

# Topological Characterization of Haze Episodes Using Persistent Homology

Nur Fariha Syaquina Zulkepli<sup>1\*</sup>, Mohd Salmi Md Noorani<sup>1</sup>, Fatimah Abdul Razak<sup>1</sup>, Munira Ismail<sup>1</sup> and Mohd Almie Alias<sup>1</sup>

<sup>1</sup> School of Mathematical Sciences, Faculty of Science and Technology, Universiti Kebangsaan Malaysia 43600 Bangi, Selangor, Malaysia.

## Abstract

Haze phenomenon is one of the major environmental issues that have continuously vexed countries worldwide, including Malaysia, for the last three decades. Therefore, this study aims to investigate the differences between topological features of months with and without haze episodes at air quality monitoring stations located in the areas of Jerantut, Klang, Petaling Jaya, and Shah Alam. This is achieved by opting for persistent homology, which is a method in topological data analysis (TDA) that analyses data in a qualitative (topological) sense via a focus on topological features, such as connected components and holes of the data. Topological features from particulate matter (PM<sub>10</sub>) data with summary statistics are subsequently utilized to summarize them. The results have consequently shown drastic changes present in the summary statistics of the lifetimes of topological features during haze episodes. These characteristics have been consistently observed in each air quality monitoring station involved in this study. Thus, this finding highlights the potential application for the development of an early warning system for haze detection based on the topological approach.

**Keywords:** Haze; Particulate matter; Persistent homology; Time delay embedding; Topological data analysis.

---

\* Corresponding author. Tel: +6012-543-2006

E-mail address: farihasyaqina@yahoo.com

## 30 INTRODUCTION

31

32 Haze phenomenon is one of the major environmental issues that occur continuously in  
33 countries across the world; Malaysia included (Wen *et al.*, 2016). The country is also afflicted,  
34 with the Klang Valley region being one of the areas that are particularly affected by such  
35 phenomenon (Latif *et al.*, 2018). The severity of this issue can be attributed to the large-scale  
36 forest and plantation fires that occur in Indonesia's Sumatra Island. Nevertheless, it is worsened  
37 further by local emissions from domestic industries, vehicle usage, and open burning activities  
38 (Afroz *et al.*, 2003; Wen *et al.*, 2016). In Malaysia, the chronological history of haze episodes can  
39 be underlined with the severe incidents recorded in the year 2005, 2013, 2014 and 2015  
40 accordingly (DOE, 2018a). This issue has continually emerged year after year in Malaysia and  
41 the neighboring countries of Brunei and Singapore, whereby the prolonged duration of the  
42 problem has spurred the study into investigating haze episodes.

43 Particulate matter ( $PM_{10}$ ) are small particles floating around in the air in the form of smoke,  
44 dirt, and dust that originate from factories, vehicles and farming activities (Schwartz *et al.*, 1996;  
45 Payus *et al.*, 2013). During haze episodes,  $PM_{10}$  is typically highlighted as a major air pollutant  
46 in the Southeast Asian regions, particularly the Klang Valley, Malaysia (Afroz *et al.*, 2003; Azmi  
47 *et al.*, 2010). Acknowledged as a harmful pollutant, inhalation of the material leads to  
48 diminishing lung function and causes various respiratory diseases, especially acute exacerbation  
49 of asthma (Schwartz *et al.*, 1996). According to the Malaysian Ambient Air Quality Guidelines

50 (MAAQG), an average of  $PM_{10}$  concentration over 24-hour that exceeds  $150 \mu\text{g m}^{-3}$  is  
51 considered to be unhealthy for human health. Therefore, a standard technique used to identify  
52 haze phenomenon is by observing the concentration of  $PM_{10}$ , whereby haze emergency is  
53 declared when it reaches the emergency level ( $> 500 \mu\text{g m}^{-3}$ ) (DOE, 2018b).

54 Several studies have explored into a comparison of  $PM_{10}$  concentration with MAAQG  
55 standards, undertaken by Azmi *et al.* (2010), Abdullah *et al.* (2012), Rahman *et al.* (2015) and  
56 Ling *et al.* (2010) respectively. They have consequently concluded that the concentration in the  
57 air quality monitoring stations located in Klang Valley exceeded the acceptable level  
58 recommended by the MAAQG during haze episodes. Various differences have also been revealed  
59 between the concentration of  $PM_{10}$  according to the locations of air quality monitoring stations,  
60 encompassing rural, urban and industrial areas respectively. Moreover, Azmi *et al.* (2010) and  
61 Abdullah *et al.* (2012) have also revealed that urban areas logged higher  $PM_{10}$  concentration  
62 compared to rural areas. A study by Yusof *et al.* (2010) has provided further analysis using a  
63 statistical model, Lognormal and Weibull distributions to investigate the relationship between  
64  $PM_{10}$  concentration and monsoon season. Similarly, other methods like chemometric analysis  
65 (Azid *et al.*, 2015), fuzzy comprehensive evaluation method (Zhao *et al.*, 2010) and chaotic  
66 approach (Hamid and Noorani, 2014) have also been utilized in assessing data on air pollutants.  
67 Furthermore, haze detection can also be found to be fast gaining scholarly traction. Previously,

68 works focusing on the topic are geared towards analyzing satellite imagery data using remote  
69 sensing methods (Makarau, 2014). Meanwhile, recent research is specifically related to haze  
70 periods and focused on analyzing particulate matter concentration to study its morphology during  
71 haze episodes (Zeb *et al.*, 2018). Yu *et al.* (2018) have also contributed by investigating human  
72 health risk in an indoor and outdoor environment that is exposed to the phenomenon, whereas  
73 another study has looked into mitigating severe urban haze episodes (Sharma and  
74 Balasubramanian, 2018). It should be noted that these literature has analyzed  $PM_{10}$  in a  
75 quantitative manner, whereas to the best of our knowledge, no effort has been expended on its  
76 qualitative aspects. Therefore, this paper is addressing the research gap by providing an analysis  
77 of the qualitative aspects and structures of  $PM_{10}$ , particularly on their topological features via  
78 persistent homology.

79 Persistent homology is a relatively new method that is robust under perturbations of input data,  
80 independent of coordinates and dimensions alike, and offers a solid representation of qualitative  
81 features of the input data (Otter *et al.*, 2017). These particular features of the data are captured  
82 accordingly as specific parameters change. As the approach is fundamentally based on the  
83 mathematical field of topology, the qualitative features in question encompass topological  
84 features like connected components, holes, voids and more. The qualitative approach of persistent  
85 homology is also particularly useful in handling complexity of data due to noise, high

86 dimensionality or incomplete structure. This poses great challenges to researchers dealing with  
87 real world data since data cleaning is required in order to remove the noise and which might  
88 result in information lost during the process (Jorquera *et al.*, 2000; Elangasinghe *et al.*, 2014). By  
89 contrast, persistent homology retains all information from data. The noise that may occur in  
90 multiple scale-levels is filtered out by persistent homology and significant features are captured  
91 (Ghrist, 2008). Furthermore, its robustness has led other researchers to explore the method further  
92 in various fields, such as Gidea and Katz (2018) effort on the capability of the holes (1-  
93 dimensional features) to detect financial crisis in stock market data. Meanwhile, Emrani *et al.*  
94 (2014) have investigated wheeze signal detection via the persistency of topological features  
95 between wheeze and non-wheeze signals. Moreover, the summary statistics of topological  
96 features undertaken by Mittal and Gupta (2017) has shown that persistent homology is usable in  
97 early detection of bifurcations and chaos in complex system. The explorations on persistent  
98 homology have also been tackled in various fields, encompassing the classification of breast  
99 cancer (Dewoskin *et al.*, 2010), viral evolution (Chan *et al.*, 2013) and protein structure (Gameiro  
100 *et al.*, 2015). A good and comprehensive review regarding the current applications of persistent  
101 homology may be sourced from Otter *et al.* (2017). To the best of our knowledge, there is a  
102 negligible amount of research output on environmental issues by persistent homology to date.  
103 Therefore, this study is conducted to fill the research gap by exploring the effectiveness of

104 persistent homology in characterizing and detecting one of the pressing environmental issues,  
105 namely haze phenomenon.

