

---

# Clustering (ab-)normal pressure signals of pipelines

---

CUI Yiran

Student ID. 20789781

ycuiat@connect.ust.hk

Room 3588

Department of Civil and Environmental Engineering, HKUST

## Abstract

1      Leakages of urban water supply systems cost considerable economic damage and  
2      labor costs per year and many methods are developed to dress problems. Transient-  
3      based methods are based on human-controllable transient waves inside the pipeline  
4      to detect the location of leaks and they have many advantages such as fast, low-  
5      cost, and accurate. However, to implement these methods preliminary knowledge  
6      of pipelines is necessary, including the number of leaks, and the existence of  
7      leaks. in order to facilitate preliminary knowledge gathering of transient-based  
8      methods, the author tried to cluster data from background noises of intact pipes,  
9      background noises of leaking pipes, transient waves generated on intact pipes,  
10     and that on leaking pipes. To explore the effects of datasets, four datasets are  
11     created and they tailored measurements, response functions of the system from  
12     measurements, sorted measurements, and Truncated Fast Fourier Transform (FFT)  
13     of measurements. The dataset of truncated FFT is more efficient than others in  
14     clustering since high-frequency noises are truncated. PCA, LTSA, Laplacian LLE,  
15     and MNcut are applied here and it turns out that four clusters can be found correctly  
16     by properly choosing the dataset and method. Noise signals are significantly  
17     different from transient signals in all four datasets, but two different noises are  
18     easier to cluster by the truncated FFT dataset and they are more entangled in the  
19     dataset of response functions.

20    

## 1 Introduction

21    With the rapid growth of urban scale and populations, water demand in urban areas is expanding.  
22    Many new supply pipes are built or connected to existing distribution networks. However, due to  
23    improper design, aging, unexpected demands and relatives, pipes may keep leaking undiscovered.  
24    The leak rate in many developed cities is around 20% to 30%, and it will cost considerable economic  
25    costs. Since most of the distribution pipes are buried in the soil, labor costs of replacement and repair  
26    are high, especially if the exact location of leakage is unknown. To locate leaks, many methods are  
27    developed, such as Smart Ball, accelerometers, transient-based methods, and so on. The transient-  
28    based methods (also known as water hammer) use the property that leaking points will reflect transient  
29    waves that can be generated by a valve closure or mechanical devices. By monitoring the reflected  
30    waves and the wave speed of the pipe, leakages can be located properly. However, due to noises inside  
31    pipes, the reflected signals are easily buried by noises. To find the reflection, it is usually implemented  
32    the cross-correlations to find the most mismatch of measured (containing leakage) signal with the  
33    intact signal. Unfortunately, in the case of unsure the existence of leakages, the cross-correlation  
34    technique can still provide a meaningless point. Therefore, it is important to study the existence of

35 leakages to gather more preliminary knowledge about the system before implementing leak-locating  
36 techniques.

37 In this article, the author tried to explore the dataset to make clusters to separate signals from cases  
38 with and without leaks, so the main question this report trying to answer is whether that is a way  
39 to cluster signals with and without leaks. It is starting from classical principle component analysis  
40 (PCA), to manifold methods such as local tangent space alignment (LTSA), Laplacian locally linear  
41 embedding (Laplacian LLE), and then to NCut methods. It is found that different datasets will  
42 influence the performances of the above-mentioned methods significantly. Besides, by choosing a  
43 dataset carefully with the proper embedding method, good clusters can be found.

44 The second section introduces datasets used in this article with an experimental description of  
45 collecting them. Besides, a brief introduction of the methods applied is given. Then section three  
46 will show the results with discussions. The fourth section provides conclusions and further questions.  
47 Besides, all data and codes used in this report can be found on the link at the end of this report.

## 48 2 Data description and methodology

49 The data used is sampled from Water-Lab in HKUST [1] and then they are modified and arranged  
50 into four datasets. Four datasets are named as 'Tailored dataset  $\mathbb{M}_T$ ', 'Response function dataset  
51  $\mathbb{M}_R$ ', 'Sorted dataset  $\mathbb{M}_S$ ', and 'Truncated FFT dataset  $\mathbb{M}_{FFT}$ '. Expect the 'Truncated FFT dataset  
52  $\mathbb{M}_{FFT}$ ', all other three are of dimensions of  $\mathbb{R}^{1001 \times 80}$ , but the 'Truncated FFT dataset  $\mathbb{M}_{FFT}$ ' is of  
53  $\mathbb{R}^{300 \times 80}$ .

### 54 2.1 Data description

55 All of the four datasets consist of four types of data with different modifications. They are background  
56 noises inside an intact pipe, background noises inside a leaking pipe, transient wave signals of an  
57 intake pipe, and that of a leaking pipe. Transient signals are generated by a manual valve closure. All  
58 four datasets have 80 columns. The first 20 columns are time series recordings of noises of an intact  
59 pipe, the second 20 columns are noises of a leaking pipe, the third 20 columns are transient signals of  
60 intact pipes, and the last 20 columns are transient signals of leaking pipes. The sampling frequency is  
61 1000 Hz.

