



Alternating diffusion maps for multimodal data fusion

Ori Katz^{*,a}, Ronen Talmon^a, Yu-Lun Lo^b, Hau-Tieng Wu^{c,d,e}

^a Viterbi Faculty of Electrical Engineering, Technion, Israel Institute of Technology, Israel

^b Department of thoracic medicine, Chang Gung Memorial Hospital, Chang Gung University, School of Medicine, Taipei, Taiwan

^c Department of Mathematics, University of Toronto, Ontario, Canada

^d Department of Mathematics and department of statistical science, Duke university, Durham, NC, USA

^e Mathematics Division, National Center for Theoretical Sciences, Taipei, Taiwan

ARTICLE INFO

Keywords:

Common manifold learning
Multimodal sensor fusion
Nonlinear-filtering
Diffusion maps
Alternating diffusion maps

ABSTRACT

The problem of information fusion from multiple data-sets acquired by multimodal sensors has drawn significant research attention over the years. In this paper, we focus on a particular problem setting consisting of a physical phenomenon or a system of interest observed by multiple sensors. We assume that all sensors measure some aspects of the system of interest with additional sensor-specific and irrelevant components. Our goal is to recover the variables relevant to the observed system and to filter out the nuisance effects of the sensor-specific variables. We propose an approach based on manifold learning, which is particularly suitable for problems with multiple modalities, since it aims to capture the intrinsic structure of the data and relies on minimal prior model knowledge. Specifically, we propose a nonlinear filtering scheme, which extracts the hidden sources of variability captured by two or more sensors, that are independent of the sensor-specific components. In addition to presenting a theoretical analysis, we demonstrate our technique on real measured data for the purpose of sleep stage assessment based on multiple, multimodal sensor measurements. We show that without prior knowledge on the different modalities and on the measured system, our method gives rise to a data-driven representation that is well correlated with the underlying sleep process and is robust to noise and sensor-specific effects.

1. Introduction

Often, when measuring a phenomenon of interest that arises from a complex dynamical system, a single data acquisition method is not capable of capturing its entire complexity and characteristics, and it is usually prone to noise and interferences. Recently, due to technological advances, the use of multiple types of measurement instruments and sensors have become more and more popular; nowadays, such equipment is smaller, less expensive, and can be mounted on every-day products and devices more easily. In contrast to a single sensor, multimodal sensors may capture complementary aspects and features of the measured phenomenon, and may enable us to extract a more reliable and detailed description of the measured phenomenon.

The vast progress in the acquisition of multimodal data calls for the development of analysis and processing tools, which appropriately combine data from the different sensors and handle well the inherent challenges that arise. One particular challenge is related to the heterogeneity of the data acquired in the different modalities; datasets acquired from different sensors may comprise different sources of variability, where only few are relevant to the phenomenon of interest.

This particular challenge as well as many others have been the subject of many studies. For a recent comprehensive reviews, see [1–3].

In this paper we consider a setting in which a physical phenomenon is measured by multiple sensors. While all sensors measure the same phenomenon, each sensor consists of different sources of variability; some are related to the phenomenon of interest, possibly capturing its various aspects, whereas other sources of variability are sensor-specific and irrelevant. We present an approach based on manifold learning, which is a class of nonlinear data-driven methods, e.g. [4–7], and specifically, we use the framework of diffusion maps (DM) [8]. On the one hand, manifold learning is particularly suitable for problems with multiple modalities since it aims to capture the intrinsic geometric structure of the underlying data and relies on minimal prior model knowledge. This enables to handle multimodal data in a systematic manner, without the need to specially tailor a solution for each modality. On the other hand, applying manifold learning to data acquired in multiple (multimodal) sensors may capture undesired/nuisance geometric structures as well. Recently, several manifold learning techniques for multimodal data have been proposed [9–12]. In [9], the authors suggest to concatenate the samples acquired by different sensors

* Corresponding author.

E-mail address: orikats@tx.technion.ac.il (O. Katz).

into unified vectors. However this approach is sensitive to the scaling of each dataset, which might be especially diverse among datasets acquired by different modalities. To alleviate this problem, it is proposed in [10] to use DM to obtain “standardized” representation of each dataset separately, and then to concatenate these “standardized” representations into the unified vectors. Despite handling better multimodal data, this concatenation scheme does not utilize the mutual relations and co-dependencies that might exist between the datasets.

While methods such as those presented in [9,10,12] take into account all the measured information, the methods presented in [11,13–15] use local kernels to implement nonlinear filtering. Specifically, following a recent line of study in which multiple kernels are constructed and combined [16–19], in [13,14], it was shown that a method based on alternating applications of diffusion operators extracts only the common source of variability among the sensors, while filtering out the sensor-specific components. Therefore we choose to establish our method based on DM which relies on those theoretical foundations. For other nonlinear methods, such as local CCA [11] and kernel CCA [15], more efforts are needed to better understand their theoretical foundation, yet, they may also be used as an alternative to DM and alternating diffusion (AD) and empirically tested as well. The shortcoming of alternating applications of diffusion operators arises when having a large number of sensors; often, sensors that measure the same system capture different information and aspects of that system. As a result, the common source of variability among all the sensors captures only a partial or empty look of the system, and important relevant information may be undesirably filtered out.

Here, we address the tradeoff between these two approaches. That is, we aim to maintain the relevant information captured by multiple sensors, while filtering out the nuisance components. Since the relevance of the various components is unknown, our main assumption is that the sources of variability which are measured only in a single sensor, i.e., sensor-specific, are nuisance. Conversely, we assume that components measured in two or more sensors are of interest. Importantly, such an approach implements implicitly a smart “sensor selection”; “bad” sensors that are, for example malfunctioned and measure only nuisance information, are automatically filtered out. These assumptions stem from the fact that the phenomenon of interest is global and not specific to one sensor. We propose a nonlinear filtering scheme, in which only the sensor-specific sources of variability are filtered out while the sources of variability captured by two or more sensors are preserved.

Based on prior theoretical results [13,14], we show that our scheme indeed accomplishes this task. We illustrate the main features of our method on a toy problem. In addition, we demonstrate its performance on real measured data in an application for sleep stage assessment based on multiple, multimodal sensor measurements. Sleep is a global phenomenon with systematic physiological dynamics that represents a recurring non-stationary state of mind and body. Sleep evolves in time and embodies interactions between different subsystems, not solely limited in the brain. Thus, in addition to the well-known patterns in electroencephalogram (EEG) signals, its complicated dynamics are manifested in other sensors such as sensors measuring breathing patterns, muscle tones and muscular activity, eyeball movements, etc. Each one of the sensors is characterized by different structures and affected by numerous nuisance processes as well. In other words, while we could extract the sleep dynamics by analyzing different sensors, each sensor captures only part of the entire sleep process, whereas it introduces modality artifacts, noise, and interferences. We show that our scheme allows for an accurate systematic sleep stage identification based on multiple EEG recordings as well as multimodal respiration measurements. In addition, we demonstrate its capability to perform sensor selection by artificially adding noise sensors.

The remainder of the paper is organized as follows. In Section 2 we present a formulation for the common source extraction problem and present an illustrative toy problem. In Section 3, a brief review for the

method proposed in [13,14] is outlined, and then, a detailed description and interpretation of the proposed scheme are presented. In Section 4, we first demonstrate the capabilities of the proposed scheme on the toy problem introduced in Section 2. Then, in Section 5, we demonstrate the performance in sleep stage identification based on multimodal measured data recorded in a sleep clinic. Finally, in Section 6, we outline several conclusions.

2. Problem setting

Consider a system driven by a set of K hidden random variables $\Theta = \{\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}\}$, where $\theta^{(k)} \in \mathbb{R}^{d_k}$. The system is measured by M observable variables $s^{(m)}$, $m = 1, \dots, M$, where each sensor has access to only a partial view of the entire system and its driving variables Θ . To formulate it, we define a “sensitivity table” given by the binary matrix $\mathbf{S} \in \mathbb{Z}_2^{K \times M}$, indicating the variables sensed by each observable variable. Specifically, the (k,m) th element in \mathbf{S} indicates whether the hidden variable $\theta^{(k)}$ is measured by the observable variable $s^{(m)}$. It should be noted that this binary notation is a simplification that does not take into account “soft degrees of observability”, e.g., the nonlinearity of the observation function, measurement noise, etc. Theoretically, when we have sufficient data, the algorithm is guaranteed to work for any bilipschitz observation function. Yet, when the amount of data is limited, these degrees of observability are dominated by the Lipschitz constants of the observation function and by the signal-to-noise ratio. Some of these observability aspects were addressed in [20], and further quantification is left for future work. The observable variables are therefore given by

$$s^{(m)} = h_m(\Theta^{(m)}, \mathbf{n}^{(m)}) \in \mathbb{R}^{D_m} \quad (1)$$

where $h_m(\cdot)$ is a bilipschitz observation function, $\mathbf{n}^{(m)} \in \mathbb{R}^{p_m}$ are hidden random variables captured only by the m th observable variable, and $\Theta^{(m)}$ is the subset of driving hidden variables of interest sensed by $s^{(m)}$, given by

$$\Theta^{(m)} = \{\theta^{(k)} \mid \forall k, S_{k,m} = 1\} \subseteq \Theta, m = 1, \dots, M \quad (2)$$

The random hidden variables $\mathbf{n}^{(m)}$ are *sensor-specific* (associated only with the m th observer). They are conditionally independent given the hidden variables of interest and will be assumed as noise/nuisance variables. We further assume that each random hidden variable in Θ is measured by at least two observable variables, such that $\sum_{m=1}^M S_{k,m} \geq 2$ for each $k = 1, \dots, K$. As a result, we refer to the hidden variables $\theta^{(k)}$ in Θ as *common variables*.