106 This paper intends to investigate haze episodes using persistent homology by analyzing  
107 the topological features of  $PM_{10}$  data in Malaysia. The proposed technique has been applied to  
108 daily average  $PM_{10}$  data from four air quality monitoring stations of Jerantut, Klang, Petaling  
109 Jaya and Shah Alam from the year 2000 until 2015. The objective of this study is to investigate  
110 topological features for months with and without haze episodes, thus rendering the data on  $PM_{10}$   
111 being partitioned according to the respective month. The topological features are then extracted  
112 for connected components (0-dimensional features) and holes (1-dimensional features), with the  
113 resulting information represented in persistence diagrams. Next, the summary statistics of the  
114 topological features, like an average of all lifetimes of the connected components and maximum  
115 lifetimes of all holes, are calculated for each persistence diagram. In this paper, the content has  
116 been organized accordingly, whereby the subsequent section consists of a discussion on data  
117 involved in this study. Then, the data transformation processes and the method of persistent  
118 homology will be explained in next two sections. Extensive summary statistics of topological  
119 features would be introduced in the section before its resulting outcomes is discussed.

120

121

122

## 123 DATA

124

125 The daily average data for PM<sub>10</sub> encompassing the year 2000 to 2015 for four air quality  
126 monitoring stations (i.e. Jerantut, Klang, Petaling Jaya, and Shah Alam) has been obtained from  
127 the Department of Environment (DOE), Malaysia. Klang and Shah Alam are categorized as urban  
128 areas, while Petaling Jaya and Jerantut are categorized as respectively of the industrial and rural  
129 areas, accordingly (DOE, 2018c). Due to its rural location, Jerantut, in particular, has been chosen  
130 as the background station for comparative purpose (Azmi *et al.*, 2010; Banan *et al.*, 2013).  
131 Missing data has been treated with the mean substitution method (Pigott, 2001), whereas the  
132 MAAQG standards are necessary to represent the safety level that causes no adverse health  
133 effects to human. Therefore, this study has adhered to the safe level as recommended by  
134 MAAQG to act as the benchmark for air quality level of each station.

135 Fig. 1 shows the time series of the daily average for PM<sub>10</sub> from 1 January 2000 until 31  
136 December 2015, and the MAAQG for a 24-hour average PM<sub>10</sub> at 150 µg m<sup>-3</sup>. The presence of  
137 several peaks (rectangles) is indicative of the values of PM<sub>10</sub> concentration exceeding the  
138 MAAQG during haze episodes that have occurred in August 2005, June 2013, March 2014,  
139 September 2015, and October 2015 accordingly (DOE, 2018a). Thus, these months have been  
140 highlighted as the main focus of this research. It is worth noting that based on the chronology of  
141 haze episodes (DOE, 2018a), the selected months have been reported of experiencing severe haze.

142 The descriptive statistics of the months are shown in Table 1 accordingly, with Klang, Petaling  
143 Jaya and Shah Alam stations showing higher concentration of  $PM_{10}$  compared to Jerantut station  
144 during the haze episodes. This is attributable to the different locations of the air quality  
145 monitoring stations, with Klang and Shah Alam being located in an urban area, Petaling Jaya in  
146 an industrial area and Jerantut (background station) in a rural area (DOE, 2018c).

## 147 148 **TIME DELAY EMBEDDING**

149  
150 Prior to the implementation of persistent homology, the time series data sets must be  
151 transformed into point cloud data, which is achieved via the Takens method. The higher  
152 dimensional data sets will allow us to look into higher dimensional topological features, such as  
153 holes and voids. Basically topological features are features that are invariant after deformations  
154 such as stretching, splitting and cutting (Hatcher, 2002; Edelsbrunner and Harer, 2010; Ghrist,  
155 2014). The idea of using a combination of Takens method and persistent homology is discussed  
156 in the work of Perea and Harer (2015) and references therein. The Takens method (Takens, 1981)  
157 stated that a time series  $x_0, x_1, \dots, x_{n-1}$  can be reconstructed in a phase space of dimension  $m$ ,  
158 where each point in the phase space is given by the vector  $x_n(m, \tau) = x_n, x_{n+\tau}, \dots, x_{n+(m-1)\tau}$ ,  $\tau$  is  
159 the time delay and  $m$  is the embedding dimension. The two parameters  $\tau$  and  $m$  require careful  
160 selection to ensure clear extraction of the desired topological features. In this study, a comparison  
161 made between the months has indicated that the different settings of time delay and embedding



162 dimension will affect the results. The two parameters have been fixed as  $\tau = 1$  and  $m = 3$ . The  
163 time delay is chosen as trivial time delay, whereas the embedding dimension chosen is 3 as the 1-  
164 dimensional topological features are best visualized in dimension 3 throughout the study. This  
165 decision is supported by Umeda (2017), whose work has fixed both parameters as  $\tau = 1$  and  
166  $m = 3$ , while others like Khasawneh *et al.* (2018) and Khasawneh and Munch (2014) chose  
167  $m = 3$  to visualize the 1-dimensional features. Similarly, previous studies by Pereira and Mello  
168 (2015) and Maletić *et al.* (2016) have also fixed  $m$ , whereas in other areas of research (not  
169 related to persistent homology), Sivakumar (2003) and Sivakumar (2002) each have chosen  $\tau = 1$   
170 to achieve better results for their respective research.

## 171 172 **PERSISTENT HOMOLOGY**

173  
174 Computations of simplicial complexes are compulsory in extracting topological features  
175 from point clouds data. The 0-simplices represent vertices or points, 1-simplices represent edges  
176 or lines, 2-simplices represent triangles and 3-simplices represent tetrahedra, and so on. A  
177 simplicial complex is built by a combination of these simplices (See Fig. 2) accordingly, whereby  
178 its complex formation is from data points and thus indicating the dependency of topological  
179 features on a scaling parameter (filtration value),  $\varepsilon$ . Fig. 3(a) in particular shows an example of a  
180 simplicial complex formation. The formation commenced at  $\varepsilon = 0$ , where the four 0-simplices  
181 formed represents four points in a point cloud. As the value  $\varepsilon$  increases to 0.5, four circles are

182 formed with each 0-simplex acting as the centre and  $\varepsilon$  as the radius of the circles. The circles  
183 keep growing as the value  $\varepsilon$  increases until any pair of the circles intersects each other. From this  
184 intersection, a 1-simplex is formed by connecting two 0-simplices by a line (edge). As seen in Fig.  
185 3(a), four 1-simplices have been formed at  $\varepsilon = 1.4$ , and the stage is marked with the formation of  
186 a simplicial complex via the combination of the 0-simplices and 1-simplices. The  $\varepsilon$  value then  
187 further increased to 2, with more circles intersecting each other and resulting in the appearance of  
188 two 2-simplices (triangles). At this stage, a new simplicial complex is formed (see Fig. 3(a)). It  
189 should be noted that the simplicial complex at  $\varepsilon = 1.4$  is contained in the simplicial complex at  
190  $\varepsilon = 2$ , which is also known as filtered simplicial complexes. In this work, the constructions of  
191 simplicial complexes are done using Vietoris-Rips simplicial complex, or otherwise also known  
192 as Rips complex. Rips complex is a set of  $k$ -simplices, such that the distance of any two points in  
193  $k$ -simplices is less than or equal to  $2\varepsilon$  (Hatcher, 2002; Edelsbrunner and Harer, 2010; Ghrist,  
194 2014).