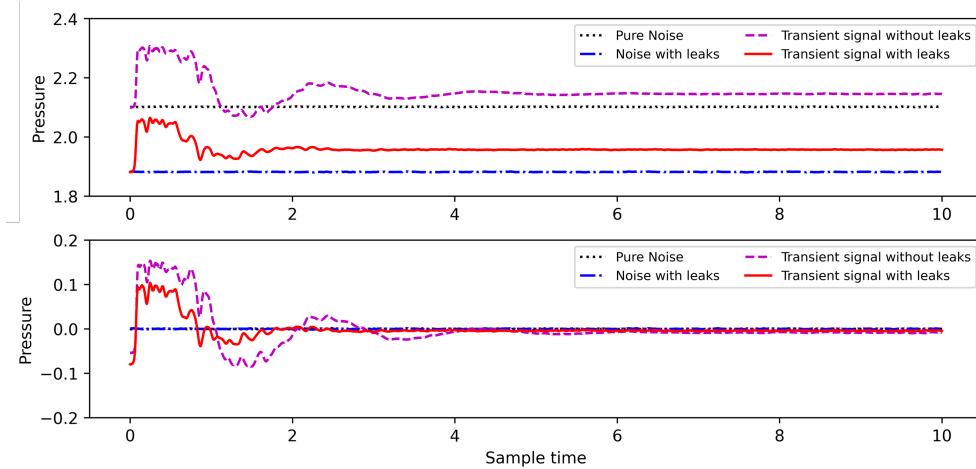


Figure 1: Example of four types of data, upper: un-normalized tailored samplings, lower: normalized  
tailored samplings. Sample time in unit of  $10^{-1}s$

62 **Tailored dataset  $\mathbb{M}_T$ :** Since the valve closure is done by a friend of the author, the reaching time of  
63 the transient signal (the time from the starting recording to the transient wave reaches the sensor)  
64 will be flexible. This means that the same transient can be represented by different time series just

65 by different reaching times. Therefore, the maximum pressure value is found in each time series  
 66 (transient waves will produce a pressure jump increasingly), and then tailor the time series from  
 67 the record that is 40 records before the maximum value point to the record that is 1000 behind the  
 68 starting record. Therefore each time series is tailored to a length of 1001 records and the 41<sup>th</sup> point  
 69 is the maximum value. Noise records were just randomly sampled for 1 second. Figure 1 upper one  
 70 shows an example of four different types of tailored data. However, it is noticeable that the static  
 71 state pressure (equilibrium state value) of the four cases is significantly different from each other.  
 72 Even this feature will make the cluster much easier, but these differences often do not exist in the real  
 73 system. Therefore, all records are normalized by subtracting the average value, and the lower part of  
 74 figure 1 shows examples of the final shape of four types of data making up the tailored dataset.

75 **Response function dataset  $\mathbb{M}_R$ :** The response function is a name borrowed from the match field  
 76 processing method since the same sampling modification method is applied. To obtain each modified  
 77 time series, the original time series is 'shifted' downward and then subtracted by the original time  
 78 series. To reduce the effects of different transient reaching times, the same procedures are done as  
 79 introduced in  $\mathbb{M}_T$  after the subtraction. Figure 2 gives an example of shapes of response function in  
 the dataset  $\mathbb{M}_R$ . However, 40 columns of noise are the same as in the tailored dataset.

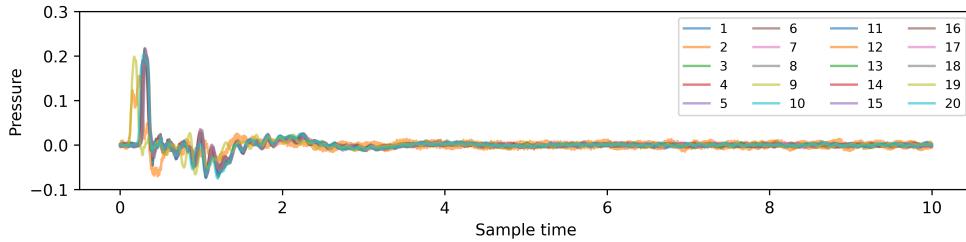


Figure 2: Example of modified transient time series in the dataset named response function

80  
 81 **Sorted dataset  $\mathbb{M}_S$ :** The sorted dataset  $\mathbb{M}_S$  is generated from Tailored dataset  $\mathbb{M}_T$ . Since the  
 82 maximum value of a transient wave will be influenced by noise and manual valve closure is also a  
 83 part of noise influencing the shape of signals, the tailored data still contain much noise. However, it is  
 84 sure that the maximum value is reached at the transient wave reaching, the minimum value is reached  
 85 at the first end-reflected transient wave reaching, and pressure goes to the equilibrium state with  
 86 fluctuating between the maximum and minimum values. Therefore, the tailored records are sorted  
 87 from small to large to make the new sorted dataset  $\mathbb{M}_S$ . In other words, the  $\mathbb{M}_S$  and  $\mathbb{M}_T$  contain the  
 88 exactly same values but the  $\mathbb{M}_S$  is the ordered version of  $\mathbb{M}_T$ .

