

A time series distance measure for efficient clustering of input/output signals by their underlying dynamics*

Oliver Lauwers¹ and Bart De Moor¹

Abstract—Starting from a dataset with input/output time series generated by multiple deterministic linear dynamical systems, this paper tackles the problem of automatically clustering these time series. We propose an extension to the so-called Martin cepstral distance, that allows to efficiently cluster these time series, and apply it to simulated electrical circuits data.

Traditionally, two ways of handling the problem are used. The first class of methods employs a distance measure on time series (e.g. Euclidean, Dynamic Time Warping) and a clustering technique (e.g. k-means, k-medoids, hierarchical clustering) to find natural groups in the dataset. It is, however, often not clear whether these distance measures effectively take into account the specific temporal correlations in these time series. The second class of methods uses the input/output data to identify a dynamic system using an identification scheme, and then applies a model norm-based distance (e.g. H_2 , H_∞) to find out which systems are similar. This, however, can be very time consuming for large amounts of long time series data.

We show that the new distance measure presented in this paper performs as good as when every input/output pair is modelled explicitly, but remains computationally much less complex. The complexity of calculating this distance between two time series of length N is $\mathcal{O}(N \log N)$.

I. INTRODUCTION

Time series clustering is an important topic in modern research. State-of-the-art clustering methods of other data types are often not suited for this high-dimensional, temporally correlated data structure. Clustering is the task of finding groups with similar elements in a dataset and consists of three components: a similarity measure based on relevant data features, a clustering algorithm and an evaluation criterion. While the latter two components might carry over, defining a good distance measure is a difficult problem, especially if one is interested in the dynamics of the generating dynamical system of the time series.

Representing the time series as single-input single-output (SISO) linear time invariant (LTI) deterministic dynamical systems further generates problems of its own, as the contributions of the input signal and the impulse response of the system are convolved in the time domain. It is thus

*This work was supported by Belgian Federal Science Policy Office: IUAP P7/19 (DYSCO, Dynamical systems, control and optimization, 2012-2017) Flemish Government: IWT: TBM IETA(130256); PhD grants Industrial Research fund (IOF): IOF Fellowship 13-0260 VLK Stichting E. van der Schueren: rectal cancer EU H2020-SC1-2016-2017 Grant Agreement No.727721: MIDAS Meaningful Integration of Data, Analytics and Services KU Leuven Internal Funds C16/15/059, C32/16/013 KIC EIT Health: New MOOC - Data Analytics in Health; EIT Health Summer School Innovation on Big Data for Healthy Living imec strategic funding 2017. Oliver Lauwers is supported by an SB-grant of the FWO (formerly IWT).

¹Oliver Lauwers and Bart De Moor. Stadius, Department of Electrical Engineering (ESAT), KU Leuven, 3000 Leuven, Belgium {oliver.lauwers, bart.demoor}@esat.kuleuven.be

not intuitively clear how these two contributions can be separated, for example when one is interested only in the dynamics of the system and not in the specific input signal.

This problem grows ever more relevant as large scale big data time series problems grow more prevalent in areas like finance, medicine, or the industrial internet of things, where clustering is important in tasks like anomaly detection [7], [12]. A typical industrial problem contains several hundred sensors per machine, tens of machines per plant, and several plants per industrial player, collecting data every few seconds, for months or even years of operation time. This results in datasets of several million time points for thousands of series. Clustering techniques should thus scale well.

In Section II we look at state-of-the-art clustering methods for time series from two perspectives, starting from a dataset containing input/output time series pairs, generated by different SISO LTI dynamical systems. From a machine learning point of view, we use an automated clustering method with an off-the-shelf time series distance such as the Euclidean distance or Dynamic Time Warping (DTW). From a system identification point of view, we apply norms such as the H_2 or H_∞ norm to compare systems estimated from the data. We find that these techniques either are very fast, but give poor results, or perform well, but are computationally expensive.