In order to simplify the notation, we denote the subset of all hidden variables (both common and sensor-specific) measured by the m th observable by $\mathcal{S}^{(m)} = \{\Theta^{(m)}, \mathbf{n}^{(m)}\}$. Furthermore, we assume that the dimensions of the observations and the hidden variables satisfy

$$D_m \geq \sum_{\theta^{(k)} \in \Theta^{(m)}} (d_k + p_k), k = 1, 2, \dots, M \quad (3)$$

i.e., the observations are in higher dimension than the hidden common and nuisance variables.

An observation of the system denoted as $(s_i^{(1)}, s_i^{(2)}, \dots, s_i^{(M)})$ is associated with a realization of the hidden variables $\Theta_i = (\theta_i^{(1)}, \dots, \theta_i^{(K)})$ and realizations of the M hidden nuisance variables $(\mathbf{n}_i^{(1)}, \dots, \mathbf{n}_i^{(M)})$. Given N observation samples $\{(s_i^{(1)}, s_i^{(2)}, \dots, s_i^{(M)})\}_{i=1}^N$, our goal is to obtain a parametrization for the underlying realizations of the common hidden random variables $\{(\theta_i^{(1)}, \dots, \theta_i^{(K)})\}_{i=1}^N$ while filtering out the nuisance variables $\{(\mathbf{n}_i^{(1)}, \dots, \mathbf{n}_i^{(M)})\}_{i=1}^N$. We note that the observations index i may represent the time index in case of time series.

2.1. Illustrative toy problem

We illustrate the problem setting using the following toy example. Consider six rotating arrows captured in simultaneous snapshots by

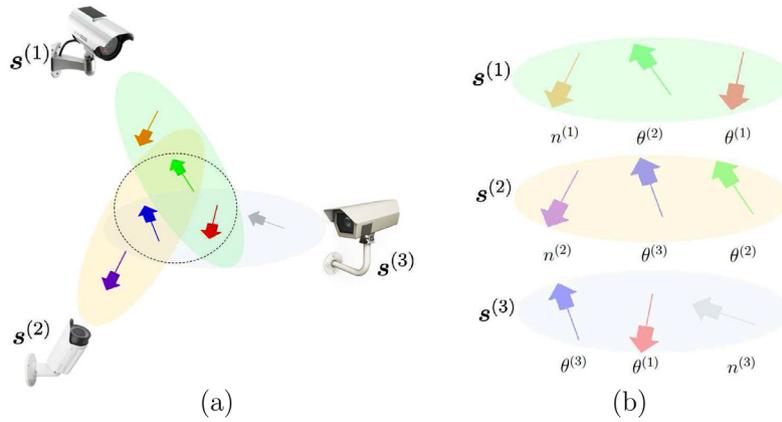


Fig. 1. Toy problem setup. (a) The coverage area of each camera, the system's range of interest is marked by the dashed circle. (b) Sample snapshot taken by each camera.

three different cameras. We assume that each arrow rotates at different speed, and that each camera can capture only a partial image of the entire system. The partial view of each camera is depicted in Fig. 1. Thus, overall, each camera captures a sequence of snapshots (a movie) of three rotating colored arrows. Further illustration of the entire system and of the captured images by each camera can be seen in the following link <https://youtu.be/a-yb7Scdnna>.

In this problem setting, the hidden variables are the six rotation angles of the arrows: the common variables $\Theta = \{\theta^{(1)}, \theta^{(2)}, \theta^{(3)}\}$ are the three rotation angles of the centred arrows, which are marked by the dashed circle in Fig. 1, and the nuisance variables $\{n^{(1)}, n^{(2)}, n^{(3)}\}$ are the three rotation angles of the peripheral arrows, since each is captured only by a single camera. It should be noted that none of the arrows is common to all of the cameras, meaning that the set of common components within the entire set of observables is empty.

In order to identify the hidden variables, we use different colors for the arrows. The arrows rotating according to the common variables $\Theta = \{\theta^{(1)}, \theta^{(2)}, \theta^{(3)}\}$ are colored in red, green and blue, respectively, and the arrows rotating according to the nuisance variables $\{n^{(1)}, n^{(2)}, n^{(3)}\}$ are colored in orange, purple and gray, respectively. The hidden variables measured by each camera are $\mathcal{S}^{(1)} = \{\theta^{(1)}, \theta^{(2)}, n^{(1)}\}$, $\mathcal{S}^{(2)} = \{\theta^{(2)}, \theta^{(3)}, n^{(2)}\}$ and $\mathcal{S}^{(3)} = \{\theta^{(3)}, \theta^{(1)}, n^{(3)}\}$. Our goal is to obtain a parametrization of the rotation angles of the three common arrows $\Theta = \{\theta^{(1)}, \theta^{(2)}, \theta^{(3)}\}$ given the three movies of the cameras, without any prior knowledge on the system and the problem structure. In the sequel, we will use this toy problem for demonstrating important aspects and how our method accomplishes this task.

3. Nonlinear filtering scheme

3.1. Diffusion maps

DM is a non-linear data-driven dimensionality reduction method [8]. Assume we have N high-dimensional data-points $\{s_i\}_{i=1}^N$. The DM method begins with the calculation of a pairwise affinity matrix based on a local kernel, often using some metric within a gaussian kernel, i.e.,

$$W_{i,j} = \exp\left(-\frac{d_M(s_i^{(1)}, s_j^{(1)})^2}{\varepsilon}\right), \quad (4)$$

where $\varepsilon > 0$ is a tuneable kernel scale and $d_M(\cdot, \cdot)$ is a metric. The choice of the metric $d_M(\cdot, \cdot)$ depends on the application; common choices are the Euclidean and the Mahalanobis distances [8,21–24]. This construction implicitly defines a weighted graph, where the data samples $\{s_i\}_{i=1}^N$ are the nodes of the graph, and $W_{i,j}$ is the weight of the edge connecting node s_i and node s_j . The next step is to normalize the affinity matrix and then to build the diffusion operator $\mathbf{K} \in \mathbb{R}^{N \times N}$, e.g., by:

$$Q_{i,i} = \left(\sum_{l=1}^N W_{i,l}\right)^{-1}; \mathbf{K} = \mathbf{Q}\mathbf{W}, \quad (5)$$

where \mathbf{Q} is a diagonal matrix used for normalization, such that in this case \mathbf{K} is row-stochastic. Hence, \mathbf{K} can be viewed as the transition matrix of a Markov chain defined on the graph. Accordingly, for $t > 0$, \mathbf{K}^t is the transition probability matrix of t consecutive steps, and $(\mathbf{K}^t)_{i,j}$ is the probability to jump from node s_i to node s_j in t steps. Let $d_t(i, j)$ be the diffusion distance [8] between the i th and the j th data samples, i.e. d_t is a function defined by

$$d_t(i, j) = \sqrt{\sum_{l=1, \dots, N} \frac{((\mathbf{K}^t)_{i,l} - (\mathbf{K}^t)_{j,l})^2}{\phi_0(l)}} \quad (6)$$

where $\phi_0(\cdot)$ is the stationary distribution of the Markov chain. The diffusion distance has been shown to be a powerful metric for measuring geometrical similarities between data-points [8]. While the Euclidean distance compares two individual data-points and might be affected by distortions and noise, the diffusion distance introduces much more noise-robust affinities since it relies on the connectivity between the two data-points using the entire data-set [8,25].

However, the direct computation of the diffusion distance is cumbersome. An efficient calculation is attainable via the spectral decomposition of \mathbf{K} . Let $\{\lambda_l\}_{l=0}^{N-1}$ and $\{\psi_l\}_{l=0}^{N-1}$ be the sets of eigenvalues and right eigenvectors of \mathbf{K} , where the eigenvalues are in descending order. Define a new representation (embedding) of the data-points:

$$\Psi_t(i): s_i \mapsto [\lambda_1^t \psi_1(i), \lambda_2^t \psi_2(i), \dots, \lambda_{N-1}^t \psi_{N-1}(i)], \quad (7)$$

where $\psi_l(i)$ denotes the i th element of ψ_l . The obtained embedding provides a new representation of the data, referred to as DM, in which the Euclidean distance between two embedded data-point is equal to the diffusion distance [8], i.e.:

$$d_t^2(i, j) = \|\Psi_t(i) - \Psi_t(j)\|^2 = \sum_{l \geq 1} \lambda_l^{2t} (\psi_l(i) - \psi_l(j))^2. \quad (8)$$

In order to achieve a compact representation in reduced dimensionality, DM is often redefined by keeping only the first L components (i.e., the L eigenvalues and eigenvectors corresponding to the largest L eigenvalues):

$$\Psi_t(i): s_i \mapsto [\lambda_1^t \psi_1(i), \lambda_2^t \psi_2(i), \dots, \lambda_L^t \psi_L(i)], \quad (9)$$

where L is usually determined by the eigenvalues decay. For more details and full analysis of this algorithm see [8,26]. The entire DM method is outlined in Algorithm 1.