89 **Truncated FFT dataset  $\mathbb{M}_{FFT}$ :** The Truncated FFT dataset is generated by performing Fast Fourier  
 90 Transformation onto the original sampled data (Sampling frequency: 1000 Hz, Sampling time: 10 s)  
 91 and make truncation to keep only the 1-300 Hz amplitudes, the 0-frequency amplitude representing  
 92 the constant value is dropped.

## 93 2.2 Experiment description

94 Experiments were contacted at Water-Lab HKUST. Figure 3 shows a brief sketch of setups of pipes.  
 95 A transient generation valve is near the buffer tank and can be used to generate transient waves  
 96 manually. A controllable leak labeled as LKV and at 98.42 m away from the transient generation  
 97 valve. In this report, data are collected at MV4. Transient waves are generated by a friend of the  
 98 author manually. In cases of a leak that exists, the leak point is fully open otherwise that point is  
 99 fully closed. Once starting recording, the transient is generated manually by a valve closure. More  
 100 details about the lab can be found on the Smart-UWSS website [1] and more details of experimental  
 101 procedures can be found in [2].

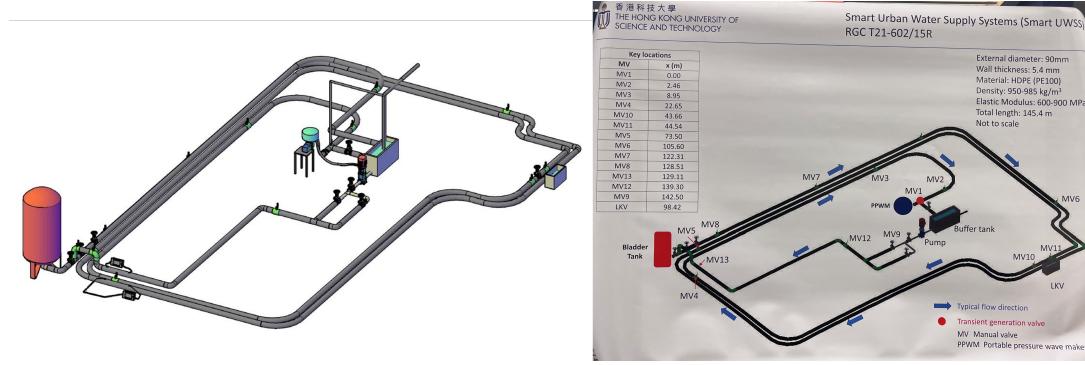


Figure 3: Pipes setups of Water-Lab in HKUST

### 102 2.3 Study questions

103 This report is trying to answer one question: Is it a way to separate/cluster signals from a dataset  
104 containing signals with and without leaks?

### 105 2.4 Methodology of analysis

106 PCA, LTSA, Laplacian-LLE, the second eigenvalue analysis (Fiedler's theorem) and eigenvector  
107 analysis (Cheeger's Inequality) are applied in this report. Here just brief introductions of them are  
108 shown, more details can be found in [3]. PCA is a classical but powerful tool in data dimension  
109 reduction.

110 Principle component analysis (PCA) looks for the variation direction of a dataset. In general, the first  
111 component represents the direction in which series of data have the most variation and the second  
112 component represents the direction that has the second largest variation. Besides, PCA is also a tool  
113 that can be used to find the local basis of high dimension data, since PCA captures the variation  
114 information of data.

115 Local tangent space alignment (LTSA) is a method of data dimension reduction by manifold learning.  
116 It uses local PCA to find the local approximate tangent space and then alignment high dimension  
117 data into lower dimension while keeping as much tangent space information of each point as possible.  
118 Since it is a manifold method, different clusters of data should be far away from each other in order  
119 to make them have different local tangent spaces to separate them. In mathematical words, the local  
120 PCA of different clusters is expected to be different from each other.

121 Laplacian locally linear embedding (Laplacian LLE) is a method that defines a normalized or un-  
122 normalized graph Laplacian to find the (approximately) optimum cut among the dataset. It can be  
123 seen that the two types of graph Laplacian are related and the normalized one is also related to row  
124 Markov matrix which is related to the random walk on a graph. Through the Fiedler's theorem the  
125 connectivity of a graph can be evaluated. It gives us a way to find the number of components in a  
126 graph (notice that four different types of data are expected as four different components). Besides,  
127 we can make an approximately optimal graph cut by Cheeger's Inequality, so it gives us a way to  
128 separate a graph even when they are fully connected.

## 129 3 Results and discussions

130 In this section, the results of the analysis are shown with discussions, starting from PCA to LTSA  
131 to Laplacian LLE. Analysis methods are applied to the four aforementioned datasets. The main  
132 conclusion that can be drawn is that there is a way to cluster different types of data but the difficulty  
133 is highly influenced by the dataset applied. Truncated FFT dataset truncated FFT dataset  $\mathbb{M}_{FFT}$  has  
134 the best performance but the response function dataset  $\mathbb{M}_R$  performs the worst.

135 **3.1 PCA**

136 Figure 4 shows the PCA analysis of the four datasets. It is clear that noise signals can be separated  
 137 from transient signals in all datasets. However, the two types of noise signals are more similar to each  
 138 other than the two types of transient signals. It is because the reflected transient signals are stronger  
 139 than noises induced by leaks in magnitudes. However, the performance of the dataset  $\mathbb{M}_R$  is worst  
 140 than others, but it is expected because the subtraction will reduce both noises and needed information.