Next, in Section III, we look at the Martin cepstral distance [3], [8], which combines insights from systems theory into a distance measure that can be computed on the raw data. This metric was defined for SISO ARMA models (i.e. LTI models that use white noise as an input signal).

The main contribution of this paper is an extension of the cepstral distance measure, that incorporates deterministic input signals, and allows to calculate distances between a broader class of SISO LTI dynamical systems. It thus allows to cluster time series by dynamics, but remains computationally much simpler than explicitly estimating models.

Subsequently, we apply this new distance measure in Section IV to an application on electrical circuits, where we generate a dataset consisting of input/output signal pairs, and the problem is to identify which data belong to which generating system. Finally, we conclude the paper and provide some paths for future research in Section V.

II. EXISTING METHODS

Existing methods to cluster time series employ a clustering technique, together with some distance measure. The author of [6] discerns three types of distance measures: measures based on raw data, measures based on features of the time series and measures based on models. For the scope of this

paper, we will focus on the first and the latter (as the distance measure we propose combines elements of these two broad classes). We present two raw data distance measures, the Euclidean metric and the Dynamic Time Warping metric [5], and two model-based distance measures, connected to the H_2 -norm and the H_∞ -norm. In the next section, we will introduce and extend the cepstral distance [3], [8], which combines the efficiency of the raw data distance measures with the insight in generative dynamics of the model norms, and thus has representations both as a raw data distance and as a model-based one.

A. Raw Data Distance Measures

In what follows we will define u_m to be the input signal of the m -th element of a dataset, y_m is the corresponding output signal and $u_m(k)$ or $y_m(k)$ is the value at timepoint k of respectively the input and output of the m -th element of the input/output dataset. Time series from element m start at $k = 0$ and end at $k = N_m$. The system that generated an output from a given input will be called the generating (dynamical) system.

1) Euclidean Distance:

Definition 1. *The Euclidean distance, $d_E(\cdot, \cdot)$ treats the time series as a vector, and applies the element-wise Euclidean vector distance between two time series of same length N_m , defined as*

$$d_E(y_m, y_n) = \sqrt{\sum_{k=0}^{N_m} (y_m(k) - y_n(k))^2}. \quad (1)$$

Advantages

- The Euclidean distance is easy to calculate, allowing for very efficient computation and clustering.
- No system identification step is needed.

Disadvantages

- There is no clear link between this distance measure and the generating system.
- This measure treats the time series as a vector, and ignores the temporal correlations in the data.
- This measure does not allow to compute distances between time series of different length.
- This measure does not take the input into account.

2) Dynamic Time Warping:

Dynamic Time Warping (DTW) [5], [11] is an algorithm that tries to locally align time series, by *warping* them such that the Euclidean distance between the warped time series is minimal. Mathematically, this *warping*, and the measure that is found in this way, can be described as follows.

Given two output signals, y_1 and y_2 , of length N_1 and N_2 respectively, a matrix M is constructed, where the (l, m) -th element of M is defined as $M_{(l,m)} = (y_1(l) - y_2(m))^2$. A warping path, $W = w_1, w_2, \dots, w_k, \dots, w_K$ is then defined,

with each $w_k = (M_{(l,m)})_k$ an element of matrix M and $\max(N_1, N_2) \leq K < N_1 + N_2 - 1$.

The path is subject to the boundary conditions $w_1 = M_{1,1}$ and $w_K = M_{N_1, N_2}$ (i.e. the path starts in one corner of the matrix and ends in the opposite one), has to be continuous, in such a way that two consecutive elements w_k and w_{k+1} are maximally one column and one row apart, and has to be monotonously increasing in its indices, i.e., that in going from w_k to w_{k+1} , column nor row number can decrease.

Definition 2. *We are now interested in the warping path W_{DTW} that minimizes the cost function*

$$d_{DTW}(y_1, y_2) = \min \left\{ \sqrt{\sum_{k=1}^K w_k} \right\}. \quad (2)$$

The sum over this path is then the DTW distance between the time series.

