the software has high
 602 requirements for data and the universality of the algorithm
 603 is not high. In addition, based on the traditional backward
 604 trajectory algorithm, the output result of the software is
 605 single. The software can only provide users with trajectory
 606 map and contribution heat map and lacks the evaluation of
 607 influence factors.

608 For the problem driven visualization system, it is not
 609 feasible to use qualitative indicators to prove the effectiveness
 610 of the system. Therefore, based on above discussion, we
 611 recruited twenty volunteers, including ten students majoring
 612 in environment who have used MeteoInfo software and ten
 613 students majoring in computer science without professional
 614 knowledge in the field of air pollution transmission. In the
 615 survey, they were asked to complete the following three tasks
 616 and explain their experience:

- 617 • Task 1: Find the source station with the most wide-
 618 spreading transport of pollutants;
- 619 • Task 2: Make a knowledge of and compare the basic
 620 information of cities involved in the propagation process;
- 621 • Task 3: Determine the main influencing factor of pollution
 622 process in different stages.

623 After learning about our algorithm and system, volun-
 624 teers began to explore MeteoInfo software and our system
 625 for the above tasks. We summarize the feedback and draw
 626 the following conclusions:

627 For Task 1, all volunteers can clearly find the source
 628 station with the most wide-spreading transport of pollutants
 629 through the path map and the corresponding pollution con-
 630 tribution graphic. However, the results of traditional algo-
 631 rithms are more limited than our algorithm. The traditional
 632 algorithm obtains the only station strictly based on data.
 633 According to parameter feedback of users, the results ob-
 634 tained by our algorithm are more consistent with the reality
 635 and also provide a more comprehensive technical reference
 636 for pollution control. At the same time, four volunteers
 637 mentioned that from the research objectives they selected,
 638 when the pollution involves multiple stations and the color
 639 of the heat map becomes complex, the heavily polluted areas
 640 will interfere with the research of other cities. However,
 641 the reachability matrix thermal diagram and the influence

642 degree sankey chart provided by the visualization system do
 643 not have this problem.

644 For Task 2, we find that the MeteoInfo software can not
 645 understand the basic information of the city specified by the
 646 user. Therefore, volunteers can not complete the comparison
 647 and further exploration of urban information through Me-
 648 teoInfo software. However, through our system, when users
 649 are curious about the information of a city in the pollution
 650 track or try to explore the change trend of stations, the trend
 651 comparison chart provides a comprehensive reference.

652 For Task 3, the HYSPLIT4 with Concentration Weighted
 653 Trajectory (CWT) algorithm adopted by MeteoInfo software
 654 lacks consideration of wind speed and direction. At the same
 655 time, the black box technology used in the software can not
 656 reasonably analyze the influencing factors of pollution for
 657 users. Therefore, all volunteers said that, the software could
 658 not complete Task 3. However, our algorithm comprehen-
 659 sively considers geographical and meteorological factors to
 660 realize distance correction. The pie distribution next to the
 661 station in the pollution propagation path map can make users
 662 clearly understand the distribution of pollution influencing
 663 factors in each path. The distance weight distribution chart
 664 also shows the influencing factors in the time dimension.
 665 The volunteers said they could clearly observe the change
 666 of the influence degree of humidity or wind and guess the
 667 occurrence of rainfall or extreme weather conditions.

668 Based on the above user study, our algorithm and system
 669 have obvious effectiveness in analyzing the causes of pollu-
 670 tion transmission and comparing the trends of different re-
 671 gions. Nevertheless, MeteoInfo software is still a mature and
 672 classic analysis platform, which can solve many backward
 673 trajectory problems. Our system aims to become an auxiliary
 674 tool for algorithm output and user satisfaction to facilitate
 675 subsequent analysis and research more conveniently.

676 7. Conclusion and future work

677 The prevention and control of air pollutants is an impor-
 678 tant issue in environmental governance. In order to assist
 679 meteorologists in process of governance, a topology-based
 680 air pollution tracing algorithm is proposed in this paper.
 681 Combined with Chapman-Enskog binary diffusion theory
 682 and Gaussian diffusion model, the algorithm fully considers
 683 the spatial distance and meteorological factors. Compared
 684 with the traditional pollution propagation models, the advan-
 685 tages of our algorithm are as follows:

- 686 • This paper focuses on the backward trace of air pollutants.
 687 Compared with the forward air pollution propagation
 688 models, this algorithm is more suitable for the accurate
 689 location of pollutant sources. Therefore, the algorithm
 690 is more suitable for the environmental governance. Our
 691 work also provided a new idea for the prevention and
 692 control of air pollutants.
- 693 • By introducing the mathematical method of reachability
 694 matrix, the algorithm steps are reasonable and effective
 695 in the field of mathematical geometry. Topological ge-
 696 ometry takes the positional relationship between objects