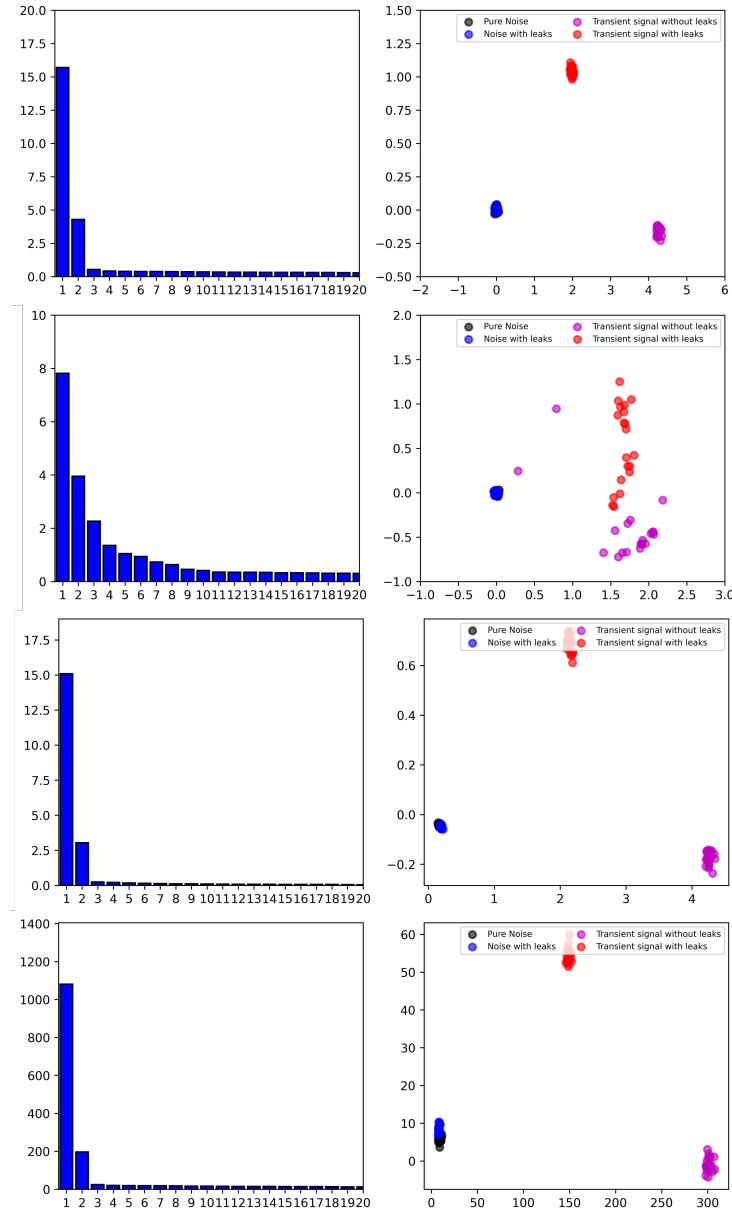


Figure 4: Principle eigenvalues of each dataset with the embedding by the first two eigenvectors. The first one on  $\mathbb{M}_T$ , the second one on  $\mathbb{M}_R$ , the third one on  $\mathbb{M}_S$  and the last one on  $\mathbb{M}_{FFT}$

141

142 Besides, it can be noticed that in the  $\mathbb{M}_{FFT}$  dataset, two types of noise are separated. Therefore,  
 143 noise data are extracted from the  $\mathbb{M}_{FFT}$  dataset to perform PCA on them, as shown in figure 5. Noise  
 144 data in other datasets are also extracted and analyzed by PCA, but two types of noise are entangled

145 with each other, so their results are not omitted here. Since the  $\mathbb{M}_{FFT}$  is a truncated dataset where  
 146 high-frequency information is deleted, this implies that leak noise mainly exists in the low-frequency  
 147 range (in this case the frequency below 300 Hz), and it is also reasonable because high-frequency  
 waves die out faster than low-frequency waves.

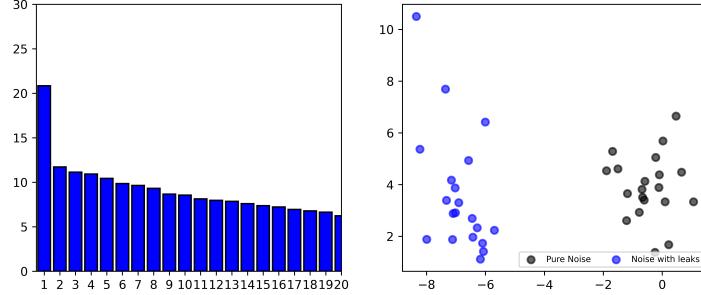


Figure 5: Principle eigenvalues of 40 noise data series from  $\mathbb{M}_{FFT}$

148

### 149 3.2 LTSA

150 In general, PCA results imply that two types of transient data are distinguished from each other  
 151 and also distinguished from noise data. However, noise data are still entangled. Therefore, LTSA  
 152 is performed mainly to explore the ability of the manifold method to separate noise data. Actually,  
 LTSA is also performed on not only noise data but on the whole dataset due to curiosities.