Though this algorithm is computationally expensive due to the combinatorial nature of the problem, several lower bounds have been devised that can be implemented efficiently. In what follows, we use the Keogh Lower Bound [5] as an efficient approximation to the DTW distance.

Advantages

- The DTW distance takes into account (part of) the local temporal correlations.
- No system identification step is needed.
- Lower bounds on the distance are reasonably efficient.
- This measure allows to calculate distances between time series of different length.

Disadvantages

- There is no clear link between this distance measure and the generating system.
- The DTW distance as such is expensive to calculate.
- This measure does not take the input into account.

B. Model-based Distance Measures

We use the same notation as in subsection II-A. The generating system of the input/output pair (u_m, y_m) will be denoted by M_m , and its corresponding transfer function will be written \mathcal{H}_m . Based on a model norm $\|\cdot\|$, the distance between two models M_i and M_j is defined as $\|\mathcal{H}_i - \mathcal{H}_j\|$.

1) H_2 -norm:

Definition 3. *The H_2 -norm, $\|\mathcal{H}\|_2$, of a discrete-time system M with transfer function \mathcal{H} is defined as*

$$\|\mathcal{H}\|_2 = \sqrt{\frac{1}{2\pi} \int_{-\pi}^{\pi} \text{Tr} \{ \mathcal{H}^H(e^{i\omega}) \mathcal{H}(e^{i\omega}) \} d\omega}, \quad (3)$$

where $\text{Tr}\{\}$ denotes the trace, the superscript \cdot^H denotes the Hermitian conjugate and i denotes the imaginary unit.

The H_2 -norm can be seen as the root-mean-square of the system response to a normalized white noise input. It is thus a measure of the power, or steady-state variance of this response. The H_2 -norm will be infinite for unstable systems.

Advantages

- The H_2 -norm provides a physically interpretable way to characterize underlying dynamics of time series.
- This norm allows to calculate distances between time series of different length.
- This norm takes the input data into account.

Disadvantages

- A system identification procedure is needed, which is both difficult to automate and often computationally expensive (at least more expensive than the raw data measures).

2) H_∞ -norm:

Definition 4. The H_∞ -norm, $\|\mathcal{H}\|_\infty$, of a discrete-time system M with transfer function \mathcal{H} is calculated as

$$\|\mathcal{H}\|_\infty = \max_{\omega \in [0, \pi[} |\mathcal{H}(e^{i\omega})|. \quad (4)$$

This norm thus measures the maximal gain of the frequency response and is called the *gain* of the system. It becomes infinite for systems with poles on the unit circle.

Advantages

- The H_∞ -norm provides a physically interpretable way to characterize underlying dynamics of time series.
- This norm allows to calculate distances between time series of different length.
- This norm takes the input data into account.

Disadvantages

- A system identification procedure is needed, which is both difficult to automate and often computationally expensive (at least more expensive than the raw data measures).

III. CEPSTRAL DISTANCE

In this section we take a closer look at an insightful distance measure on ARMA models, which can be interpreted both as a raw data distance measure and as a model norm: the Martin cepstral norm [3], [8]. We first give a very concise review of the cepstral norm in the stochastic case, then proceed with an extension that allows us to incorporate information about the deterministic input signal.

A. Original Cepstral Norm

Based on the power spectral density, Φ_y , of a signal y , we can define its power cepstrum, c_y as

$$c_y = \mathcal{F}^{-1}(\log(\Phi_y)), \quad (5)$$

where \mathcal{F}^{-1} denotes the inverse Fourier transform. This produces a series of coefficients, $c_y(k)$, with integer $k \in [0, N]$, where N denotes the length of time series y .