195 The birth and death points of topological features are captured depending on the  
196 formation of simplicial complexes and subsequently recorded in a diagram known as barcode  
197 (Fig. 3(b)). Each feature is represented by a horizontal line in the barcode, with the left end point  
198 of the line being the birth point and the right end point of the line considered as the death point.  
199 One can recognize any noise in the barcode by looking at the short lines, while the significant

200 features are signified by the long lines (Ghrist, 2008). Furthermore, the black lines in the barcode  
201 represent the connected components, whereas the red line represents the hole (see Fig. 3(b)).  
202 Moreover, persistence diagram (Fig. 3(c)) is yet another representation of barcode that serves to  
203 summarize the birth and death of topological features in  $\mathbb{R}^2$ . The  $x$ -axis and  $y$ -axis of the  
204 persistence diagram is the representation of the birth and death points respectively for each of the  
205 topological features. Besides, a point  $(b, d)$  with multiplicity  $q$  in a persistence diagram  
206 represents  $q$  features that have same birth,  $b$  and death,  $d$  points. Another feature that persists  
207 longer through the filtration stage is located way above the diagonal line in the persistence  
208 diagram. Additionally, the persistency of a feature is measured by the difference between the  
209 death and birth points,  $d - b$  and subsequently known as the lifetime of the feature (Otter *et al.*,  
210 2017).

211 Fig. 3 shows an example of connected components (0-dimensional) and hole (1-  
212 dimensional) extracted based on the formation of simplicial complexes. At the starting point, four  
213 connected components have appeared with a birth value at  $\varepsilon = 0$ , and as the filtration value  
214 increases the features have remained until 1-simplices appeared at  $\varepsilon = 1.4$ . It should be noted that  
215 the four connected components have collapsed into one at  $\varepsilon = 1.4$ . Interestingly, a new feature  
216 has appeared at this stage, represented by a hole appearance. As the values of  $\varepsilon$  have continued  
217 to increase to 2, two 2-simplices are subsequently formed and closed the hole. This causes the

218 death of the hole and renders only one connected component to persist by the end of the filtration.  
219 Persistence diagram in Fig. 3(c) is another representation of these topological features,  
220 simplifying the barcode in Fig. 3(b).

## 221 222 **SUMMARY STATISTICS OF TOPOLOGICAL FEATURES** 223

224 A persistence diagram consists of  $k$ -dimensional features with 0-dimensional features  
225 representing the connected components, 1-dimensional features representing holes, and 2-  
226 dimensional features representing voids etc. Meanwhile, a persistence diagram  $\omega_m$  consists of  $n$   
227 features  $\omega_{m_i} = (b_i, d_i)$  with  $b_i$  and  $d_i$  ( $i = 1, 2, \dots, n$ ) indicating their birth and death points  
228 respectively. All of the features are summarized using summary statistics and described below:  
229 The first summary statistic is the sum of all lifetimes (Eq. (1)) (Pereira and Mello, 2015) of  $k$ -  
230 dimensional features. If the value of the sum is close to zero, the persistence diagram has  
231 practically short-lived features for each particular dimension,  $k$ .

$$232 \quad \text{sum}_k = \sum_{i=1}^n (d_i - b_i). \quad (1)$$

233  
234 The second summary statistic is the average of all lifetimes (Pereira and Mello, 2015; Mittal and  
235 Gupta, 2017) of  $k$ -dimensional features, as described in Eq. (2). Likewise, a small value of  
236 average indicates that for a particular dimension,  $k$  of the data set has mostly short-lived features  
237 and vice versa.

238

$$\text{avg}_k = \frac{\sum_{i=1}^n (d_i - b_i)}{n}. \quad (2)$$

239 Next, the third summary statistic is the maximum of all lifetimes (Pereira and Mello, 2015; Mittal  
240 and Gupta, 2017) of  $k$ -dimensional features. For each dimension  $k$ , the value of the maximum  
241 lifetimes of all features is defined as

242

$$\text{max}_k = \max_{\omega_{m_i}} (d_i - b_i), \quad (3)$$

243

244 which will determine the most significant  $k$ -dimensional feature.

245 In this study, the summary statistics of 0-dimensional and 1-dimensional features have been  
246 calculated using Eqs. (1-3). The changes of values in each summary statistics have been observed  
247 to track the evolution of the topological features for months with and without haze accordingly.

248 **Remark.** We acknowledge that according to Cohen-Steiner *et al.* (2007) and Cohen-Steiner *et al.*  
249 (2010), the summary statistic,  $\text{sum}_k$  is stable whereas the stability of the  $\text{avg}_k$  is still in doubt.

250 However, there may be merit in taking into account  $\text{avg}_k$  as considered by Pereira and Mello  
251 (2015) and Mittal and Gupta (2017) but in general it should be done with care. Based on our  
252 observations using our data sets, small variation of the input data leads to small variation of  $\text{avg}_k$

253 values.

254

255

## 256 RESULTS

257

258 This study has analyzed the time series of daily average for  $PM_{10}$  that is partitioned  
259 according to month for four air quality monitoring stations located in Jerantut, Klang, Petaling  
260 Jaya, and Shah Alam for 16 years (2000-2015). The time delay embedding with time delay  $\tau = 1$   
261 and embedding dimension  $m = 3$  have been applied for each month of data on  $PM_{10}$ , allowing  
262 each month to generate a point cloud data. Persistent homology is then applied to each point  
263 cloud with maximum filtration value,  $\varepsilon_{\max} = 700$ . The evolution of topological features based on  
264 Rips complexes is next observed in the range of filtration value,  $\varepsilon = [0, 700]$ , starting from the  
265 simplest form which is under-approximation until it displays over-estimation (i.e. one large Rips  
266 complex). The computation of persistent homology in this study has been completed using the R-  
267 package TDA (Fasy *et al.*, 2017).

268 All extracted topological features have been represented in persistence diagrams,  
269 following which the summary statistics for each persistence diagrams have been calculated. Fig.  
270 4 has displayed a comparison between persistence diagrams generated in August 2005 where  
271 severe haze has been reported (DOE, 2018a), in contrast with December 2005 where no haze has  
272 occurred. The point clouds generated using time delay embedding which respect to the  
273 persistence diagrams are also shown in Fig. 4. The month of December has been selected due to  
274 the lesser likelihood for haze to occur during the north-east monsoon (November to March)

275 season, as Malaysia receives more rainfall during the period and increases the removal rate of  
276  $PM_{10}$  (Yusof *et al.*, 2010).

277 From the persistence diagrams shown in Fig. 4, the topological features (i.e. black dots  
278 represent connected components, red triangles represent holes) for a month with haze episode  
279 recorded (i.e. August 2005) is spread away from the origin. In contrast, a month without haze (i.e.  
280 December 2005) has shown that the features accumulate close to the origin. For each station, a  
281 hole is present among the holes (red triangles) as in Fig. 4(a), which is located the furthest from  
282 the origin and the highest value of birth and death points compared to other holes. It should be  
283 noted the corresponding feature in Jerantut is not too far from the origin compared to other  
284 stations, implying that the stations in Klang, Shah Alam, and Petaling Jaya have experienced haze  
285 of higher severity compared to Jerantut. Summary statistics are calculated to summarize the  
286 persistence diagrams, whereby the resulting outcomes are shown in Fig. 5 and Fig. 6.

287 An observation of the various peaks in Fig. 5(a) has revealed drastic increments of the  
288 sum of all lifetimes,  $sum_0$  for connected components (0-dimensional features) during haze  
289 episodes (rectangles). This is indicative of the connected components persisting longer during  
290 haze as the filtration value varied compared to the normal months. This is further elucidated by  
291 the average of all lifetimes,  $avg_0$  as shown in Fig. 5(b), which suggests that the persistence  
292 diagrams for months with haze consist of mostly long-lived features, hence producing the peaks.