- 697 into sufficient consideration. Therefore, it is more suitable
 698 for tracing the source of pollutants between regions and
 699 under global. Our algorithm and system provide technical
 700 support for collaborative governance between regions.
- 701 • The algorithm goes beyond the ordinary Gaussian diffusion
 702 model and considers the transmission factors of pol-
 703 lutants more comprehensively. With full consideration of
 704 geographical and meteorological factors, the accessibility
 705 of pollution transmission path becomes more reasonable
 706 and accurate.
- 707 Combined with the designed air pollution tracing visual-
 708 ization system, this paper conducts the case analysis of
 709 the air pollution situation around 2014. Finally, the tracing
 710 results of this paper are basically consistent with the local
 711 investigation, which verifies that the algorithm provided in
 712 this paper is true and effective. However, there are still some
 713 improvements in the algorithm. For example, we can further
 714 consider the different transmission modes of six pollutants
 715 because of the difference of chemical composition. Due to
 716 the lack of description of extreme weather in data set, we
 717 do not consider different extreme weather when using this
 718 algorithm and system.
- 719 After perfecting the algorithm and system, we will pro-
 720 mote the online release of the system, so that more users can
 721 use our system to trace air pollutants.
- ## 722 References
- 723 Ahmad, M., Chandio, A.A., Solangi, Y.A., Shah, S.A.A., Shahzad, F.,
 724 Rehman, A., Jabeen, G., 2021. Dynamic interactive links among sus-
 725 tainable energy investment, air pollution, and sustainable development
 726 in regional china. Environmental Science and Pollution Research 28,
 727 1502–1518.
- 728 Begum, B.A., Kim, E., Jeong, C.H., Lee, D.W., Hopke, P.K., 2005. Evalua-
 729 tion of the potential source contribution function using the 2002 quebec
 730 forest fire episode. Atmospheric Environment 39, 3719–3724.
- 731 Carvalho, J.C., De Vilhena, M.T.M., 2005. Pollutant dispersion simulation
 732 for low wind speed condition by the ils method. Atmospheric Environ-
 733 ment 39, 6282–6288.
- 734 Chapman, S., Cowling, T.G., 1990. The mathematical theory of non-
 735 uniform gases: an account of the kinetic theory of viscosity, thermal
 736 conduction and diffusion in gases. Cambridge university press.
- 737 Cotrozzi, L., 2020. Leaf demography and growth analysis to assess the
 738 impact of air pollution on plants: A case study on alfalfa exposed to
 739 a gradient of sulphur dioxide concentrations. Atmospheric Pollution
 740 Research 11, 186–192.
- 741 Demuzere, M., Trigo, R., Vila-Guerau de Arellano, J., Van Lipzig, N., 2008.
 742 The impact of weather and atmospheric circulation on o3 and pm10
 743 levels at a mid-latitude site. Atmos. Chem. Phys. Discuss 8, 21037–
 744 21088.
- 745 Deng, Z., Weng, D., Chen, J., Liu, R., Wang, Z., Bao, J., Zheng, Y., Wu,
 746 Y., 2019. Airvis: Visual analytics of air pollution propagation. IEEE
 747 transactions on visualization and computer graphics 26, 800–810.
- 748 Deng, Z., Weng, D., Liang, Y., Bao, J., Zheng, Y., Schreck, T., Xu, M., Wu,
 749 Y., 2021. Visual cascade analytics of large-scale spatiotemporal data.
 750 IEEE Transactions on Visualization and Computer Graphics .
- 751 Fan, H., Zhao, C., Yang, Y., 2020. A comprehensive analysis of the spatio-
 752 temporal variation of urban air pollution in china during 2014–2018.
 753 Atmospheric Environment 220, 117066.
- 754 Guo, F., Gu, T., Chen, W., Wu, F., Wang, Q., Shi, L., Qu, H., 2019.
 755 Visual exploration of air quality data with a time-correlation-partitioning
 tree based on information theory. ACM Transactions on Interactive
 Intelligent Systems (TiiS) 9, 1–23.
- Huang, Liu, Sheng, 2018. On adjoint method based atmospheric emission
 source tracing (in chinese). Chin Sci Bull 63, 1594–1605.
- Jin, L., Godri Pollitt, K.J., Liew, Z., Rosen Vollmar, A.K., Vasiliou, V.,