Definition 5. The cepstral norm, $\|\mathcal{H}\|_C$, of model M with transfer function \mathcal{H} , and output y is defined as

$$\|\mathcal{H}\|_C = \sum_{k=0}^N k (c_y(k))^2. \quad (6)$$

For ARMA models it was proven in [3] that there are multiple methods to calculate this norm: it can be derived

from the subspace angles of the output Hankel matrices of the generating system, from the mutual information of the output space of a system, and from a combination of poles and zeros of the transfer function of the model. Moreover, equation (6) allows us to calculate the norm straight from raw data, without the need to identify the underlying systems. We can thus connect the cepstral norm to a raw data distance measure in the following sense:

Definition 6. The cepstral distance, $d_C(y_i, y_j)$, between two time series, y_i and y_j , is defined as

$$d_C(y_i, y_j) = \sum_{k=0}^{\max\{N_i, N_j\}} k (c_{y_i}(k) - c_{y_j}(k))^2, \quad (7)$$

where $\max\{N_i, N_j\} - \min\{N_i, N_j\}$ zeros are added at the end of the cepstrum of length $\min\{N_i, N_j\}$.

Advantages

- The cepstral distance has an interpretation in terms of the generating model of the time series.
- The cepstral distance is easy to calculate, allowing for very efficient computation and clustering.
- No system identification step is needed.
- This measure allows to calculate distances between time series of different length.

Disadvantages

- This distance measure can only take information coming from a stochastic input into account.

B. Extended Cepstral Distance

The cepstrum, defined in the previous section, finds its roots in homomorphic signal processing [9, Chapter 10]. In this type of processing, the original time series data, which often involves complex multiplicative operators like convolutions, is mapped, through a non-linear mapping, to a different domain, that allows for linear filtering. The cepstrum, as in equation (5), is a good example. The convolution in the time domain changes into a multiplication by calculating the power spectral density. Applying a logarithmic transformation then turns the multiplication in frequency domain into an addition. Finally, the inverse Fourier transform takes the problem back to (a transformed version of) the time domain. Equation (5) is thus effectively a method to transform the convolution into an addition.

This allows us to take the output, and separate the contributions from the input signal (which was the main disadvantage left in the cepstral distance, see subsection III-A) and the impulse responses of the system. Indeed, defining the cepstrum coefficients of the input signal u as $c_u(k)$, and the contribution to the cepstrum coefficients of the transfer function \mathcal{H} as $c_h(k)$, we can write

$$c_y(k) = c_u(k) + c_h(k). \quad (8)$$

Based on input/output signal pairs, we now have a measure of the underlying generating system dynamics by looking at $c_h(k) = c_y(k) - c_u(k)$.

Definition 7. The extended cepstral distance, $d_{C_e}((y_i, u_i), (y_j, u_j))$, between two input/output pairs of time series, (y_i, u_i) and (y_j, u_j) , with respective transfer functions \mathcal{H}_i and \mathcal{H}_j , is defined as

$$d_{C_e}((y_i, u_i), (y_j, u_j)) = \sum_{k=0}^{\min\{N_i, N_j\}} k (c_{h_i}(k) - c_{h_j}(k))^2. \quad (9)$$

Note that, for now, this distance measure does not have the whole theoretic framework with connections to subspace angles, mutual information and generating system parameters.¹ However, it is clear that the $c_h(k)$ can only come from the generating system dynamics, and thus the distance measure tells us something about these systems, even if it is still unclear what exactly is measured.

We propose this extended cepstral distance as a way to efficiently cluster input/output data by their generating dynamics.

Advantages

- The extended cepstral distance is linked to the generating model of the time series.
- The extended cepstral distance is easy to calculate, allowing for very efficient computation and clustering.
- No system identification step is needed.
- This measure allows to calculate distances between time series of different length.
- This measure takes the input into account.

Disadvantages

- The interpretation of the measure in terms of system parameters and properties is not immediately clear, thus the theoretical framework of the original cepstral distance does not carry over trivially.