The term “diffusion distance” in (6) suggests that $d_t(i, j)$ induces a reasonable notion of distance. Recall the definition of a distance.

Definition 1. Let X be a set. A *distance* (or *metric*) on X is a function

Input: High-dimensional samples from an observable variables: $\{s_i\}_{i=1}^N$.

Output: L dimensional representation of the data-set $\{\Psi_t(i)\}_{i=1}^N$ where $\Psi_t(i) \in \mathbb{R}^L$.

1. Calculate the affinity matrix \mathbf{W} :

$$W_{i,j} = \exp\left(-\frac{d_M(s_i, s_j)^2}{\varepsilon}\right)$$

2. Compute the diffusion operator (transition matrix) \mathbf{K} :

$$Q_{i,i} = \left(\sum_{l=1}^N W_{i,l}\right)^{-1}; \mathbf{K} = \mathbf{QW},$$

3. Calculate the spectral decomposition of \mathbf{K} and obtain its eigenvalues $\{\lambda_t\}_{t=0}^{N-1}$ and eigenvectors $\{\psi_t\}_{t=0}^{N-1}$.
4. Define a new embedding for the data-points:

$$\Psi_t(i) : s_i \mapsto [\lambda_1^t \psi_1(i), \lambda_2^t \psi_2(i), \dots, \lambda_L^t \psi_L(i)]$$

where $t > 0$ is a selected number of steps and $\psi_t(i)$ denotes the i th element of ψ_t .

Algorithm 1. Diffusion Maps.

$d: X \times X \rightarrow \mathbb{R}_+$ such that for all $x, y, z \in X$:

1. $d(x, y) = 0$ if and only if $x = y$,
2. $d(x, y) = d(y, x)$,
3. $d(x, z) \leq d(x, y) + d(y, z)$.

The following proposition states that the “diffusion distance” is really a metric defined on the nodes of the graph.

Proposition 2. *If \mathbf{K} is full rank, then d_t is a distance function.*

Since we could not find a proof in the literature, for the sake of self-containment, we provide a proof that summarizes the discussion in [27].

Proof. We prove that (1)–(3) hold. Define $\tilde{\mathbf{K}}^t = \Phi^{-1}\mathbf{K}^t$ where Φ is a diagonal matrix such that $\Phi_{k,k} = \sqrt{\phi_0(k)}$. Since \mathbf{K}^t is full rank and since Φ is non-degenerate by the construction of the weighted graph, $\tilde{\mathbf{K}}^t$ is full rank. Denote the i th row of $\tilde{\mathbf{K}}^t$ as \mathbf{v}_i . Accordingly, $d_t(i, j)$ can be expressed as the Euclidean distance between the i th and the j th rows of $\tilde{\mathbf{K}}^t$:

$$d_t(i, j) = \sqrt{\sum_{l=1, \dots, N} ((\tilde{\mathbf{K}}^t)_{i,l} - (\tilde{\mathbf{K}}^t)_{j,l})^2} = \|\mathbf{v}_i - \mathbf{v}_j\|_{\mathbb{R}^N} \quad (10)$$

The properties of the Euclidean distance in (10) imply that (2) and (3) hold. If $i = j$, then $d_t(i, j) = 0$. If $d_t(i, j) = 0$, then $\|\mathbf{v}_i - \mathbf{v}_j\|^2 = 0$, implying that $\mathbf{v}_i = \mathbf{v}_j$. Since $\tilde{\mathbf{K}}^t$ is full rank, there are no identical columns. In other words, no two different samples \mathbf{v}_i and \mathbf{v}_j for $i \neq j$ have identical affinities to all other samples, i.e., $\mathbf{v}_i \neq \mathbf{v}_j$. Therefore, if $\mathbf{v}_i = \mathbf{v}_j$, then $i = j$. \square

3.2. Alternating diffusion

Consider a system similar to the one described in Section 2, with only $M = 2$ observable variables and $K = 1$ common variable. The AD algorithm, outlined in Algorithm 2, builds from the observations an AD operator that is equivalent to a simple diffusion operator (as described in Section 3.1) that would have been computed if we had a direct access to samples of the common hidden variables. This operator enables to capture only the structure of the common variables while ignoring the nuisance (sensor-specific) variables. For more details and full analysis of this algorithm see [13,14]; here, we only bring a brief review of the method and the construction of the AD operator.

Assume we have N aligned samples (realizations) from 2 observable variables: $\{(\mathbf{s}_i^{(1)}, \mathbf{s}_i^{(2)})\}_{i=1}^N$. For each observation we build an affinity matrix $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$ as follows:

$$W_{ij}^{(1)} = \exp\left(-\frac{d_{M_1}(\mathbf{s}_i^{(1)}, \mathbf{s}_j^{(1)})^2}{\varepsilon^{(1)}}\right); W_{ij}^{(2)} = \exp\left(-\frac{d_{M_2}(\mathbf{s}_i^{(2)}, \mathbf{s}_j^{(2)})^2}{\varepsilon^{(2)}}\right) \quad (11)$$

for all $i, j = 1, \dots, N$, where $\varepsilon^{(1)}$ and $\varepsilon^{(2)}$ are the tuneable kernel scales and $d_{M_1}(\cdot, \cdot)$ and $d_{M_2}(\cdot, \cdot)$ are the chosen metrics for each set of observations. Based on the affinity matrix, we calculate the diffusion operators $\mathbf{K}^{(1)}$ and $\mathbf{K}^{(2)}$ according to:

$$\mathbf{Q}_{i,i}^{(1)} = \left(\sum_{l=1}^N W_{i,l}^{(1)}\right)^{-1}; \mathbf{Q}_{i,i}^{(2)} = \left(\sum_{l=1}^N W_{i,l}^{(2)}\right)^{-1}$$

$$\mathbf{K}^{(1)} = \mathbf{Q}^{(1)}\mathbf{W}^{(1)}; \mathbf{K}^{(2)} = \mathbf{Q}^{(2)}\mathbf{W}^{(2)}$$

where $\mathbf{Q}^{(1)}$ and $\mathbf{Q}^{(2)}$ are diagonal matrices used for normalization. Next, we define $\mathbf{K}^{(1)\cap(2)} = \mathbf{K}^{(1)}\mathbf{K}^{(2)}$ as the AD operator. Note that $\mathbf{K}^{(1)\cap(2)}$ is row-stochastic, and hence, can be considered as a transition probability matrix of a new Markov chain that alternates between the two data sets. Namely, each step of this alternating process consists of a propagation step using $\mathbf{K}^{(1)}$ followed by a propagation step using $\mathbf{K}^{(2)}$.

Broadly, in each propagation step, the Markov chain jumps with high probability to neighboring samples that are similar in terms of the

kernel. Combining alternating steps results in consecutive jumps according to similarities in the first set and then according to similarities in the second set. Overall, only similarities in terms of the common components among the two views are maintained.

Formally, we define the diffusion distance between the i th and the j th sample based on the AD operator as the following Euclidean distance

$$d_t^{(1)\cap(2)}(i, j) = \sqrt{\sum_{l=1, \dots, N} \frac{(((\mathbf{K}^{(1)\cap(2)})^t)_{i,l} - ((\mathbf{K}^{(1)\cap(2)})^t)_{j,l})^2}{\phi_0^{(1)\cap(2)}(l)}} \quad (12)$$

where $\phi_0^{(1)\cap(2)}$ is the stationary distribution of $\mathbf{K}^{(1)\cap(2)}$ and $t > 0$ is the number of alternating steps. The following corollary is an immediate result of Proposition 2.

Corollary 3. *If $\mathbf{K}^{(1)\cap(2)}$ is full rank, then $d_t^{(1)\cap(2)}$ is a distance function.*

It can be shown that this distance is equivalent to the diffusion distance that would have been computed if we had a direct access to observable variables that see only the common variable [13].

3.3. Common graph

AD provides us with an access to the common variables between a pair of observable variables. By using AD as a building block, we propose a generalization for a set of multiple observable variables. Consider the system described in Section 2 with aligned samples from M observable variables: $\{(\mathbf{s}_i^{(1)}, \mathbf{s}_i^{(2)}, \dots, \mathbf{s}_i^{(M)})\}_{i=1}^N$. The observable variables are driven by a set of K hidden random variables $\Theta = (\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)})$ and contaminated by a set of M nuisance sensor-specific variables $(\mathbf{n}^{(1)}, \dots, \mathbf{n}^{(M)})$. Our goal is to obtain a parametrization of the common hidden random variables Θ from the observations.