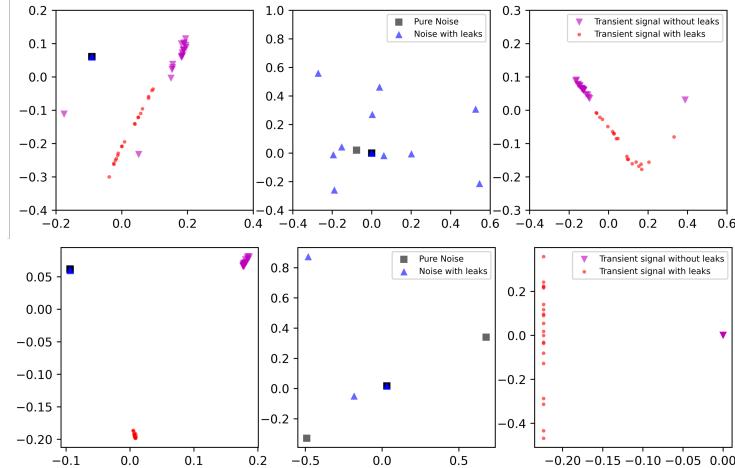


Figure 6: LTSA embedding of the whole dataset, separated noise data and separated transient data, the upper one is based on  $\mathbb{M}_R$  and the lower one is based on  $\mathbb{M}_{FFT}$

153

154 Only LTSA results based on  $\mathbb{M}_R$  (noise data is actually the same as  $\mathbb{M}_T$ ) and on  $\mathbb{M}_{FFT}$  are shown  
 155 here in figure 6, the upper one and the lower one, respectively. It can be found that LTSA can separate  
 156 noise signals from transient signals as well, but two types of noise signals are still entangled. This is  
 157 because the variation inside the transient data group is higher than that inside the noise data group,  
 158 hence when performing the global alignment small variations of noises data are dropped, as we only  
 159 use the smallest 2 eigenvectors in this case. Besides, the results of LTSA on noise data separated  
 160 from  $\mathbb{M}_R$  show that pure noise (the case without leak) is more concentrated than of with leak, but  
 161 in the case of  $\mathbb{M}_{FFT}$  both two types of noises are concentrated into one point expect four outliers.  
 162 This means that with the leak existing background noise has more variations than that without leak,  
 163 and those variations mainly exist in high frequency that is truncated in the  $\mathbb{M}_{FFT}$ . The last thing

164 that should be mentioned is that the embedding of transient signals without leaks from  $\mathbb{M}_{FFT}$  is also  
 165 shrinking into a point. It means that variations inside transient signals with leak are much higher than  
 166 those inside transient signals without leak.

### 167 3.3 Laplacian LLE

168 The Laplacian LLE is performed and the second eigenvalues of normalized graph Laplacian and  
 169 eigenvectors are also calculated to further explore the Laplacian LLE results. Figure 7 shows the  
 170 results of Laplacian embedding results of a simple connected graph defined by sampled data with  
 171 different settings of the number of neighbors. It can be found that two types of transient data can be  
 172 well separated, except in the case of the graph generated by response function (upper-right figure)  
 173 one data series of transient signals without leaks is embedded closer to data series of transient signals  
 174 with leaks. This may be because of an improper subtraction of the original data. Except for noise  
 175 data from  $\mathbb{M}_{FFT}$ , all other embeddings of noise data show that two types of noise are entangled. The  
 176 successful embedding of noise data from  $\mathbb{M}_{FFT}$  implies that low-frequency parts of two types of  
 177 noise are distinguished but high-frequency parts are similar. Besides, results of embeddings of the  
 178 whole dataset show that two types of noises are more similar to each other than two types of transient  
 signals.