IV. APPLICATION ON ELECTRICAL CIRCUITS

A. Simulation Set-Up

To test the proposed techniques, we simulate data coming from electrical circuits. We start out by modelling two circuits with the same topology, but different values for the R, L, and C components. The topology was taken from a course on linear physical systems analysis [2]. The network topology and the values of the components are shown in Figure 1. The input of the system is the current i_u , the output is the voltage over L_2 , e_y . State-space models of order 3 are then written down for these networks.

We provide both systems with 200 different input signals (100 outputs of LTI models of order 15, 50 multisine waves corrupted by Gaussian white noise with standard deviation of 0.1 and 50 white noise signals), and measure the output signals. This generates a dataset of 400 input/output signal pairs (200 inputs times 2 models). The question at hand is whether we can use this input/output data, and only this

¹These theoretical equivalences will be researched and most of them proven to carry over in a forthcoming paper, where we will also try to connect the extended cepstral distance to an extended cepstral model norm.

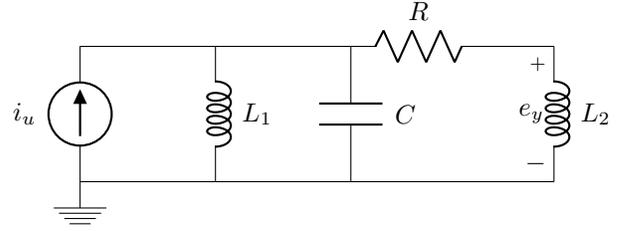


Fig. 1. Electric circuit that was used for the experiments. Two sets, S_1 and S_2 , of values were chosen for the components, namely $S_1 = \{R = 100\Omega, L_1 = 60\text{H}, L_2 = 20\text{H}, C = 50\text{F}\}$ and $S_2 = \{R = 100\Omega, L_1 = 160\text{H}, L_2 = 200\text{H}, C = 75\text{F}\}$. These two electrical circuits were used to perform the simulations in Section IV.

data, to determine which pairs were generated by the same system, i.e. cluster the dataset in two groups, defined by the generative dynamics.

We will do this using the distance measures defined in section II and subsection III-A, keeping in mind that we use the Keogh Lower Bound [5] as an efficient approximation to DTW. We then compare to the technique developed in subsection III-B. There, the power spectral density is estimated by using Welch's method [13], which provides a stable approximation² of the Fourier transform for short time series. In the Appendix, we give a pseudo-code overview of how the distance measure is calculated, as well as a link to a minimal working example of the simulations discussed and a complexity analysis of the algorithm.

The performance of these simulations will be measured by the Adjusted Rand Index (ARI) [4], [10], which is a similarity measure between partitions. The ARI compares two partitions, S_1 and S_2 , by calculating the ratio of pairs that have the same partitioning status (i.e. belonging to the same partition or not) in both S_1 and S_2 to the total amount of data pairs, then adjusting the resulting ratio by subtracting the expected value, to account for guessing (i.e. a partitioning that is the result of random guessing is assigned an ARI of 0). An ARI of 1 corresponds to perfectly similar partitions.

We compare the partitions generated by a hierarchical clustering method, cut-off at two clusters, using distance matrices generated by the different distance measures of section II and section III versus the ground truth (i.e. the time series was generated by the system with parameters S_1 or with parameters S_2 , as in Figure 1).

B. Results

The results for the set-up in the previous subsection are shown in Figure 2, which shows the average and standard deviation for the ARI of the simulation results, and Figure 3, which shows the average and standard deviation for the execution time of the simulations.

It is clear that the extended cepstral distance gives the best results. In fact, it manages to cluster the simulated input/output pairs perfectly every time. This is, of course, to be

²Note that, for longer time series (i.e. 2^{10} and beyond), the Fast Fourier Transform [1] provides a clean enough output to work on. We could thus speed up the algorithm even further for longer series.