293 Fig. 5(a) and 5(b) also show that the sum and average of the lifetimes of connected components  
294 extracted in the background station, Jerantut being the lowest compared to other stations. It  
295 should be noted that the maximum of all lifetimes are not calculated for connected components as  
296 the formation of simplicial complexes (Rips complexes) for  $PM_{10}$  persists until it becomes a  
297 large simplicial complex with birth,  $b = 0$  and death,  $d = 700$  (since maximum filtration value,  
298  $\varepsilon_{\max} = 700$ ) resulting in the same values of maximum for all lifetimes ( $\max_0 = 700$ ) according  
299 to all selected months.

300 Furthermore, the summary statistics for holes (1-dimensional features) are shown in Fig. 6  
301 and revealed drastic changes in the values of the sum, average and maximum of all lifetimes  
302 denoted as  $\text{sum}_1$ ,  $\text{avg}_1$  and  $\text{max}_1$  respectively (refer Fig. 6(a-c)) during haze episodes (rectangles).  
303 As seen in Fig. 6(c), the peaks (in rectangles) are the values of a maximum of the lifetimes of all  
304 holes which imply that there is a hole with its lifetime higher compared to other holes in the  
305 month that experienced severe haze. Hence, Fig. 5 and Fig. 6 have revealed the behavior of the  
306 topological features extracted by persistent homology to be consistent with real haze phenomenon.  
307 When the haze occurred, persistent homology has identified this occurrence by showing the  
308 strong rise of the lifetimes for related topological features.

309 The rectangles in Fig. 5 and Fig. 6 indicate the particular month recorded with severe haze  
310 in the selected years, with a clear and drastic increase in the values of summary statistics in the



311 month. These characteristics have also been shown by each station involved, substantiating the  
312 consistency of the topological characterization for months with and without haze accordingly.  
313 Overall, the values (i.e. summation, average, maximum) of all lifetimes for the topological  
314 features of the background station (Jerantut) are the lowest compared with the others. Thus, it is  
315 clear that the topological features extracted via persistent homology are capable of distinguishing  
316 the months with severe haze and without haze. While the existing methods (Azmi *et al.*, 2010;  
317 Ling *et al.*, 2010; Abdullah *et al.*, 2012; Rahman *et al.*, 2015) have served to directly quantify  
318  $PM_{10}$  concentration without providing a qualitative understanding, this paper has successfully  
319 filled in research gap by revealing the changes of topological features seen between months with  
320 and without haze accordingly.

## 321 322 **CONCLUSIONS**

323  
324 This study has opted for the method from topological data analysis (TDA) known as  
325 persistent homology so as to qualitatively analyze the topological features embedded within the  
326 data on the concentration of a major pollutant,  $PM_{10}$ . The changes of topological features have  
327 been found in the months that recorded severe haze episodes, which have been consistently  
328 shown in each station involved in this study. Despite existing methods is capable of quantifying  
329 and analyzing the concentration of  $PM_{10}$  directly in determining the months affected during haze  
330 episodes, this study has proposed a different approach for analysis based on the topological

331 viewpoint. The topological characterization of haze episodes complement the existing methods  
332 by providing its qualitative structures which is robust to the effect of noise. The summary  
333 statistics of the topological features (i.e. connected components and holes) have indicated a  
334 strong rising of the sum, average and maximum of all lifetimes of the topological features during  
335 haze episodes. Furthermore, topological features  $PM_{10}$  extracted from different locations of air  
336 quality monitoring stations have succeeded in characterizing Jerantut as the background station to  
337 have a low level of air pollution. This is seen by the lowest values of summary statistics  
338 compared to other stations (i.e. Klang, Petaling Jaya, and Shah Alam). Thus, this study proposed  
339 a new approach to analyze haze episodes using a qualitative approach, specifically by extracting  
340 topological features that can characterize haze episodes. This is done by observing the summary  
341 statistics of such topological features. Hence, this finding has displayed its potential as an  
342 application for the development of an early warning system for haze detection, based on the  
343 topological approach. We particularly believe that early warning signals can be determined by  
344 refining the way data is analyzed and the progressive efforts spearheaded by this work is making  
345 waves in the particular area. We also suggest future research in addressing the stability of  
346 summary statistic,  $avg_k$  since the referee has pointed out that it's stability is still unknown.

347

348

349

350 **ACKNOWLEDGMENTS**

351

352           The authors would like to express their utmost gratitude to Universiti Kebangsaan  
353 Malaysia for the Research University Grant (DIP-2017-011, GGPM-2016-001), Ministry of  
354 Education Malaysia Grant FRGS/1/2017/STG06/UKM/01/1, and to Department of the  
355 Environment (DOE) for providing the air quality data. We also thank the reviewers for their  
356 helpful and critical comments.

357

358

359

360

361

362

363

364

365

366

367

368

369

370 **REFERENCES**

371

372 Abdullah, A.M., Samah, M.A.A. and Jun, T.Y. (2012). An overview of the air pollution trend in  
373 Klang Valley, Malaysia. *Open Environmental Sciences*. 6(1): 13–19.

374 Afroz, R., Hassan, M.N. and Ibrahim, N.A. (2003). Review of air pollution and health impacts in  
375 Malaysia. *Environ. Res.* 92(2): 71–77.

376 Azid, A., Juahir, H., Ezani, E., Toriman, M.E., Endut, A., Rahman, M.N.A., Yunus, K., Nordin,  
377 M., Kamarudin, M.K.A., Hasnam, C.N.C., Saudi, A.S.M. and Umar, R. (2015). Identification  
378 source of variation on regional impact of air quality pattern using chemometric. *Aerosol Air*  
379 *Qual. Res.* 15: 1545–1558.

380 Azmi, S.Z., Latif, M.T., Ismail, A.S., Juneng, L. and Jemain, A.A. (2010). Trend and status of air  
381 quality at three different monitoring stations in the Klang Valley, Malaysia. *Air Qual. Atmos.*  
382 *Health.* 3(1): 53–64.

383 Banan, N., Latif, M.T., Juneng, L. and Ahamad, F. (2013). Characteristics of surface ozone  
384 concentrations at stations with different backgrounds in the Malaysian Peninsula. *Aerosol Air*  
385 *Qual. Res.* 13(3): 1090-1106.

386 Chan, J.M., Carlsson, G. and Rabadan, R. (2013). Topology of viral evolution. Proceedings of the  
387 National Academy of Sciences. 110(46): 18566–18571.

388 Cohen-Steiner, D., Edelsbrunner, H., and Harer, J. (2007). Stability of persistence

389 diagrams. *Discrete Comput. Geom.* 37(1): 103-120.

390 Cohen-Steiner, D., Edelsbrunner, H., Harer, J. and Mileyko, Y. (2010). Lipschitz functions have  
391  $L_p$ -stable persistence. *Found. Comput. Math.* 10(2): 127-139.

392 Dewoskin, D., Climent, J., Cruz-White, I., Vazquez, M., Park, C. and Arsuaga J. (2010).  
393 Applications of computational homology to the analysis of treatment response in breast cancer  
394 patients. *Topol. Appl.* 157(1): 157–164.

395 DOE (2018a). Chronology of haze episodes in Malaysia. Department Of Environment Malaysia.  
396 [https://www.doe.gov.my/portalv1/en/info-umum/info-kualiti-udara/kronologi-episod-jerebu-](https://www.doe.gov.my/portalv1/en/info-umum/info-kualiti-udara/kronologi-episod-jerebu-di-malaysia/319123)  
397 [di-malaysia/319123](https://www.doe.gov.my/portalv1/en/info-umum/info-kualiti-udara/kronologi-episod-jerebu-di-malaysia/319123). Last Access: 10 July 2018.

398 DOE (2018b). A guide to air pollutant index (API) in Malaysia. Department Of Environment  
399 Malaysia. Enviro Knowledge Centre. <https://enviro.doe.gov.my/>. Last Access: 10 August 2018.

400 DOE (2018c). Environmental quality report 2015. Department Of Environment Malaysia. Enviro  
401 Knowledge Centre. <https://enviro.doe.gov.my/> . Last Access: 10 July 2018.