Proposition 3 provides the analytic foundation and justification to the method presented in this paper. More specifically, in the context of our problem, consider a pair of observable variables $\mathbf{s}^{(m)}$ and $\mathbf{s}^{(n)}$. Applying AD to $\mathbf{s}^{(m)}$ and $\mathbf{s}^{(n)}$ yields the common hidden variables measured by the two. Therefore, its operation can be written as

$$\mathcal{S}^{(m)\cap(n)} = \Theta^{(m)\cap(n)} \quad (19)$$

In other words, AD captures only a subset of the common hidden variables $\Theta^{(m)\cap(n)}$, and in addition, filters out the nuisance variables $\mathbf{n}^{(m)}$ and $\mathbf{n}^{(n)}$, which are specific to each observation.

The main idea in our method is based on the fact that the desired set of variables Θ can be derived from the union of the pairwise intersections between all pairs, meaning that:

$$\Theta = \bigcup_{m \neq n} (\Theta^{(m)\cap(n)}). \quad (20)$$

A direct implementation of the scheme in (14) is not feasible, since the pairwise intersections of $\Theta^{(m)}$ and $\Theta^{(n)}$ are not accessible to us. However, note that by substituting (19) in (20), we get

$$\Theta = \bigcup_{m \neq n} (\mathcal{S}^{(m)\cap(n)}) \quad (21)$$

meaning that Θ can be expressed using the accessible observations sets $\mathcal{S}^{(m)}$ through the union of the intersections of all possible pairs. Thus, this scheme for recovering Θ can be implemented by multiple applications of AD to all possible pairs of observable variables.

The union is implemented through the formulation of a new kernel in which the affinity between each pair of samples is given by the sum of the diffusion distances over all pairs of observations. Therefore for each kernel resulting from an application of AD to a single pair of observations, we compute the following diffusion distance $d_t^{(m)\cap(n)}$, similarly to (12)

$$d_t^{(m)\cap(n)}(i, j) = \sqrt{\sum_{l=1}^N \frac{(((\mathbf{K}^{(m)\cap(n)})^t)_{i,l} - ((\mathbf{K}^{(m)\cap(n)})^t)_{j,l})^2}{\phi_0^{(m)\cap(n)}(l)}} \quad (22)$$

Input: Aligned samples from 2 observable variables: $\{(\mathbf{s}_i^{(1)}, \mathbf{s}_i^{(2)})\}_{i=1}^N$.

Output: Diffusion distances $d_t^{(1) \cap (2)}$.

1. Calculate two pairwise affinity matrices $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$ based on a gaussian kernel as follows:

$$W_{i,j}^{(1)} = \exp\left(-\frac{d_M(\mathbf{s}_i^{(1)}, \mathbf{s}_j^{(1)})^2}{\varepsilon^{(1)}}\right); \quad W_{i,j}^{(2)} = \exp\left(-\frac{d_M(\mathbf{s}_i^{(1)}, \mathbf{s}_j^{(1)})^2}{\varepsilon^{(2)}}\right)$$

for all $i, j = 1, \dots, N$, where $\varepsilon^{(1)}$ and $\varepsilon^{(2)}$ are the kernel scales and $d_M(\cdot, \cdot)^2$ is the chosen metric.

2. Create two diffusion operators $\mathbf{K}^{(1)}$ and $\mathbf{K}^{(2)}$:

$$Q_{i,i}^{(1)} = \left(\sum_{j=1}^N W_{i,j}^{(1)}\right)^{-1}; \quad Q_{i,i}^{(2)} = \left(\sum_{j=1}^N W_{i,j}^{(2)}\right)^{-1}$$

$$\mathbf{K}^{(1)} = \mathbf{Q}^{(1)} \mathbf{W}^{(1)}; \quad \mathbf{K}^{(2)} = \mathbf{Q}^{(2)} \mathbf{W}^{(2)}$$

3. Build the alternating-diffusion kernel:

$$\mathbf{K}^{(1) \cap (2)} = \mathbf{K}^{(1)} \mathbf{K}^{(2)}$$

4. Compute the alternating-diffusion distance between each two points (i, j)

$$d_t^{(1) \cap (2)}(i, j) = \sqrt{\frac{\sum_{l=1, \dots, N} \left(((\mathbf{K}^{(1) \cap (2)})^t)_{i,l} - ((\mathbf{K}^{(1) \cap (2)})^t)_{j,l} \right)^2}{\phi_0^{(1) \cap (2)}(l)}}$$

where $\phi_0^{(1) \cap (2)}$ is the stationary distribution of $\mathbf{K}^{(1) \cap (2)}$ and $t > 0$ is a tuneable parameter.

Algorithm 2. Alternating Diffusion.

Table 1
List of important notation.

Nomenclature	
K	Number of common hidden variables
$\theta^{(k)}$	k th common hidden variable
d_k	Dimension of $\theta^{(k)}$
Θ	Set of all common hidden variables
M	Number of observable variables
$s^{(m)}$	m th observable variable
D_m	Dimension of $s^{(m)}$
S	Sensitivity table
$h_m(\cdot)$	A bilipschitz m th observation function
$\mathbf{n}^{(m)}$	m th sensor-specific hidden (nuisance) variables
p_m	Dimension of $\mathbf{n}^{(m)}$
$\Theta^{(m)}$	Subset of Θ sensed by $s^{(m)}$
$\mathcal{S}^{(m)}$	Subset of all hidden variables measured by $s^{(m)}$

where $\phi_0^{(m)\cap(n)}$ is the stationary distribution of $\mathbf{K}^{(m)\cap(n)}$ and $t > 0$ is a tuneable parameter indicating the number of AD steps. We then define the *common diffusion distance* $d_t^{(u)}$ as a summation over the alternating diffusion distances (22) resulting from applications to all possible pairs of observations, according to

$$d_t^{(u)}(i, j) = \sum_{1 \leq m, n \leq M, m \neq n} d_t^{(m)\cap(n)}(i, j) \quad (23)$$

where $i, j = 1, \dots, N$. We now show that $d_t^{(u)}$ is a metric.

Proposition 4. *Let X be a set and consider two distance functions $d_1, d_2: X \times X \rightarrow \mathbb{R}_+$. Define $d(x, y) = d_1(x, y) + d_2(x, y)$ for all $x, y \in X$. Then d is a distance function as well. In particular, if $\mathbf{K}^{(m)\cap(n)}$ are full rank for all $m, n = 1, \dots, M, m \neq n$, then $d_t^{(u)}$ is a distance function.*

Proof. By definition d is $d: X \times X \rightarrow \mathbb{R}_+$. We prove that properties (1)–(3) in Definition 1 hold. Using the symmetry property of d_1 and d_2 we have that $d(x, y) = d(y, x)$. Consider $x, z \in X$, using property (3) of d_1 and d_2 , for any $y \in X$ $d(x, z) = d_1(x, z) + d_2(x, z) \leq d_1(x, y) + d_1(y, z) + d_2(x, y) + d_2(y, z) = d(x, y) + d(y, z)$. If $x = y$ then $d(x, y) = 0$. If $d(x, y) = 0$, using the non-negativity property (1) of d_1 and d_2 we have that $d_1(x, y) = 0$ and $d_2(x, y) = 0$. From property (1) we obtain that $x = y$.

Now, by Corollary 3, if $\mathbf{K}^{(m)\cap(n)}$ is full rank, then, $d_t^{(m)\cap(n)}$ is a distance function, and therefore, by a straight-forward generalization, it follows that $d_t^{(u)}$ is a distance function. \square

Based on the common diffusion distance $d_t^{(u)}$, then we calculate an affinity matrix

$$W_{ij}^{(u)} = \exp\left(-\frac{d_t^{(u)}(i, j)}{\varepsilon^{(u)}}\right), \quad (24)$$

where $\varepsilon^{(u)} > 0$ is the chosen kernel scale. Next we normalize the affinity matrix and build the common diffusion operator $\mathbf{K}^{(u)} \in \mathbb{R}^{N \times N}$

$$Q_{i,i}^{(u)} = \left(\sum_{l=1}^N W_{i,l}^{(u)}\right)^{-1}; \quad \mathbf{K}^{(u)} = \mathbf{Q}^{(u)}\mathbf{W}^{(u)} \quad (25)$$

In conclusion, the new graph with kernel $\mathbf{K}^{(u)}$ consists of two main components. First, the intersections between any pair of observations $\mathcal{S}^{(m)} \cap \mathcal{S}^{(n)}$ are implemented using AD that provides the extraction of the common hidden variables $\Theta^{(m)\cap(n)}$. Second, the union $\bigcup_{m,n} (\Theta^{(m)\cap(n)})$ is implemented via the summation of the resulting diffusion distances from the AD applications. By construction, in the kernel $\mathbf{K}^{(u)}$, the connectivity between the i th and the j th data samples is proportional to the intrinsic distance $\|\Theta_i - \Theta_j\|$. This means that the common global diffusion kernel $\mathbf{K}^{(u)}$ can be used for obtaining a low-dimensional representation of Θ . The common graph algorithm described in this section is summarized in Algorithm 4.

Four final remarks follow. First, it should be noted that there is a

theoretical gap between the desired union described in (21) and its implementation via the summation of the common diffusion distances based on the metric $d_t^{(u)}$ in (23). This particular implementation was used since the obtained embeddings based on the metric $d_t^{(u)}$ were empirically found to be well correlated with the expected results from the union scheme; the supporting empirical results are presented in Section 4, whereas rigorous analysis of the union scheme will appear in future work.