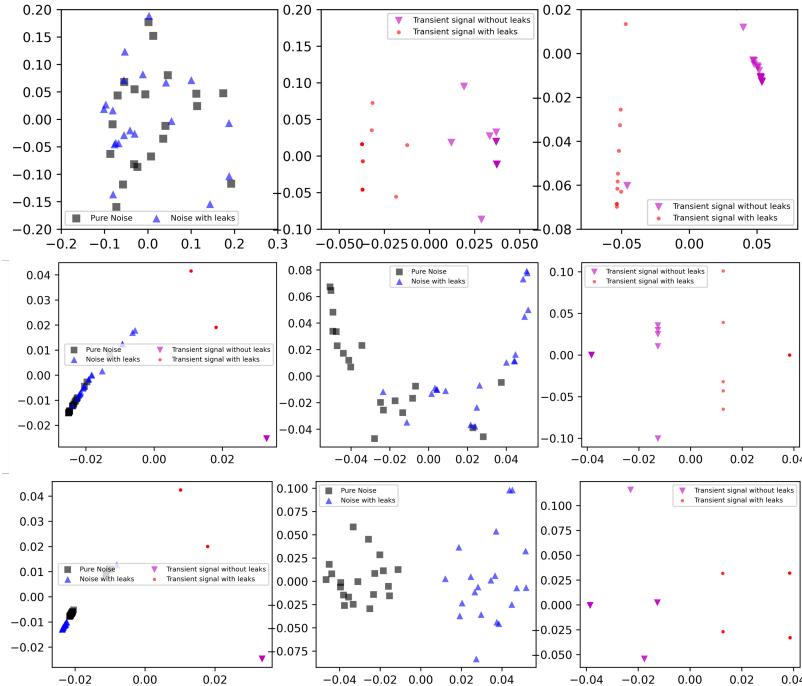


Figure 7: Laplacian LLE results of simple connected graph generated from different datasets. The first one is based on noise data from  $\mathbb{M}_T$  with 5 neighbors (left), transient data from  $\mathbb{M}_T$  with 22 neighbors (middle), and transient data from  $\mathbb{M}_R$  with 22 neighbors (right). The second one is based on the whole data from  $\mathbb{M}_S$  with 25 neighbors (left), noise data from  $\mathbb{M}_S$  with 25 neighbors (middle), and transient data from  $\mathbb{M}_S$  with 25 neighbors (right). The last one is based on the whole data from  $\mathbb{M}_{FFT}$  with 25 neighbors (left), noise data from  $\mathbb{M}_{FFT}$  with 25 neighbors (middle), and transient data from  $\mathbb{M}_{FFT}$  with 25 neighbors (right)

179

180 Since the above results are generated from simple connected graphs, they can be further improved by  
 181 adding weights to generate weighted graphs. By assigning different denominators to the power of  
 182 exponential, neighbors who are far away from each other can be further distinguished. Parts of the  
 183 results of weighted Laplacian graphs LLE are shown in figure 8. It can be seen that in this case the  
 184 noise data from  $\mathbb{M}_S$  are separated better than in the case of simple connected graph. However, noise

185 data from  $\mathbb{M}_T$  are still entangled. Adding weights to neighbors can make non-similar neighbors has  
 186 less weight but if different data are mixed up with each other (a point with both the same type and  
 187 different types of neighbors with similar distance), this method will fail. By rearranging the number  
 of neighbors and carefully setting  $\gamma$ , there do has a hope to cluster them properly.

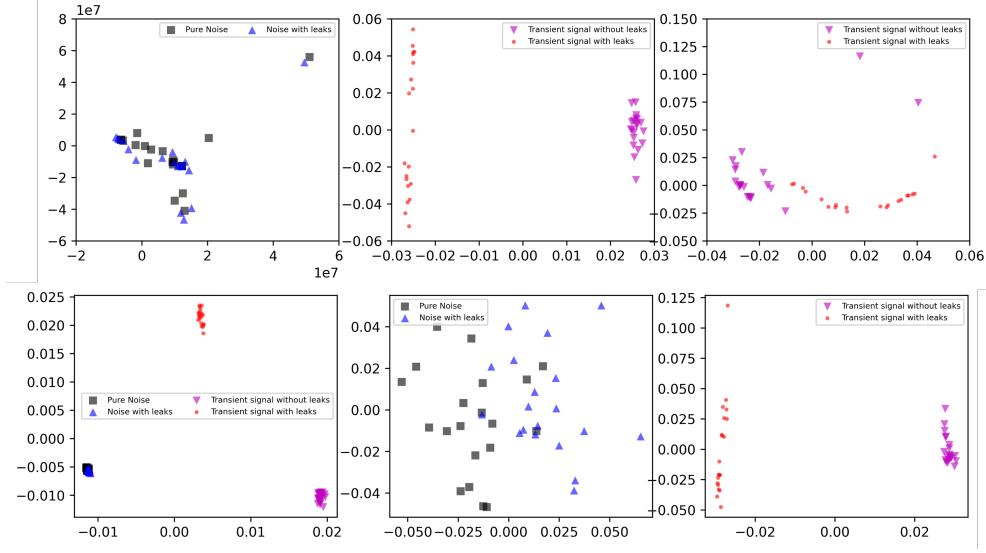


Figure 8: Laplacian LLE results of weighted connected graph generated from different datasets. The upper left: noise data from  $\mathbb{M}_T$  with  $\gamma = 520$ , the upper middle: transient data from  $\mathbb{M}_T$  with  $\gamma = 0.01$ , the upper right: transient data from  $\mathbb{M}_R$  with  $\gamma = 0.01$ , the lower left: the whole data from  $\mathbb{M}_S$  with  $\gamma = 10^{-10}$ , the lower middle: noise data from  $\mathbb{M}_S$  with  $\gamma = 10^{-5}$ , the lower right: transient data from  $\mathbb{M}_S$  with  $\gamma = 0.1$ , where  $\gamma$  is a parameter of calculating weights  $e^{-\gamma \|x_1 - x_2\|^2}$

188  
 189 To further explore the ability of graph Laplacian, eigenvalues of unnormalized graph Laplacian  
 190 (simple connected) are calculated. Based on Fiedler's theorem, the number of zeros is equal to the  
 191 number of components on a graph. Therefore, it is clear that in both  $\mathbb{M}_T$ ,  $\mathbb{M}_S$  and  $\mathbb{M}_{FFT}$ , two types  
 192 of noise are connected. In four datasets transient signals generated unnormalized graph Laplacians  
 193 imply that there are at least two components on the graph. However, as shown in figure 7 and 8, these  
 194 components may consists of data with different types.