402 Edelsbrunner, H. and Harer, J. (2010). *Computational topology: an introduction*. American  
403 Mathematical Society, Providence.

404 Elangasinghe, M.A., Singhal, N., Dirks, K.N., and Salmond, J.A. (2014). Development of an  
405 ANN-based air pollution forecasting system with explicit knowledge through sensitivity  
406 analysis. *Atmospheric pollution research.* 5(4): 696-708.

407 Emrani, S., Gentimis, T. and Krim, H. (2014). Persistent homology of delay embeddings and its  
408 application to wheeze detection. *IEEE Signal Process. Lett.* 21(4): 459–463.

409 Fasy, B.T., Kim, J., Lecci, F., Maria, C. and Rouvreau, V. (2017). Statistical tools for topological  
410 data analysis. arXiv: Mathematical Software: Available from [https://cran.r-](https://cran.r-project.org/web/packages/TDA/TDA.pdf)  
411 [project.org/web/packages/TDA/TDA.pdf](https://cran.r-project.org/web/packages/TDA/TDA.pdf).

412 Gameiro, M., Hiraoka, Y., Izumi, S., Kramar, M., Mischaikow, K. and Nanda, V. (2015). A  
413 topological measurement of protein compressibility. *Jpn. J. Ind. Appl. Math.* 32(1): 1–17.

414 Ghrist, R. (2008). Barcodes: The persistent topology of data. *Bulletin of the American*  
415 *Mathematical Society.* 45(1): 61–75.

416 Ghrist, R.W. (2014). *Elementary applied topology*, Seattle: Createspace.

417 Gidea, M. and Katz, Y.A. (2018). Topological Data Analysis of financial time series: landscapes  
418 of crashes. *Physica A: Statistical Mechanics and its Applications.* 491: 820-834.

419 Hamid, N.Z.A. and Noorani, M.S.M. (2014). A pilot study using chaotic approach to determine  
420 characteristics and forecasting of PM10 concentration time series. *Sains Malaysiana.* 43(3):  
421 475–481.

422 Hatcher, A. (2002). *Algebraic topology*. Cambridge University Press, Cambridge.

423 Jorquera, H., Palma, W., and Tapia, J. (2000). An intervention analysis of air quality data at  
424 Santiago, Chile. *Atmos. Environ.* 34(24): 4073-4084.

425 Khasawneh, F.A. and Munch, E. (Eds.) ASME 2014 International Mechanical Engineering  
426 Congress and Exposition 2014, American Society of Mechanical Engineers, Montreal, Canada.

427 Khasawneh, F.A., Munch, E. and Perea, J.A. (2018). Chatter Classification in Turning Using  
428 Machine Learning and Topological Data Analysis. arXiv:1804.02261 (Preprint).

429 Latif, M.T., Othman, M., Idris, N., Juneng, L., Abdullah, A.M., Hamzah, W.P., Khan, M.F.,  
430 Sulaiman, N.M.N., Jewaratnam, J., Aghamohammadi, N., Sahani, M., Chung, J.X., Ahamad,  
431 F., Amil, N., Darus, M., Varkkey, H., Tangang, F. and Jaafar, A.B. (2018). Impact of regional  
432 haze towards air quality in Malaysia: a review. *Atmos. Environ.* 177: 28-44.

433 Ling, O.H.L., Ting, K.H., Shaharuddin, A., Kadaruddin, A. and Yaakob, M.J. (2010). Urban  
434 growth and air quality in Kuala Lumpur city, Malaysia. *Environ. Asia.* 3(2): 123-128.

435 Makarau, A., Richter, R., Muller, R. and Reinartz, P. (2014). Haze detection and removal in  
436 remotely sensed multispectral imagery. *IEEE Trans. Geosci. Remote. Sens.* 52(9): 5895–5905.

437 Maletić, S., Zhao, Y. and Rajković, M. (2016). Persistent topological features of dynamical  
438 systems. *Chaos: An Interdisciplinary Journal of Nonlinear Science.* 26(5): 053105.

439 Mittal, K. and Gupta, S. (2017). Topological characterization and early detection of bifurcations  
440 and chaos in complex systems using persistent homology. *Chaos: An Interdisciplinary Journal*  
441 *of Nonlinear Science.* 27(5):051102.

442 Otter, N., Porter, M.A., Tillmann, U., Grindrod, P. and Harrington, H.A. (2017). A roadmap for  
443 the computation of persistent homology. *EPJ Data Sci.* 6(1): 17.

- 444 Payus, C., Abdullah, N. and Sulaiman, N. (2013). Airborne particulate matter and meteorological  
445 interactions during the haze period in Malaysia. *International Journal of Environmental*  
446 *Science and Development*. 4(4): 398-402.
- 447 Perea, J.A. and Harer, J. (2015). Sliding windows and persistence: An application of topological  
448 methods to signal analysis. *Foundations of Computational Mathematics*. 15(3), 799-838.
- 449 Pereira, C.M. and Mello, R.F.D. (2015). Persistent homology for time series and spatial data  
450 clustering. *Expert Syst Appl*. 42(15-16): 6026–6038.
- 451 Pigott, T.D. (2001). A review of methods for missing data. *Educ. Res. Eval*. 7(4): 353–383.
- 452 Rahman, S.R., Ismail, S.N., Ramli, M.F., Latif, M.T., Abidin, E.Z. and Praveena, S.M. (2015).  
453 The assessment of ambient air pollution trend in Klang Valley, Malaysia. *World Environment*.  
454 5(1): 1-11.
- 455 Schwartz, J., Dockery, D.W. and Neas, L.M. (1996). Is daily mortality associated specifically  
456 with fine particles? *J. Air. Waste. Manag. Assoc*. 46(10): 927–939.
- 457 Sharma, R. and Balasubramanian, R. (2018). Size-fractionated Particulate Matter in Indoor and  
458 Outdoor Environments during the 2015 Haze in Singapore: Potential Human Health Risk  
459 Assessment. *Aerosol Air Qual. Res*. 18(4): 904-917
- 460 Sivakumar, B. (2002). A phase-space reconstruction approach to prediction of suspended  
461 sediment concentration in rivers. *J. Hydrol*. 258(1-4): 149-162.
- 462 Sivakumar, B. (2003). Forecasting monthly streamflow dynamics in the western United States: a



463 nonlinear dynamical approach. *Environ. Modell. Softw.* 18(8-9):721–728.

464 Takens, F. (1981). Detecting strange attractors in turbulence. Lecture Notes in Mathematics  
465 Dynamical Systems and Turbulence, Warwick 1980: 366–381.

466 Umeda, Y. (2017). Time series classification via Topological Data Analysis. *Trans. Jpn. Soc.*  
467 *Artif. Intell.* 32(3): D-G72\_1-12.

468 Wen, Y.S., Mohd Nor, A.F., Fazilan, N.N. and Sulaiman, Z. (2016). Transboundary air pollution  
469 in Malaysia: Impact and Perspective on Haze. *Nova Journal of Engineering and Applied*  
470 *Sciences.* 5(1): 1–11.

471 Yu, S., Li, P., Wang, L., Wu, Y., Wang, S., Liu, K., Zhu, T., Zhang, Y., Hu, M., Zeng, L., Zhang,  
472 X., Cao, J., Alapaty, K., Wong, D.C., Pleim, J., Mathur, R., Rosenfeld, D. and Seinfeld J.H.  
473 (2018). Mitigation of severe urban haze pollution by a precision air pollution control approach.  
474 *Sci. rep.* 8(1): 8151.

475 Yusof, N.F.F.M., Ramli, N.A., Yahaya, A.S., Sansuddin, N., Ghazali, N.A. and Madhoun, W.A.  
476 (2010). Monsoonal differences and probability distribution of PM10 concentration. *Environ.*  
477 *Monit. Assess.* 163(1-4): 655–667.