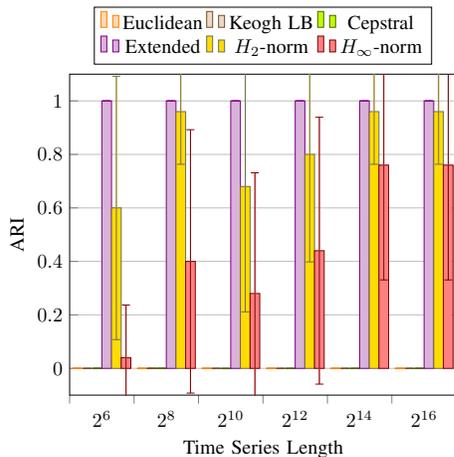


Fig. 2. Performance of the different clustering algorithms, as measured by the ARI. For each time series length, shown on the x-axis, the average ARI over 100 experiments of finding 2 clusters in 400 time series is depicted as the height of the bar. The error bars show the standard deviation for the performance on these 100 experiments. Note that the Euclidean, Keogh LB and cepstral distance have an ARI of 0, i.e., they amount to random guessing. The extended cepstral distance performs best for all series lengths. The model based distances were given a wrong model order, but still give good performance for longer time series.

expected, as this distance measure was tailored specifically to take into account the dynamics of the underlying model³, and nothing but those dynamics. The reasons why it performs better than the other measures will be explained in what follows, and we will again use the distinction between raw data and model-based distances measures from Section II.

1) Raw Data Distance Measures:

The reason why the other raw data distance measures do not perform well on the problem at hand, is because they do not take into account the information from the input signal. Indeed, the dynamics of the output are dominated by the input, due to the way the inputs were designed (i.e. the models generating the inputs are of higher order than the models describing the electrical circuits). The other distance measures are thus dominated by contributions coming from the input to cluster the time series, as they cannot separate the different contributions.

If we only use white noise inputs, we see, on the left hand side in Figure 4, that the original cepstral distance performs better.⁴ The Euclidean and DTW distances still do not deliver good results when detecting the difference in dynamics.

There is thus no hope to achieve better results by taking the input signal into account in the case of the Euclidean distance or the DTW distance. Indeed, the distances look at the shape of the signal, rather than its generative dynamics. DTW is better at this job [5], but, as we can see from Figure 3, also

³We redid the experiments for generating systems of higher order, and the extended cepstral distance still performed best. Results were omitted.

⁴In fact, the original and extended cepstral distance are equivalent in this case. Indeed, the cepstrum of white noise is only non-zero in its zeroth component, which is not taken into account in the sum in equations (7) and (9), which coincide in that case.

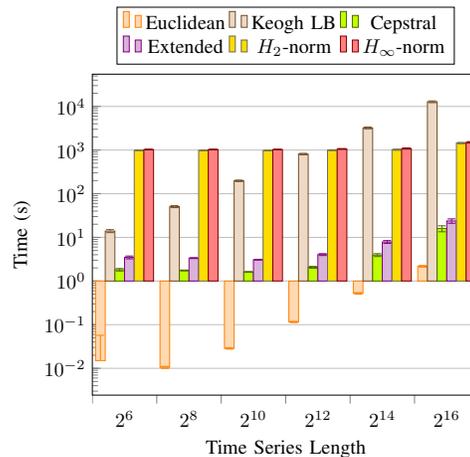


Fig. 3. Execution time of the different clustering algorithms, measured in seconds. For each time series length, shown on the x-axis, the average time over 100 experiments of finding 2 clusters in 400 time series is depicted as the height of the bar. The error bars show the standard deviation for the execution time on these 100 experiments. Note that the y-axis has logarithmic scale. The extended cepstral distance remains several orders of magnitudes faster than the model-based distances. Note that Keogh LB quickly becomes the computationally most expensive technique. The Euclidean distance is always fastest.

has a big disadvantage: it takes a lot of time to compute, especially for long time series, where it even surpasses the model-based distance measures in computation time.