 Johnson, C.H., Zhang, Y., 2021. Use of untargeted metabolomics to
 explore the air pollution-related disease continuum. Current Environmental
 Health Reports 8, 7–22.
- Kim, S.S., Kim, Y.K., Kang, Y.M., 2021. Aorm: Fast incremental arbitrary-
 order reachability matrix computation for massive graphs. IEEE Access
 9, 69539–69558.
- Kong, L., Tang, X., Zhu, J., Wang, Z., Li, J., Wu, H., Wu, Q., Chen, H., Zhu,
 L., Wang, W., et al., 2021. A 6-year-long (2013–2018) high-resolution
 air quality reanalysis dataset in china based on the assimilation of surface
 observations from cnemc. Earth System Science Data 13, 529–570.
- Li, J., Chen, S., Zhang, K., Andrienko, G., Andrienko, N., 2018. Cope:
 Interactive exploration of co-occurrence patterns in spatial time series.
 IEEE transactions on visualization and computer graphics 25, 2554–
 2567.
- Li, J., Fan, Z., Deng, M., 2017. A method of spatial interpolation of air
 pollution concentration considering wind direction and speed. J. Geo
 Inf. Sci 19, 382–389.
- Li, L., Cai, J., Zhou, M., 2015. Potential source contribution analysis of the
 particulate matters in shanghai during the heavy haze episode in eastern
 and middle china in december, 2013. Huan Jing ke Xue= Huanjing
 Kexue 36, 2327–2336.
- Liu, D., Veeramachaneni, K., Geiger, A., Li, V.O., Qu, H., 2021. Aqeyes:
 visual analytics for anomaly detection and examination of air quality
 data. arXiv preprint arXiv:2103.12910 .
- Pandey, A., Brauer, M., Cropper, M.L., Balakrishnan, K., Mathur, P., Dey,
 S., Turkoglu, B., Kumar, G.A., Khare, M., Beig, G., et al., 2021. Health
 and economic impact of air pollution in the states of india: the global
 burden of disease study 2019. The Lancet Planetary Health 5, e25–e38.
- Parajuli, R.P., Shin, H.H., Maquiling, A., Smith-Doiron, M., 2021. Multi-
 pollutant urban study on acute respiratory hospitalization and mortality
 attributable to ambient air pollution in canada for 2001–2012. Atmo-
 spheric Pollution Research 12, 101234.
- Qu, D., Lin, X., Ren, K., Liu, Q., Zhang, H., 2020. Airexplorer: Visual
 exploration of air quality data based on time-series querying. Journal of
 Visualization 23, 1129–1145.
- Qu, H., Chan, W.Y., Xu, A., Chung, K.L., Lau, K.H., Guo, P., 2007. Visual
 analysis of the air pollution problem in hong kong. IEEE Transactions
 on visualization and Computer Graphics 13, 1408–1415.
- Ren, K., Wu, Y., Zhang, H., Fu, J., Qu, D., Lin, X., 2020. Visual analytics
 of air pollution propagation through dynamic network analysis. IEEE
 Access 8, 205289–205306.
- Scarpato, N., Cilia, N.D., Romano, M., 2019. Reachability matrix ontology:
 a cybersecurity ontology. Applied Artificial Intelligence 33, 643–655.
- Sharovsky, R., César, L., Ramires, J., 2004. Temperature, air pollution,
 and mortality from myocardial infarction in sao paulo, brazil. Brazilian
 journal of medical and biological research 37, 1651–1657.
- Stein, A., Draxler, R., Rolph, G., Stunder, B., 2015. B., cohen, md, and
 ngan, f.: Noaa's hysplit atmospheric transport and dispersion modeling
 system, b. Am. Meteorol. Soc 96, 2059–2077.
- Tian, Li, Cheng, 2021. Visual analysis system for fine-grained inline
 relationship of air quality data. Journal of Computer-Aided Design
 Computer Graphics 33, 11.
- Vo, L.H.T., Yoneda, M., Nghiem, T.D., Shimada, Y., Van, D.A., Nguyen,
 T.H.T., Nguyen, T.T., 2022. Indoor pm0. 1 and pm2. 5 in hanoi: Chemical
 characterization, source identification, and health risk assessment.
 Atmospheric Pollution Research , 101324.
- Wang, Y.Q., 2014. Meteoinfo: Gis software for meteorological data
 visualization and analysis. Meteorological Applications 21, 360–368.
- Xu, X., Hong, Y., 2012. Matrix expression and reachability analysis of finite
 automata. Journal of Control Theory and Applications 10, 210–215.
- Zhou, Z., Ye, Z., Liu, Y., Liu, F., Tao, Y., Su, W., 2017. Visual analytics
 for spatial clusters of air-quality data. IEEE computer graphics and
 applications 37, 98–105.