478 Zeb, B., Alam, K., Sorooshian, A., Blaschke, T., Ahmad, I. and Shahid, I. (2018). On the  
479 morphology and composition of particulate matter in an urban environment. *Aerosol Air Qual.*  
480 *Res.* 18: 1431-1447.

481 Zhao, X., Qi, Q. and Li, R. (2010). The establishment and application of fuzzy comprehensive  
482 model with weight based on entropy technology for air quality assessment. *Procedia Eng.* 7:

483 217–222.

484

485

486

487

488

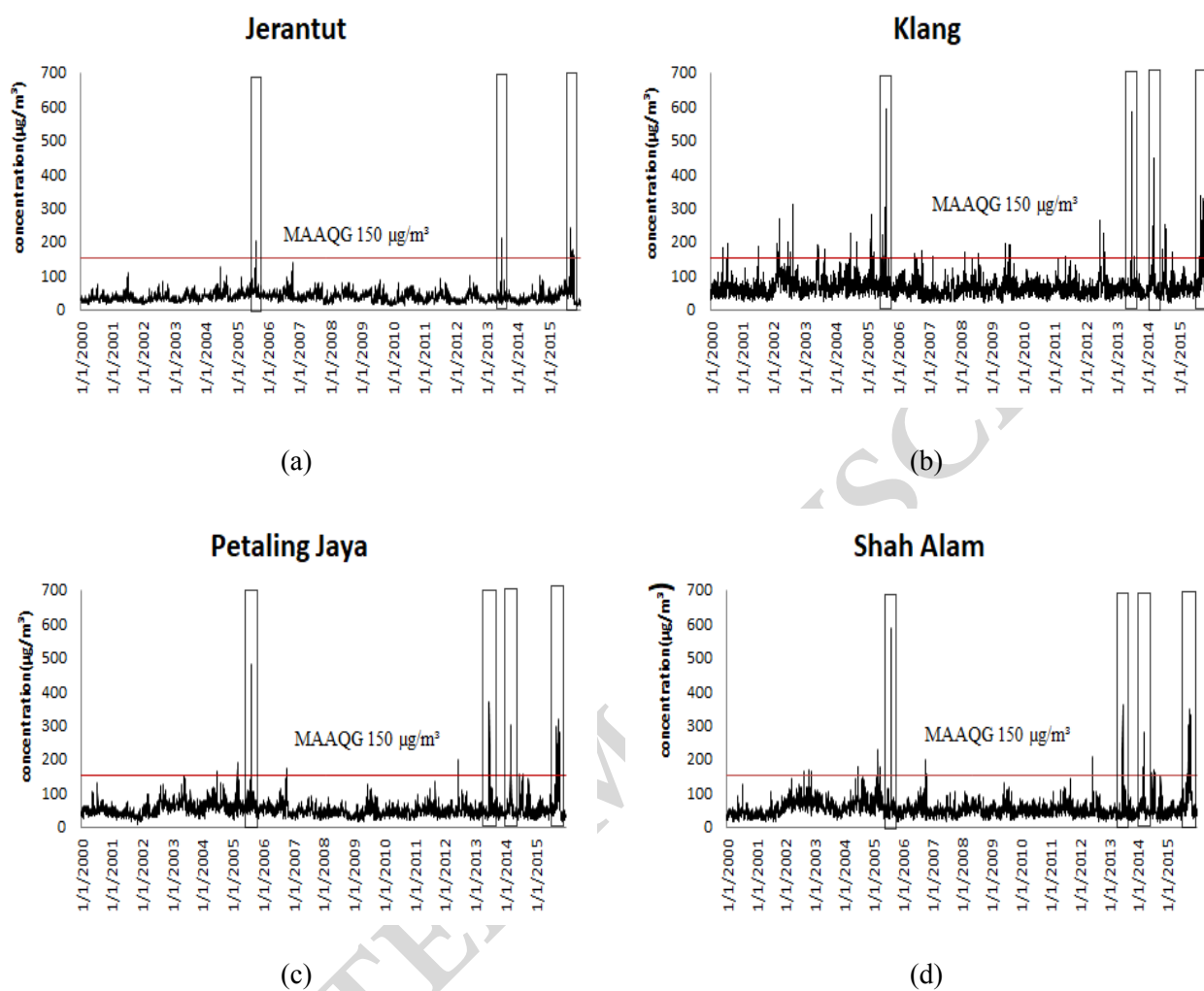
ACCEPTED MANUSCRIPT

489  
490

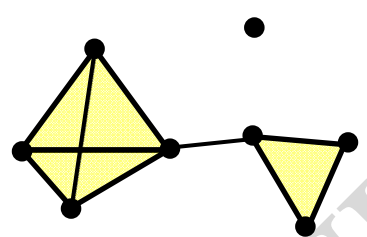
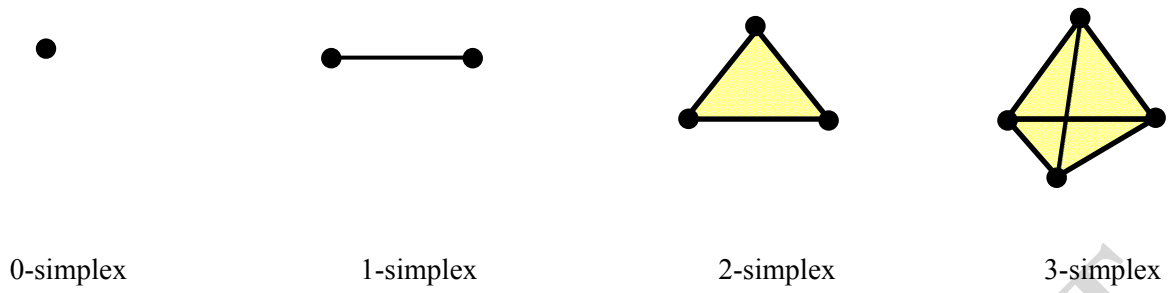
**Table 1.** Descriptive statistics for daily average of PM<sub>10</sub> for the chosen months, August 2005, June 2013, March 2014, September and October 2015.

Month	Statistic	Station			
		Jerantut	Klang	Petaling Jaya	Shah Alam
Aug-05	Min	35	36	43	26
	Max	205	590	482	587
	Mean	76.12903	139.5161	119.2581	114.8064516
	Std. Deviation	47.34888	138.0956	111.0225	135.0062762
Jun-13	Min	19	36	20	21
	Max	211	581	370	362
	Mean	56.56667	122	84.26667	83.23333333
	Std. Deviation	43.06826	125.8639	82.18731	81.23940072
Mar-14	Min	17	47	33	36
	Max	49	448	303	279
	Mean	28.67742	137.9355	94.64516	94.67741935
	Std. Deviation	7.943077	98.63364	65.17134	62.35457059
Sep-15	Min	35	59	49	49
	Max	242	337	295	301
	Mean	101.8333	141.4333	123.2	135.4
	Std. Deviation	52.76629	69.77584	64.72403	66.17458623
Oct-15	Min	13	52	24	42
	Max	176	326	320	346
	Mean	75.64516	158.6774	125.5484	147.4516129
	Std. Deviation	49.42034	76.16972	72.68509	78.30829616