Second, the proposed implementation of the union via diffusion distance summation enhances the common variables that appear multiple times in the various intersections. By doing so, we slightly abuse the definition of the union, where duplicates are all “put together”. In other words, in the strict definition of a union, in contrast to our implementation, common hidden variables related to two or more intersection results should be taken into account only once. Depending on the application at hand, this may be a desired property, and the derivation of a scheme in which each common components has a uniform gain is postponed to future work.

Third, the proposed algorithm can be viewed from a nonlinear filtering standpoint. By applying the proposed algorithm, we maintain or even enhance the common hidden variables, while filtering out the nuisance variables that are sensor/observation-specific.

While that the previous remarks addressed the proposed implementation, and the potential theoretical gaps that should be addressed in the future, the final remark deals with the practice of computing the approximation through the implementation of the metric $d_t^{(u)}$. For the application of sleep stage identification described in Section 5, we have empirically found that by a well chosen metric $d_t^{(u)}$, which is outlined next, gives rise to improved performance. This choice of metric results in a “smoother” embedding, which better represents the sleep stage. In the alternative implementation, rather than calculating $d_t^{(u)}$ as in (29), we calculate it in the following way. First, for each pair of sensors we apply the standard DM based on the pairwise kernels $\mathbf{K}^{(m)\cap(n)}$, $1 \leq m, n \leq M, m \neq n$ computed in (17). For each pair we obtain a $L^{(m)\cap(n)}$ -dimensional representation, where $L^{(m)\cap(n)}$ is a chosen parameter for the pair (m, n) , estimated using the “spectral gap” of the decay of the eigenvalues of $\mathbf{K}^{(m)\cap(n)}$. Second, we concatenate the low-dimensional representations obtained from the previous step into a single vector. In other words, we now have N concatenated L -dimensional vectors, where $L = \sum_{1 \leq m, n \leq M, m \neq n} L^{(m)\cap(n)}$, representing the N observations taken simultaneously from all M sensors. Third, we calculate the pairwise distance $d_t^{(u)}$ between the N new concatenated vectors. Broadly, this technique is similar to [10], only here we combine the already “filtered” components (the results of AD rather than DM). Since these vectors consist of components from different sensors, we chose to use a modified version of the Mahalanobis distance. This modified Mahalanobis distance was first introduced in [21], and since then, was shown to exhibit remarkable capability to standardize measurements from different sources, e.g. in [22–24,28]. In [22,23], it was shown to build intrinsic representations by revealing a hidden process driving the measurements. Recently, this technique was applied to multimodal data in [12]. Importantly, compared with AD and our proposed method, these methods [12,22,23] combine the information embodied in all the measurements and do not attempt to suppress nuisance variables or to extract only the common components.

The numerical implementation of the Mahalanobis distance deserves a remark. The computation of the Mahalanobis distance requires estimation of the local covariance matrices of the vectors, each of size $L \times L$. This computation might be computationally cumbersome when L is large, as often in our case. In order to relax the required computational load, prior to the computation of the Mahalanobis distance, one can project the concatenated samples onto a lower dimensional vector space, for example, using random projections (RPs) [29], and then, compute the Mahalanobis distance for the projected samples with reduced dimensionality. This heuristic method for calculating the common diffusion $d_t^{(u)}$ is summarized in Algorithm 3.

Input: $M(M - 1)$ alternating-diffusion operators $\mathbf{K}^{(m) \cap (n)}$, $1 \leq m, n \leq M, m \neq n$

Output: Alternating-diffusion distance $d_t^{(u)}$

1. For $1 \leq m, n \leq M, m \neq n$, calculate the spectral decomposition of each kernel $\mathbf{K}^{(m) \cap (n)}$, and obtain its eigenvalues $\{\lambda_l^{(m) \cap (n)}\}_{l=0}^{N-1}$ and eigenvectors $\{\psi_l^{(m) \cap (n)}\}_{l=0}^{N-1}$.
2. For $1 \leq m, n \leq M, m \neq n$, build an $L_{(m) \cap (n)}$ -dimensional representation using standard DM (9) for each time sample $i = 1 \dots N$.

$$\Psi_t^{(m) \cap (n)}(i) = \left[\lambda_1^i \psi_1^{(m) \cap (n)}(i), \lambda_2^i \psi_2^{(m) \cap (n)}(i), \dots, \lambda_{L_{(m) \cap (n)}}^i \psi_{L_{(m) \cap (n)}}^{(m) \cap (n)}(i) \right]$$

where $t > 0$ is a tuneable parameter.

3. For each time sample $i = 1 \dots N$, concatenate the low-dimensional representations into a single vector

$$\begin{aligned} \Psi_t^{(u)}(i) = & \left(\Psi_t^{(1) \cap (2)}(i), \Psi_t^{(1) \cap (3)}(i), \dots, \Psi_t^{(1) \cap (M)}(i), \right. \\ & \Psi_t^{(2) \cap (1)}(i), \Psi_t^{(2) \cap (3)}(i), \dots, \Psi_t^{(2) \cap (M)}(i), \\ & \dots \\ & \left. \Psi_t^{(M) \cap (1)}(i), \Psi_t^{(M) \cap (2)}(i), \dots, \Psi_t^{(M) \cap (M-1)}(i) \right) \end{aligned}$$

4. Calculate $d_t^{(u)}$ using the Mahalanobis distance:

$$d_t^{(u)}(i, j) = \|\Psi_t^{(u)}(i) - \Psi_t^{(u)}(j)\|_{\text{Mahalanobis}}$$

for $i, j = 1, \dots, N$.

Algorithm 3. Mahalanobis-based Union Scheme.

Input: Aligned samples from M sets of observations: $\{(\mathbf{s}_i^{(1)}, \mathbf{s}_i^{(2)}, \dots, \mathbf{s}_i^{(M)})\}_{i=1}^N$.

Output: Low-dimensional representation of the common hidden random variables Θ .

1. For each pair of observation sets $1 \leq m, n \leq M, m \neq n$, apply alternating diffusion (Algorithm 2), and obtain the diffusion distance $d_t^{(m) \cap (n)}$.
2. Compute the distance $d_t^{(u)}$

$$d_t^{(u)}(i, j) = \sum_{1 \leq m, n \leq M, m \neq n} d_t^{(m) \cap (n)}(i, j)$$

for $i, j = 1, \dots, N$.

3. Based on the common diffusion distance $d_t^{(u)}$ calculate an affinity matrix

$$W_{i,j}^{(u)} = \exp\left(-\frac{(d_t^{(u)}(i, j))^2}{\epsilon^{(u)}}\right)$$

4. Construct the diffusion operator $\mathbf{K}^{(u)}$:

$$Q_{i,i}^{(u)} = \left(\sum_{l=1}^N W_{i,l}^{(u)}\right)^{-1}; \mathbf{K}^{(u)} = \mathbf{Q}^{(u)} \mathbf{W}^{(u)}$$

5. Apply standard diffusion maps (steps 3 and 4 in Algorithm 1) using $\mathbf{K}^{(u)}$, and obtain an L -dimensional representation of Θ .
-

Algorithm 4. Common Graph.

4. Simulation results

Consider the toy problem described in Section 2.1. We simulate 6 hidden scalar variables: 3 common variables ($\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$) and 3 nuisance variables ($n^{(1)}, n^{(2)}, n^{(3)}$). The variables are statistically independent and uniformly distributed in $[0, 2\pi]$. We then build 3 sets of N RGB images: $\{r_i^{(1)}\}, \{r_i^{(2)}\}, \{r_i^{(3)}\}, i = 1, \dots, N$. The sensitivity table of this example is given by

$$S^T = \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \end{pmatrix}. \quad (32)$$

Each image contains 3 arrows, where each arrow is rotated according to a randomly generated angle: the angles of the arrows in $r_i^{(1)}$ are $(\theta_i^{(1)}, \theta_i^{(2)}, n_i^{(1)})$, the angles in $r_i^{(2)}$ are $(\theta_i^{(2)}, \theta_i^{(3)}, n_i^{(2)})$, and the angles in $r_i^{(3)}$ are $(\theta_i^{(3)}, \theta_i^{(1)}, n_i^{(3)})$. The dimensionality of each RGB image is $36 \times 96 \times 3$. We column-stack the RGB images, i.e., $r_i^{(1)}, r_i^{(2)}, r_i^{(3)}$ are vectors of length $J = 10368$. The proposed algorithm is data-driven, and therefore, it does not assume any prior knowledge on the nature of observations. In order to highlight this important property, we use RPs. First, RPs with sufficiently large dimension maintain the underlying geometry, yet the image appearances are lost, which shows that our algorithm does not apply any image processing. Second, in the original images, the different hidden variables are manifested in separate coordinates/pixels; RPs mix the hidden variables, enabling a more challenging extraction task. We generate $D = 1600$ orthonormal vectors $\{b_i\}_{i=1}^D$ of length J and denote by $B \in \mathbb{R}^{J \times D}$ the matrix whose columns are these random vectors. We build the data of the sensors (cameras) by RPs $s_i^{(m)} = B^T r_i^{(m)}$, where m is the camera index. In the case of data acquired by cameras, B can be viewed as the coding system in the cameras. An illustration of the images and their RPs is depicted in Fig. 2. Illustration of the “movies” of the RPs captured by each camera can be seen in the following link <https://youtu.be/91N6mhlYQYY>.