195 To say components of graphs, the normalized graph Laplacians are calculated and the first three  
 196 smallest eigenvectors (the eigenvector corresponding to eigenvalue zero is dropped hence it is  
 197 excluded here) are plotted. The results based on dataset  $\mathbb{M}_S$  and  $\mathbb{M}_{FFT}$  are shown as an example  
 198 in figure 10. It is clear that if play with the whole dataset, two types of transient data are seen as  
 199 two components which is what we want but two types of noise data are mixed up as one component.  
 200 If the normalized graph Laplacian is built on transient types of data, two components with correct  
 201 clusters can still be found. In the case of noise in  $\mathbb{M}_S$ , two components are found in the graph but it  
 202 can be seen that each component consists no longer only one type of noise. Since all settings of them  
 203 are the same but the noise in  $\mathbb{M}_{FFT}$  can be clustered nicely, it suggests that the high-frequency part  
 204 of background noise will negatively influence the performance of Laplacian LLE.

205 To further illustrate the point, figure 11 shows the approximate optimal NCut calculated from noise  
 206 data from  $\mathbb{M}_S$  and  $\mathbb{M}_{FFT}$ . Clearly, noise from  $\mathbb{M}_{FFT}$  can give a nice cut where each component  
 207 contains only one type of data but that from  $\mathbb{M}_S$  generates a relatively poor cut even though it still  
 208 cut data into two clusters. To illustrate the meaning of Neut of noise data from  $\mathbb{M}_{FFT}$ , the Cheeger's  
 209 constant calculated from the reordered node set ( $f_1 \leq f_2 \leq \dots \leq f_n$ ,  $\vec{f} = D^{-1/2}v$ , where  $D$  is the  
 210 degree matrix and  $v$  is the eigenvector corresponding to the second smallest eigenvalue) is shown in  
 211 figure 12. A clear sharp corner is found in that figure and this implies that the NCut does separate  
 212 two different types of data otherwise the curve of Cheeger constant will be flattened.

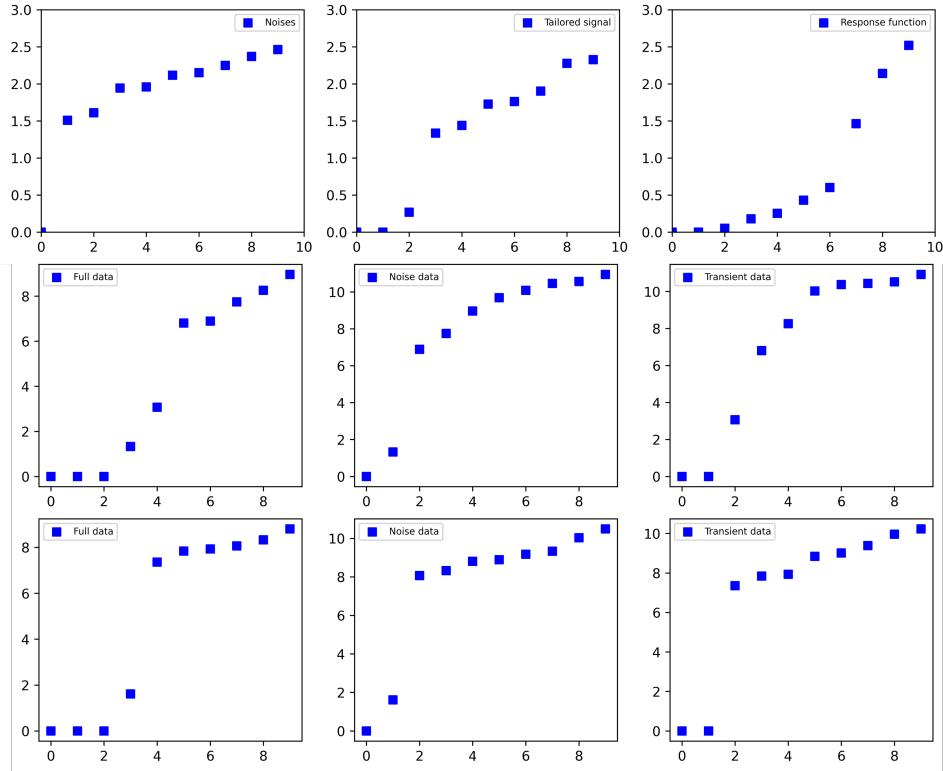


Figure 9: The ascending ordered eigenvalues of unnormalized Laplacian graph (to apply Fiedler theory). They are generated based on search neighbors with the number of 10 and calculated based on: noise data from  $\mathbb{M}_T$  (the upper left), the transient data from  $\mathbb{M}_T$  (the upper middle), the transient data from  $\mathbb{M}_R$  (the upper right), the whole dataset of  $\mathbb{M}_S$  (the middle left), the noise data from  $\mathbb{M}_S$  (the middle central), the transient data from  $\mathbb{M}_S$  (the middle right), the whole dataset of  $\mathbb{M}_{FFT}$  (the lower left), the noise data from  $\mathbb{M}_{FFT}$  (the lower middle) and the transient data from  $\mathbb{M}_{FFT}$  (the lower right)

## 213 4 Conclusions

214 This report is trying to answer the question 'Is it a way to separate/cluster signals from a dataset  
 215 containing signals with and without leaks?'. Pressure signals are sampled from Water-Lab in HKUST,  
 216 and four different datasets are generated based on them. Both PCA, LTSA, and Laplacian LLE are  
 217 done to answer the question. It turns out that there do have a way to cluster different types of signals,  
 218 but the final cluster highly depends on the quality of the dataset applied and the pre-knowledge of  
 219 the processed data. With a proper understanding of the data (at least the magnitude of different  
 220 phenomena), graph Laplacian provides a way to achieve the job.