491



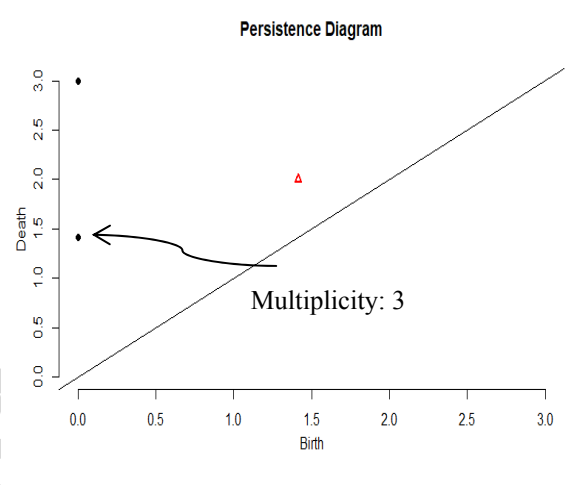
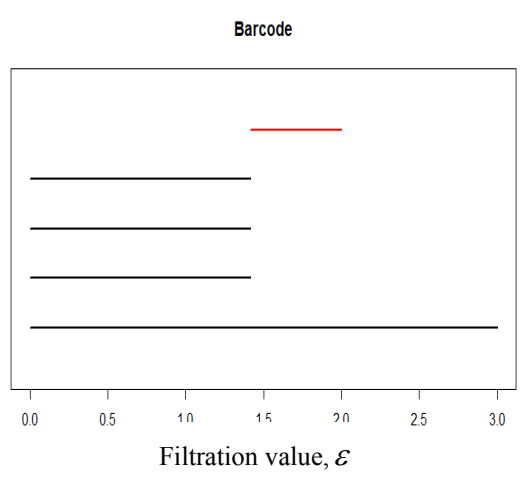
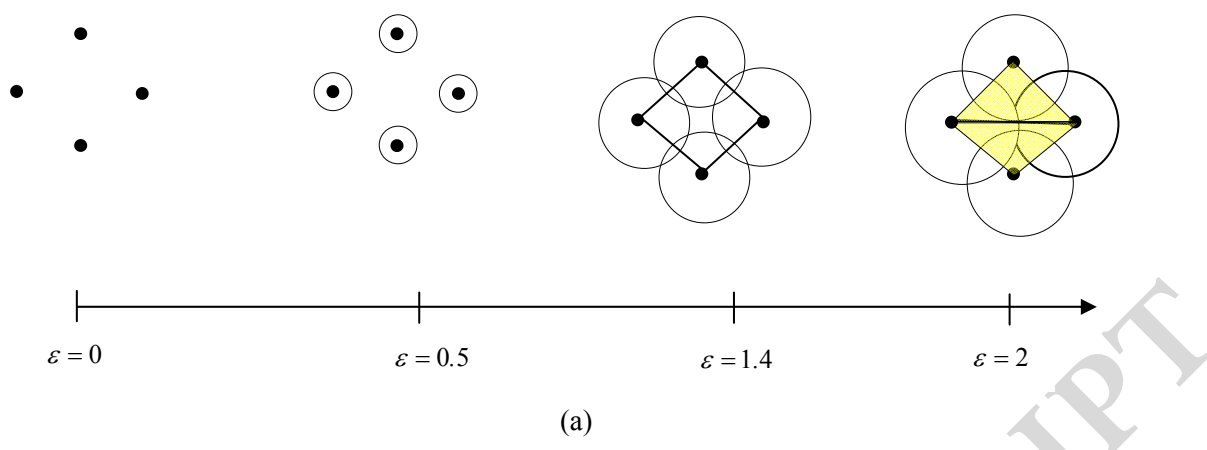
**Fig. 1.** Time series of daily average of  $PM_{10}$  for (a) Jerantut, (b) Klang, (c) Petaling Jaya and (d) Shah Alam air quality monitoring stations from 1 January 2000 until 31 December 2015.



A simplicial complex

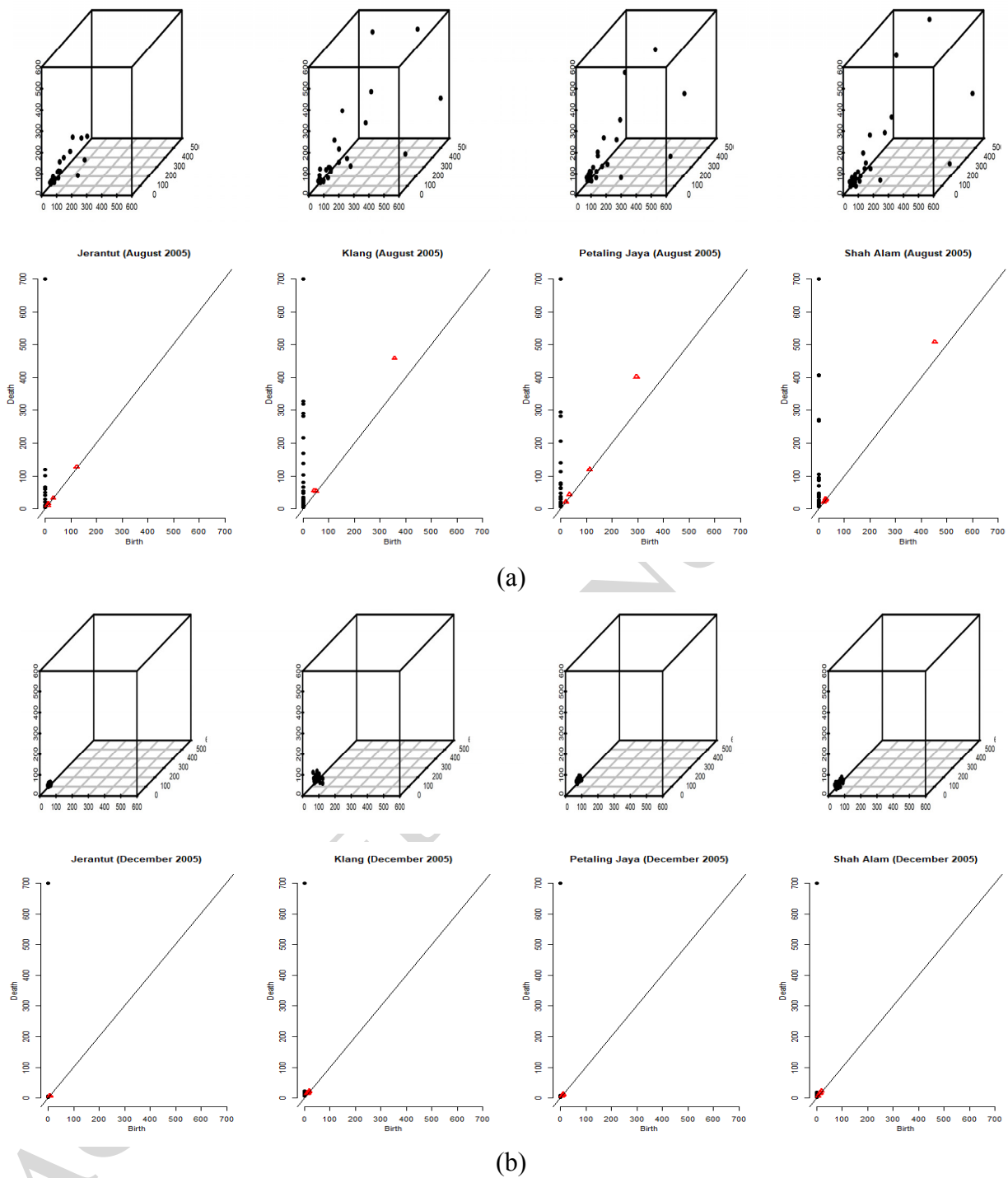
**Fig. 2.**  $k$ -simplices for  $0 \leq k \leq 3$ .

497  
498  
499  
500  
501  
502  
503  
504



**Fig. 3.** (a) Formation of simplicial complexes with respect to the filtration values,  $\varepsilon$ . (b)(c) Barcode and persistence diagram for the formation of simplicial complex illustrated in (a). In (b) and (c), the black lines and black dots represent connected components in (a), while red line and red triangle correspond to hole in (a).