We first apply DM separately to each set of observations. Fig. 3 presents 2-dimensional views of the obtained 3-dimensional embeddings. Each subfigure presents a scatter plot of embedded data-points. Each data-point is an image (a frame in the movie) captured by a certain camera after a random-projection $s_i^{(m)}$, where i is the frame index and m is the camera index. The axes of the scatter plot are the first 3 components of the obtained embedding derived from the corresponding camera. The embedded data-points are colored according to the rotating angles $(\theta^{(1)}, \theta^{(2)}, \theta^{(3)})$ and 3 noise variables $(n^{(1)}, n^{(2)}, n^{(3)})$. It should be noted that this information (the color) was added after calculating the embedding and was not taken into account in the computation of embedding. The subfigure in the l th column and in the m th row contains the embedded data-points derived from the m th camera $\{s_i^{(m)}\}_{i=1}^N$, and its data-points are colored according to the rotating angle of the l th arrow. In the 3 left columns the color coding is according to $\{\theta_i^{(1)}\}_{i=1}^N, \{\theta_i^{(2)}\}_{i=1}^N, \{\theta_i^{(3)}\}_{i=1}^N$, and in the 3 right columns the color coding is according to $\{n_i^{(1)}\}_{i=1}^N, \{n_i^{(2)}\}_{i=1}^N, \{n_i^{(3)}\}_{i=1}^N$. In other words, in each row the same scatter plot is shown, but with different color coding. The 3-dimensional scatter plots are rotated so that the obtained color gradient is best visualized from our 2-dimensional view point. For example, the subfigures in the second row are derived from the observations from the second camera $\{s_i^{(2)}\}_{i=1}^N$. The data-points in the first column are colored according to the rotation angles $\{\theta_i^{(1)}\}_{i=1}^N$, in the second column according to $\{\theta_i^{(2)}\}_{i=1}^N$, etc.

As can be seen, in each row, 3 scatter plots exhibit a smooth color

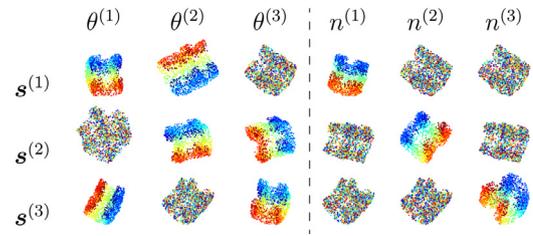
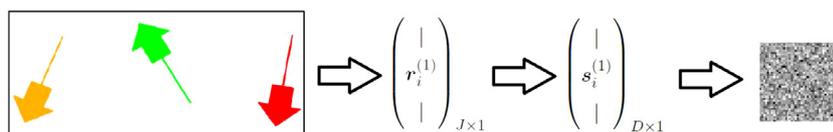


Fig. 3. 3D embedding obtained by applying diffusion map on a single observer. The subfigures are arranged such that subfigures in each row are obtained from the same observer. The data-points in each column are colored according to different arrow’s rotation angles.

gradient, 2 from the left 3 columns and 1 from the right 3 columns, corresponding to the variables sensed by the respective camera. In the 3 left columns, we see that the color gradients indicates accurate detection of the common variables according to the sensing matrix S . On the 3 right columns, only in the diagonal subfigures exhibit a smooth color gradient, indicating that each captures only its own nuisance variable, as expected. In conclusion, Fig. 3 implies that the obtained embeddings by DM provide accurate parametrizations of the hidden variables measured by each observation (camera), both the common and the nuisance variables.

The proposed algorithm is applied to the three sets of observations. The obtained embedding is depicted in Fig. 4. The same 3 dimensional scatter plot of the obtained embedding is shown with different color coding. The subfigures in the top row are colored (from left to right) according to the common variables $\{\theta_i^{(1)}\}_{i=1}^N, \{\theta_i^{(2)}\}_{i=1}^N, \{\theta_i^{(3)}\}_{i=1}^N$, while the subfigures in the bottom row are colored (from left to right) according to the nuisance variables $\{n_i^{(1)}\}_{i=1}^N, \{n_i^{(2)}\}_{i=1}^N, \{n_i^{(3)}\}_{i=1}^N$. As in Fig. 2, the 3 dimensional embedding is rotated, such that the corresponding color gradient is emphasized from the depicted 2 dimensional point of view. We can see from the obtained color gradients that the embedding provides a parametrization of only the common variables, meaning that the proposed algorithm manages to extract all 3 of the common variables $(\theta^{(1)}, \theta^{(2)}, \theta^{(3)})$ (despite having none in common to all three observations), while suppressing all 3 nuisance observation-specific variables $(n^{(1)}, n^{(2)}, n^{(3)})$. Upon publication, the Matlab code and data of this toy problem will be made available online.

5. Application to sleep stage assessment

As mentioned above, the problem of extracting the common hidden variables from multiple data sets taken by different observables can be perceived as a problem of nonlinear filtering. To demonstrate the potential of this particular nonlinear filtering scheme in processing real data, we apply the proposed algorithm to sleep data, where the ultimate goal is to devise an automatic system for sleep stage assessment.

Sleep is a global and recurrent physiological process, which is in charge of the memory consolidation, the learning redistribution, tissue regeneration, immune system enhancement, etc [30]. The sleep dynamics are characterized by particular temporal physiological features, which are intimately related to the quality of sleep. The clinically acceptable sleep stage is mainly determined by reading recorded electroencephalogram (EEG) signals based on the Rechtschaffen and Kales (R&K) criteria [31,32]. In the R&K criteria, the sleep dynamics are divided into two broad stages: rapid eye movement (REM), and non-rapid eye movement (NREM) [30]. The NREM stage is further divided into two shallow stages, which are denoted N1 and N2, and a deep sleep

Fig. 2. Random projection diagram of the i th image. Each RGB image was column stacked into a vector of length $J = 10368$. Then it was projected on a subspace of \mathbb{R}^D using an orthonormal set $\{v_i\}_{i=1}^D$. The projection is illustrated by a gray-scale 40×40 image. As can be seen the image’s property are lost through this projection.

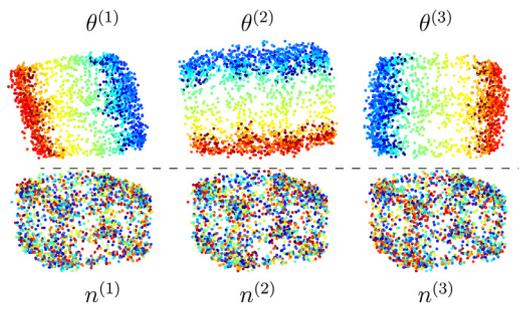


Fig. 4. 3D embedding obtained by applying the proposed algorithm on the observers set. The subplots in the first row are colored according to the common variables, the subplots in the second row are colored according to the noise variables. As can be seen, the obtained parametrization corresponds to the common variables.

stage, which is denoted N3. In addition to the interest stemming from physiological aspects, sleep stage assessment has important clinical applications. For example, REM is associated with perceptual skill improvement [33], slow wave sleep is associated with Alzheimer’s disease [34], poor sleep quality is associated with weaning failure [35], etc. Besides personal health purposes, the sleep quality is also responsible for several public catastrophes [36]. These facts indicate the importance of an accurate automatic annotation system for sleep stage assessment and its broad applications.

In the past decades, various automatic annotation methods have been proposed. Those methods mainly extract various features from the EEG recordings for the purpose of studying sleep dynamics [37], such as time domain summary statistics, spectral or coherence features, time-frequency features, and information entropy, just to name a few [38–40]. Recently, a theoretically solid approach suitable for analyzing and estimating the dynamics of the brain activity from recorded EEG signals has been proposed in [22,23]. A particular aspect of sleep