Based on these results, the extended cepstral distance is thus preferred to cluster input/output signals based on the dynamics of their generating models.

2) Model-based Distance Measures:

The model-based distance measures show better results than the raw data distance measures, and this again is to be expected. Indeed, the model-based measures take the input information into account and thus manage to peel out the information on the system that generated the input/output pair. However, since a priori we have no information on the order of the underlying system, we arbitrarily have to set a model order. In this case, we estimated transfer functions of order 5. If we share the information on the correct model order (3) with the system identification algorithm, the performance of the model norms increases, as on the right hand side of Figure 4.

There exist, of course, schemes to determine appropriate model orders, and more effort can be put in correctly identifying the underlying model. However, as can be seen from Figure 3, the model norm techniques are already several orders of magnitude slower than the extended cepstrum distance measure. For problems concerning large amounts of long input/output-pairs, as can be found in realistic problems in process industry (see, for example, [7], where more than 250 sensors make a measurement every 5 minutes for 6 months), this becomes highly impractical.

The extended cepstral distance is thus preferred over explicitly identifying systems, because of both being easier to automate, and taking less time to compute.

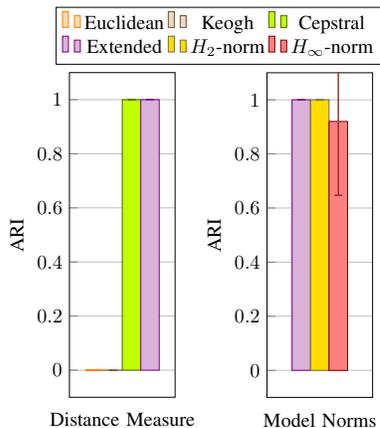


Fig. 4. On the left, the performance is shown of the different raw data distance measures, as measured by the Adjusted Rand Index (ARI), in the case of white noise as an input, and time series of length 2^{10} . Here, the average over 100 experiments with 400 output signals is shown. Note that the original cepstral distance now shows the same performance as the extended one. On the right, results of an experiment where we provided the system identification step with the correct orders of the models are shown. Here, we calculated an average over 100 experiments with 40 output signals, to reduce computation time. Again, we simulated time series of length 2^{10} , the model-based distances now show better performance.

V. CONCLUSION AND FURTHER RESEARCH

We have devised a distance measure that is as insightful as a model norm-based distance, yet remains computationally much simpler than explicitly estimating models. It allows to meaningfully cluster large input/output signal pair datasets based exclusively on the dynamics of the generating systems. We have tested it on a simulation of data coming from electrical circuits, where we started from two electrical circuits with a current as input and a voltage difference over an inductor as output. We provided both circuits with 200 different inputs, resulting in 400 input/output pairs.

We then showed that the proposed measure performs as well as model-based distances on estimates of the generative systems, but is much easier to calculate and that other distance measures (Euclidean, DTW) perform much worse.

We furthermore show that, in the stochastic input case, the extended distance proposed in this paper reduces to the original cepstrum distance, which was proven ([3], [8]) to be equivalent to a model norm. This gives hope that the extended distance could also be linked to a model norm. Research that looks into this link is currently under way and will be discussed in a forthcoming paper.

The results indicate the extended cepstral distance measure does a good job of capturing the dynamics of input/output pairs. An application to a real-life dataset is needed to validate the effectiveness in practice, but for the simulated problem at hand, the distance measure succeeded in perfectly distinguishing different dynamics based on raw data alone.