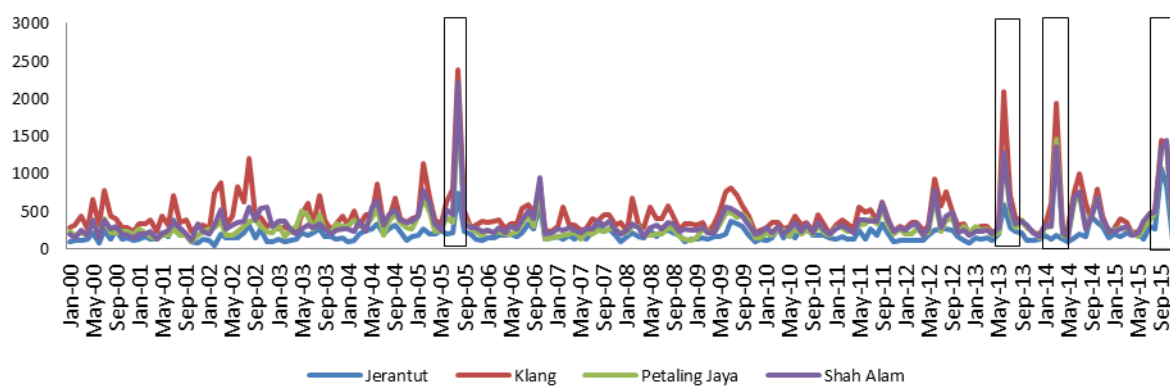
505  
506  
507  
508  
509



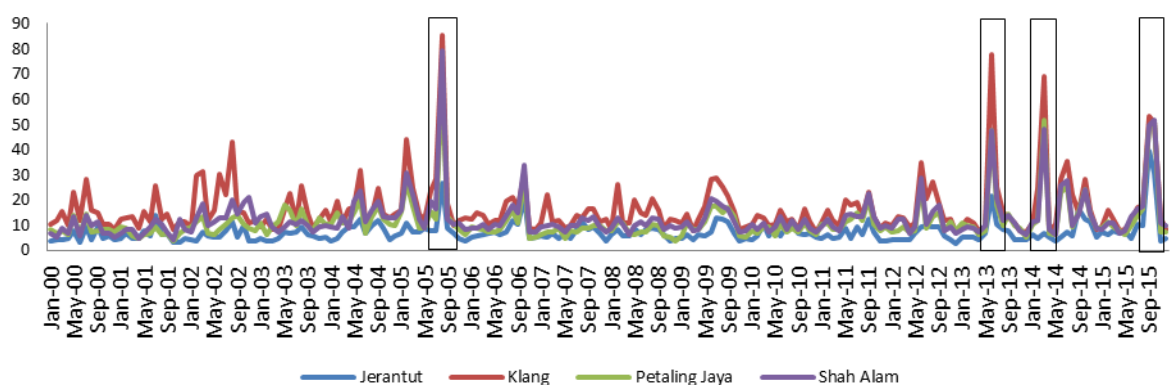
**Fig. 4.** Point clouds and persistence diagrams for the months with haze, August 2005 (a) and without haze, December 2005 (b) for Jerantut, Klang, Petaling Jaya and Shah Alam stations (left to right). The black dots and red triangles in persistence diagrams represent connected components and holes.

510

511



(a)



(b)

**Fig. 5.** (a)  $\text{sum}_0$  of all lifetimes and (b)  $\text{avg}_0$  of all lifetimes, for four air quality monitoring stations Jerantut, Klang, Petaling Jaya and Shah Alam. The rectangles indicate the months with severe haze episodes namely August 2005, June 2013, March 2014 and September and October 2015 (left to right).

512

513

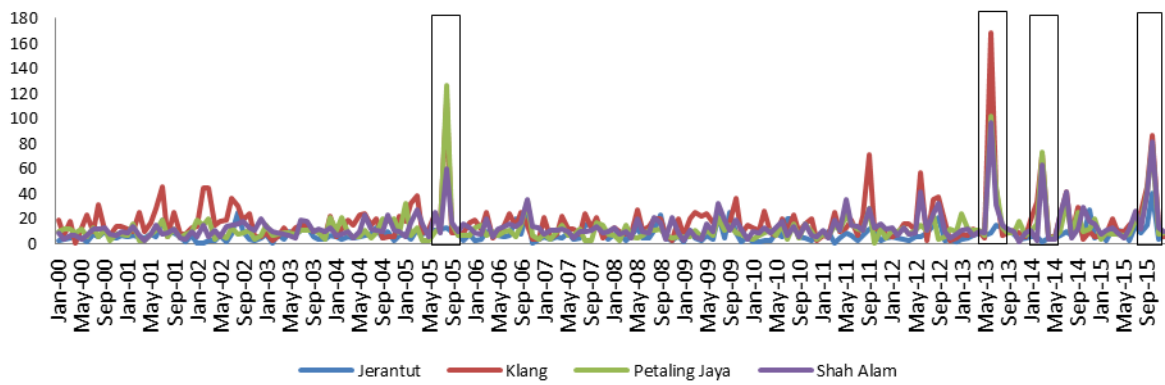
514

515

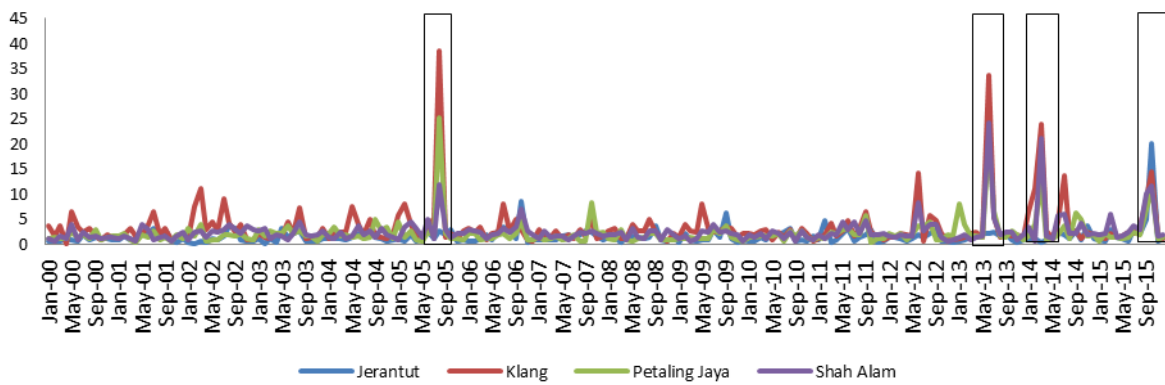
516

517

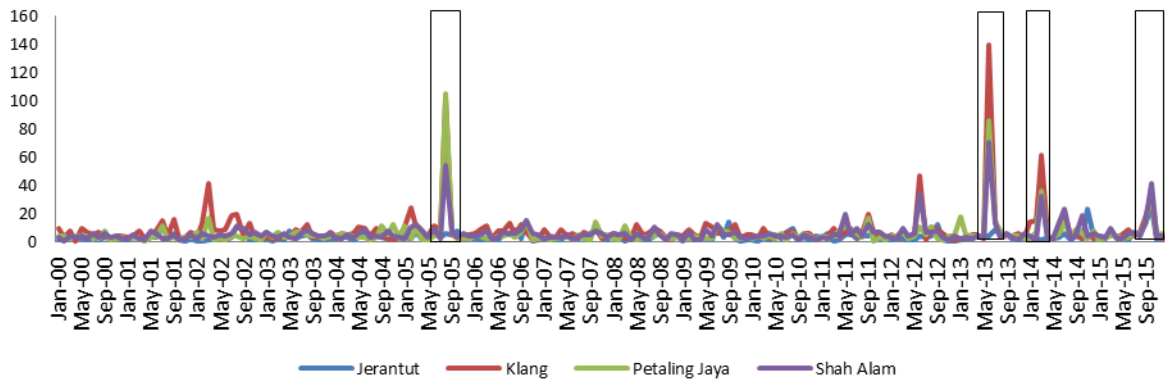




(a)



(b)



(c)

**Fig. 6.** (a)  $\text{sum}_1$  of all lifetimes (b)  $\text{avg}_1$  of all lifetimes and (c)  $\text{max}_1$  of all lifetimes for four air quality monitoring stations Jerantut, Klang, Petaling Jaya and Shah Alam.