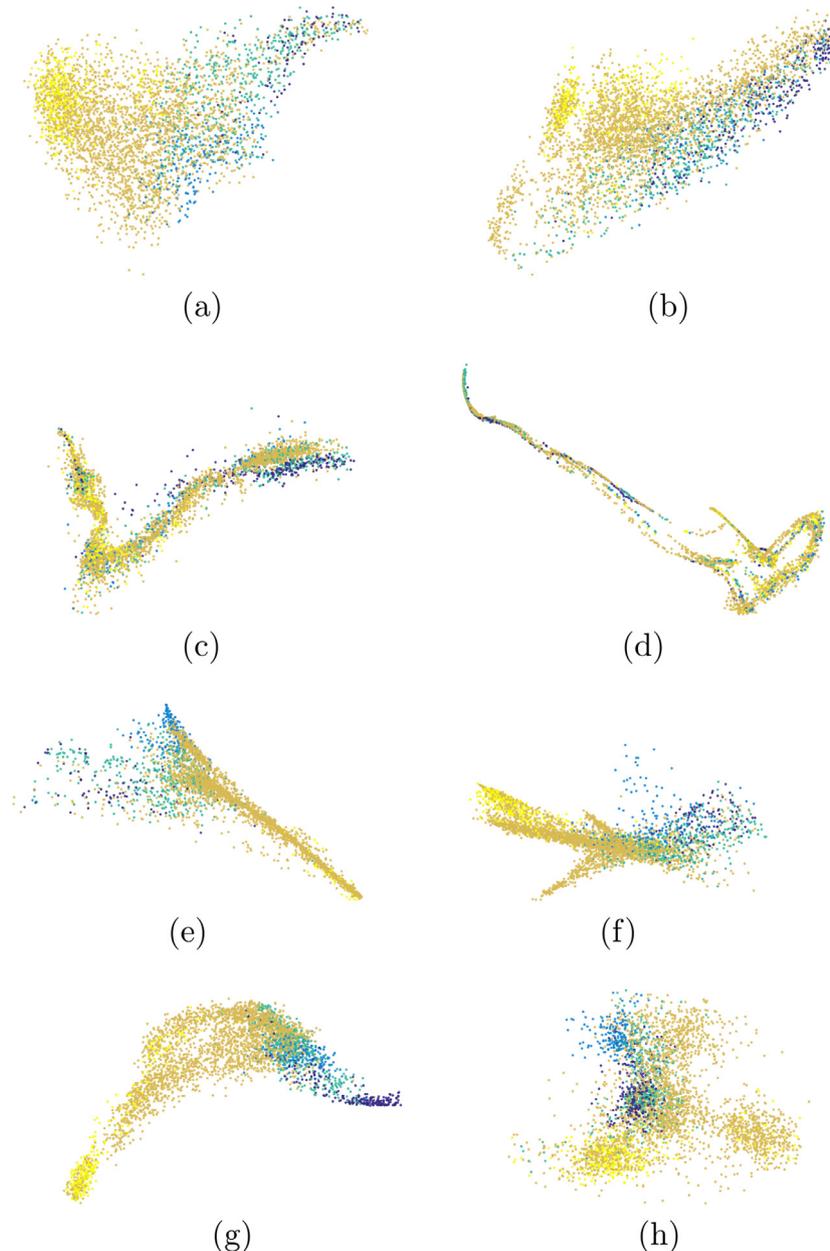


Fig. 5. The 3D RPs of the embeddings obtained by single-sensor DM (top row), the concatenation scheme (second row), the multiplication scheme (third row), and Algorithm 4 (bottom row). The points are colored according to the sleep stage. The embeddings are based on the O2A1 channel in (a) and on the airflow measurements in (b). From the second row to the bottom row, the embeddings on the left column are based on the EEG set, and on the right column are based on the respiratory set.

Table 2

Classification results using SVM. The prediction errors (standard deviations) based on the different embeddings are presented. The total error (standard deviation) is calculated by a weighted mean of the prediction errors (standard deviations) in each sleep stage. (a) The classification based on the respiratory set. (b) The classification based on the EEG set.

(a)						
Prediction errors based on the respiratory set						
Sensor	Scheme					
	Awake	REM	N1	N2	N3	Total
CFlow	0.383	0.258	0.545	0.179	0.244	0.264 (0.073)
ABD	0.299	0.154	0.426	0.162	0.185	0.209 (0.051)
THO	0.262	0.15	0.412	0.153	0.164	0.196 (0.051)
Noise	0.559	0.513	0.583	0.582	0.529	0.562 (0.039)
Concatenation scheme	0.476	0.474	0.552	0.516	0.5	0.507 (0.039)
Multiplication scheme	0.343	0.262	0.555	0.265	0.265	0.307 (0.069)
Common Graph	0.252	0.183	0.443	0.178	0.201	0.22 (0.048)

(b)						
Prediction errors based on the EEG set						
Sensor	Scheme					
	Awake	REM	N1	N2	N3	Total
O1A2	0.267	0.281	0.624	0.132	0.273	0.25 (0.056)
O2A1	0.283	0.25	0.603	0.142	0.24	0.244 (0.069)
C4A1	0.305	0.273	0.623	0.139	0.291	0.258 (0.068)
C3A2	0.298	0.276	0.619	0.132	0.289	0.254 (0.06)
Noise	0.551	0.499	0.579	0.58	0.53	0.557 (0.042)
Concatenation scheme	0.342	0.325	0.499	0.307	0.325	0.34 (0.039)
Multiplication scheme	0.219	0.191	0.47	0.238	0.2	0.252 (0.045)
Common Graph	0.227	0.153	0.425	0.133	0.179	0.188 (0.036)

dynamics, which has not gained much research attention in the line of research mentioned above, is that sleep is not localized solely in the brain and is reflected in other physiological systems as well. For example, the regulation of mechanoreceptor and the chemoreceptor leads to breathing pattern variability in the respiratory signal. We have a remarkably regular breathing during N3 stage and irregular breathing with fast varying instantaneous frequency and amplitude during REM stage. Those physiological phenomena motivated various studies to explore the relation between the sleep stage and the patterns in the respiratory signals, e.g. [41–43]. Physiologically, these variations are not originated from the same controller, and phenomenologically do not have the same patterns in the recorded time series. Thus, while we could observe the sleep dynamics via observing the characteristics of different sensors, each of them reflects only part of the sleep dynamic, and is complicated by the nature of the sensor.

Based on the above physiological facts, an automatic approach for assessing the sleep stage was presented in [28]. It relies on the assumption that there exist hidden low-dimensional physiological processes driving the sleep dynamics, and hence the accessible measured signals. However, these hidden processes may be deformed by the observation procedures; each observation (e.g., an EEG channel measuring brain activity or a chest belt measuring respiration) can be influenced by nuisance factors, which are sensor- or channel-specific (e.g., the specific type of sensors and their exact positions), yet our interest is in the intrinsic variables related to the sleep stages. In [28], empirical intrinsic geometry (EIG) method [22,23], which is based on nonlinear independent component analysis [21] and was proven to be invariant to the measurement modality, was applied to build an intrinsic representation of the measured data. In [44], this method was extended to a pair of sensors. It was shown that by analyzing the measurements taken simultaneously from 2 sensors, a more reliable intrinsic representation of the sleep dynamics can be obtained,

compared with the analysis based only on a single signal.

In this section we extend the algorithm shown in [44], and process jointly multiple channels. We show that extracting the underlying common variables from multiple data sets acquired in different channels recovers systematically a representation, which is well correlated with the sleep stage. The analogy to the setting described in this paper is as follows. We assume that the sleep dynamics are intimately related to hidden controllers that affect the respiratory as well as the brain neural system. These controllers are not accessible to us; yet, they can be recovered by analyzing observations from multiple channels/sensors, each captures different, partial yet complementary aspects of it. Under this assumption, our interest is in obtaining the intrinsic variables underlying the measurements related to these controllers. On the one hand, by analyzing multiple observation channels we can gather more information on the hidden controllers. On the other hand, observations from each channel might be deformed by the different acquisition and measurement modalities and may be affected by noise and interferences, specific to the particular (type of) sensor. In the context of this work, this tradeoff is addressed by defining the intrinsic variables (related to the hidden controllers of interest) as those which are not sensor-specific, and hence, the variables of interest are those that are common among at least two observables.

Twenty subjects without sleep apnea were chosen for this study. The demographic characteristics of these individuals fall within the normal ranges. We used recordings of 6 hours per subject, which were performed in the sleep center at Chang Gung Memorial Hospital (CGMH), Linkou, Taoyuan, Taiwan. The institutional review board of the CGMH approved the study protocol (No. 101-4968A3) and the enrolled subjects provided written informed consent. See [28] for more details regarding the experimental setting and the collected data.

We build the common graph according to Algorithm 4 for extracting the common hidden variables separately to two sets of sensors. The first set includes 3 signals: abdominal and chest motions, which are recorded by piezo-electric bands, and airflow, which is measured using thermistors and nasal pressure, all 3 at sampling rate of 100 Hz. The second set comprises recordings from 4 EEG channels: C3A2, C4A1, O1A2 and O2A1 at sampling rate of 200 Hz. The recorded respiratory signals are denoted by R_m , $m = 1, \dots, 3$ and the EEG signals are denoted by E_m , $m = 1, \dots, 4$.

Prior to the application of our method, each of the single-channel recordings was preprocessed by applying the scattering transform as in [28], which was shown to improve the regularity and stability of signals with respect to various deformations [45]. We then apply Algorithm 4 separately twice: once to the respiratory set, and once to the EEG set.

In order to demonstrate the inherent “sensor selection” capability of the proposed method, for each set of measurements we added an artificial “pure noise” sensor to simulate possible sensor failure. Because our processing pipeline begins with the application of the scattering transform, which automatically degenerates any stationary noise, the noise sensor consists of a non-stationary sequence generated by modulating a sine-wave according to

$$\begin{aligned} \phi(\tau) &= \frac{1}{2} + \frac{1}{4} \sin\left(\frac{2\pi\tau}{512 \cdot 10}\right) \\ n(t) &= \sin\left(2\pi \int_0^t \phi(\tau) d\tau\right) \end{aligned} \quad (33)$$

where $n(t)$ is the continuous time signal and $\phi(\tau)$ can be viewed as the instantaneous frequency of $n(t)$. We sample the obtained modulated sine-wave $n(t)$ at a sampling rates of 200 Hz and 100 Hz for the EEG set and for the respiratory set, respectively. It should be noted that this particular non-stationary “noise” implementation was chosen just for the sake of demonstration, and any other non-stationary sequence could be chosen instead.