APPENDIX

A pseudo-code overview of the algorithm is shown in Algorithm 1. A minimal working example of the simulations performed in Section IV is available on GitHub.⁵

Algorithm 1: Algorithm for the extended cepstral distance

input : Two input/output signal pairs, (y_1, u_1) of length N_1 , and (y_2, u_2) of length N_2
output: The extended cepstral distance $d_{C_e}((y_1, u_1), (y_2, u_2))$ between these two pairs, as defined in Subsection III-B

```

1 for  $i \leftarrow 1$  to 2 do
2    $\Phi_{u_i} \leftarrow$  Welch's Method  $u_i$ 
3    $c_{u_i} \leftarrow$  ifft( $\log(\Phi_{u_i})$ )
4    $\Phi_{y_i} \leftarrow$  Welch's Method  $y_i$ 
5    $c_{y_i} \leftarrow$  ifft( $\log(\Phi_{y_i})$ )
6   //  $c_{u_i}$  and  $c_{y_i}$  are vectors of length  $N_i$ 
7 end
8  $w = [0, 1, \dots, \max\{N_1, N_2\} - 1]$ 
9 add  $(\max\{N_1, N_2\} - \min\{N_1, N_2\})$  0's to the cepstra of the
   signal pair of length  $\min\{N_1, N_2\}$ 
10  $d_{C_e}((y_1, u_1), (y_2, u_2)) \leftarrow w * ((c_{y_1} - c_{u_1})^\top - (c_{y_2} - c_{u_2})^\top)^2$ 

```

Calculating the extended cepstral distance amounts to estimating the power spectral density of both input and output by Welch's method [13] (employing the FFT, which is of $\mathcal{O}(n \log n)$, with n the length of the windows considered in Welch's method), taking the logarithm of the resulting vector, and then applying an inverse Fourier transform (employing the IFFT, running in $\mathcal{O}(N \log N)$ time, with N the length of the time series) on them. In the end, we then apply a weighted Euclidean distance on the results.

The complexity of calculating the extended cepstral distance between two time series is thus $\mathcal{O}(N \log N)$, with N the length of the time series.

REFERENCES

- [1] E. Oran Brigham. *The Fast Fourier Transform and Its Applications*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1988.
- [2] E. Cheever. Linear physical systems analysis. Swarthmore College Department of Engineering, <http://lpsa.swarthmore.edu/>, retrieved on 04/22/2016.
- [3] K. De Cock and B. De Moor. Subspace angles between arma models. *Systems & Control Letters*, 46(4):265–270, 2002.
- [4] L. Hubert and P. Arabie. Comparing partitions. *Journal of classification*, 2(1):193–218, 1985.
- [5] E. Keogh. Exact indexing of dynamic time warping. In *Proceedings of the 28th international conference on Very Large Data Bases*, pages 406–417. VLDB Endowment, 2002.
- [6] T. W. Liao. Clustering of time series data survey. *Pattern recognition*, 38(11):1857–1874, 2005.
- [7] L. Martí, N. Sanchez-Pi, J. M. Molina, and A. C. B. Garcia. Anomaly detection based on sensor data in petroleum industry applications. *Sensors*, 15(2):2774–2797, 2015.
- [8] R. J. Martin. A metric for arma processes. *IEEE transactions on Signal Processing*, 48(4):1164–1170, 2000.
- [9] A. V. Oppenheim and R. W. Schaffer. *Digital signal processing*. Englewood Cliffs, New York, 1975.
- [10] W. M. Rand. Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical association*, 66(336):846–850, 1971.
- [11] C. A. Ratanamahatana and E. Keogh. Everything you know about dynamic time warping is wrong. In *Third Workshop on Mining Temporal and Sequential Data*. Citeseer, 2004.
- [12] T. Vafeiadis, S. Krinidis, C. Ziogou, D. Ioannidis, S. Voutetakis, and D. Tzovaras. Robust malfunction diagnosis in process industry time series. In *2016 IEEE 14th International Conference on Industrial Informatics (INDIN)*, pages 111–116, July 2016.
- [13] P. Welch. The use of fast fourier transform for the estimation of power spectra: a method based on time averaging over short, modified periodograms. *IEEE Transactions on audio and electroacoustics*, 15(2):70–73, 1967.

⁵<https://github.com/Olauwers/Extended-Cepstral-Distance>