## 221 Supplement

222 All data are collected by the author and his friends, so it is decided to publish the data for educational  
 223 purposes. All data and codes can be found here: [https://drive.google.com/drive/folders/1ygImku3AsiVB9mamjkzqyPsLHnm9X9nt?usp=share\\_link](https://drive.google.com/drive/folders/1ygImku3AsiVB9mamjkzqyPsLHnm9X9nt?usp=share_link)

## 225 References

- 226 [1] Smart UWSS. (2023, May 17). *Controlled facility 1: HKUST lab pipeline*. Smart Urban Water Supply  
 227 System. URL:[https://sites.google.com/view/smartuwss/facilities/Facilities\\_1](https://sites.google.com/view/smartuwss/facilities/Facilities_1)

- 228 [2] Wang, X., Lin, J., Keramat, A., Ghidaoui, M.S., Meniconi, S., & Brunone, B.(2019c). Matched-Field  
 229 Processing for Leak Detection in a Viscoelastic Pipe: An Experimental Study. Mechanical Systems and  
 230 Signal Processing, 124, 459-478
- 231 [3] Yao Yuan (2023). *Geometric and Topological Data Reduction: A Mathematical Introduction to Data Science*.  
 232 URL:<https://yao-lab.github.io/bookdatasci/>

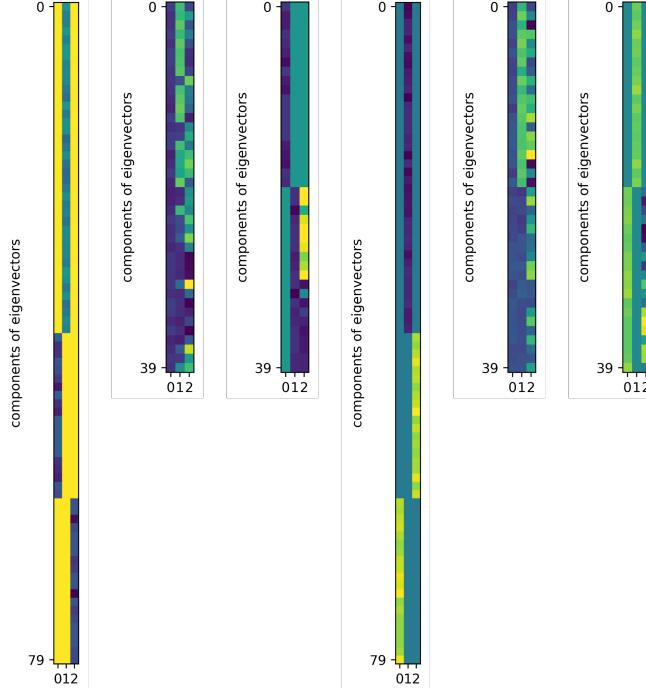


Figure 10: The first three smallest eigenvectors of normalized graph Laplacian with different colors representing values, positive number are marked as green, negative number is marked as blue and zero is marked as yellow, darker color means higher absolute value. The leftest (first) based on the whole dataset of  $\mathbb{M}_S$ , the second one from noise data from  $\mathbb{M}_S$ , the third one from transient data from  $\mathbb{M}_S$ , the fourth one from the whole dataset of  $\mathbb{M}_{FFT}$ , the fifth one from noise data from  $\mathbb{M}_{FFT}$ , and the last (sixth) one from transient data from  $\mathbb{M}_{FFT}$



Figure 11: Approximate optimal NCut based on Cheeger's constant, calculated from the noise data from  $\mathbb{M}_S$  (upper), and noise data from  $\mathbb{M}_{FFT}$  (lower), two colors represents two different clusters of data after the cut.

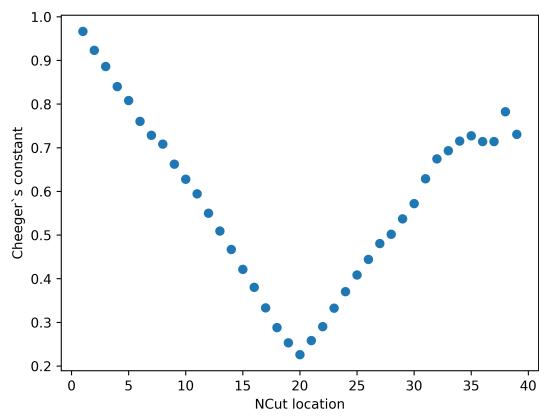


Figure 12: The Cheeger constant calculated from the reordered node set of noise data in dataset  $\mathbb{M}_{FFT}$