We compare the results of the common graph algorithm, analyzing multiple sensors, with the results attained by the standard DM applied

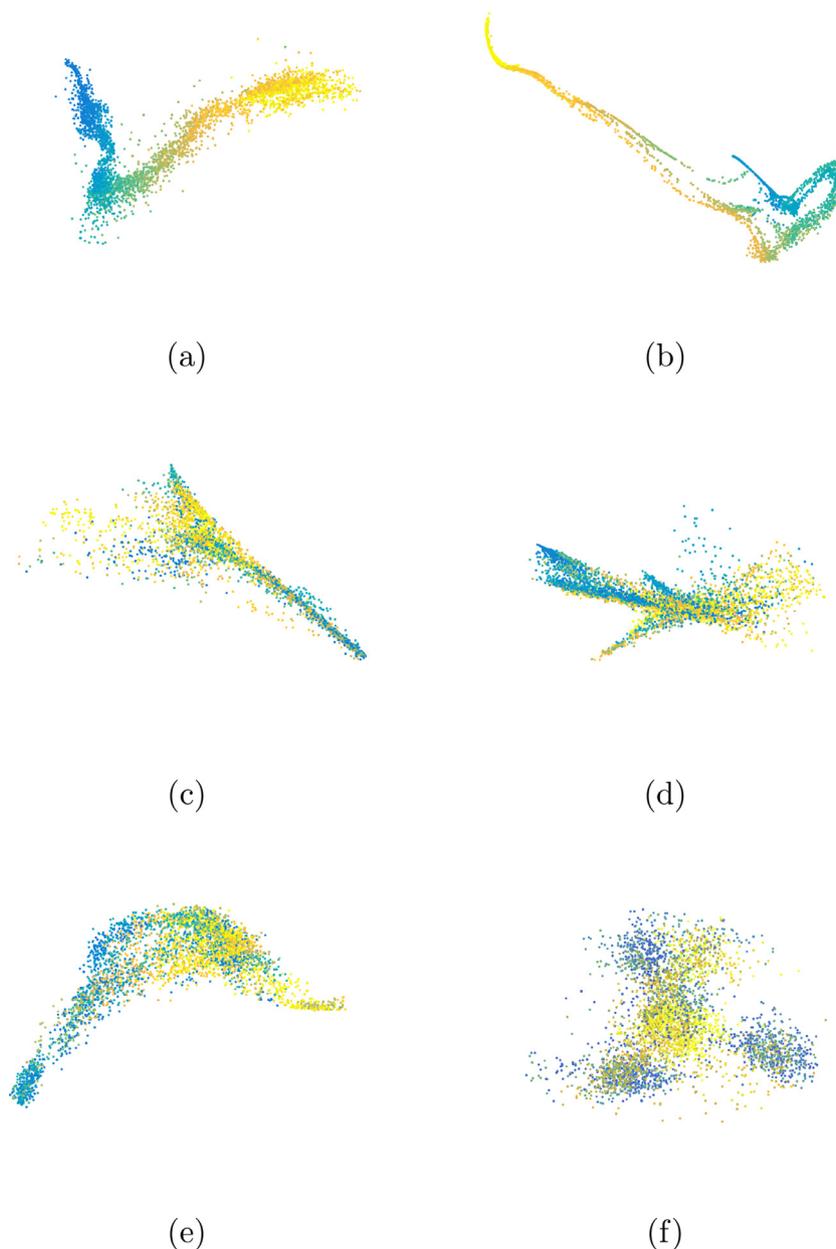


Fig. 6. The same embeddings as in Fig. 5, colored according to the instantaneous frequency of the noise sensor.

separately to each individual sensor. In addition, we compare the results to two competing schemes analyzing multiple sensors. In the first scheme, we concatenate the scattering transform components from each sensor, and then, apply the standard DM. We note that conceptually this scheme takes into account the information captured by all the sensors without any filtering. We refer to the first scheme as the *concatenation scheme*. In the second scheme, we apply AD to the entire set of sensors. Namely, we calculate the diffusion kernel $\mathbf{K}^{(m)}$ for the m th sensor, where $m = 1, \dots, M$ and build an AD kernel based on the product of all the kernels, that is, $\mathbf{K} = \mathbf{K}^{(1)}\mathbf{K}^{(2)}\dots\mathbf{K}^{(M)}$. Then, we apply DM with this AD kernel. This scheme takes into account only the information that is captured simultaneously by all of the sensors, namely $\bigcap_{m \neq n} (\mathcal{S}^{(m)} \cap \mathcal{S}^{(n)})$, thereby performing excessive filtering. We refer to the second scheme as the *multiplication scheme*.

The calculation of the affinity matrices, which is a core element in the tested methods, is carried out using the Mahalanobis distance variant presented in [21], which was discussed in Section 3.3. To be able to depict information embodied in more than three eigenvectors, we randomly project the embeddings attained by the competing algorithms

to 3 dimensions. This allows us to visually inspect the portion of the relevant information and the portion of the nuisance information manifested in the representations obtained by the different algorithms. We use the same projection in all tested methods.

The RPs of the embeddings are depicted in Fig. 5. The RPs based on the single channel DM applied to the O2A1 EEG channel and to the airflow channel are depicted in the top row. The RPs based on the concatenation scheme, multiplication scheme and the proposed algorithm are depicted in the second, third and bottom rows, respectively. The embeddings depicted on the left column are based on the EEG set, and on the right column are based on the respiratory set. Each embedded point is colored according to its respective sleep stage, as identified by a human expert. Importantly, the information on the sleep stage (e.g., the color) was not taken into account in the algorithms forming the embeddings.

Fig. 5 provides a visual illustration of the obtained parametrization with respect to the sleep stage. By comparing the Fig. 5a and b with Fig. 5g and h we can observe the improvement achieved by the additional information obtained from combining information from multiple

sensors. In addition, by comparing Fig. 5c–f with Fig. 5g and h we observe the improvement achieved by filtering out of the sensor-specific nuisance variables. In these comparisons, it can be seen that the embeddings obtained by using the proposed algorithm results in a better parametrization of the sleep stage evaluation; different sleep states appear to be more separated, especially in the case of the respiratory signals.

To objectively assess the quality of the obtained embeddings, we use multi-class support vector machine (SVM). To ensure convergence and to prevent overfitting, we process only the 15 most dominant eigenvectors from each embedding. It should be noted that due to the obtained fast decay of the eigenvalues, taking only the 15 most dominant eigenvectors preserves the geometrical structure of the data. We randomly partition the data into 2 sets – a training set (consisting of 75% of the samples) and a validation set (consisting of 25% of the samples). The validation set contains 1,250 time segments, which consist (on average) of 13.2%(165) segments labeled as awake stage, 10%(125) segments labeled as REM stage, 11.2%(140) segments labeled as N1 stage, 49.6%(620) segments labeled as N2 stage and 16%(200) segments labeled as N3 stage. The trained classifier is used to classify the sleep stage in the validation set. We repeat this classification 10 times, for different randomly chosen partitions of training and validation sets. The average classification results for each scheme are depicted in Table 2. The obtained classification results achieved by the proposed algorithm are superior compared to the obtained results from other schemes, both in the case of the EEG set and in the case of the respiratory set. In these results, the advantages of proper filtering are evident, as it can be seen that in contrast to the proposed algorithm, the concatenation scheme and the multiplication scheme attain inferior classification results, and in some cases, their results are comparable to the results achieved by processing data from only a single sensor. In the case of the multiplication scheme this may be caused by too excessively filtering. In the case of the concatenation scheme, where no filtering is applied, this may be caused by the existence of interferences and noise.

The results in Table 2 may provide additional insights related to the sleep dynamics that extend the scope of the evaluation of the algorithms. The classification results achieved by the multiplication scheme are inferior comparing to the results achieved by single-sensor schemes in the case of the respiratory set, where as in the case of the EEG set the achieved results are similar to the single-sensor schemes. This supports the hypothesis that different EEG recording exhibit more homogeneous geometrical structures, with possibly less noise and fewer distortions, compared to the data acquired through the different respiratory recordings.

The homogeneity of the EEG set might explain another interesting observation stemming from these classification results. In the case of the EEG set, we can see that combining the information acquired from multiple sensors using the proposed algorithm results with superior results, even compared to the results that would have been achieved using the best single-sensor scheme. This implies that the proposed algorithm manages to simultaneously cancel the effect of the additional noise-sensor as well as to properly integrate the information embodied in the multiple sensors. Conversely, in the case of the respiratory set, we can see that even though the proposed algorithm manages to improve the results achieved by the CFlow channel, it did not manage to improve the results achieved by the single sensor schemes based on the ABD or the THO channel. Yet, it did manage to cancel the effect of the noise-sensor, but not as successfully as in the case of the EEG set. In this regard, it is worth emphasizing that the evaluation of the results from each sensor and from each scheme are based on unknown sleep stage labelling. Thus, the “quality” of the different sensors are not known in advance, and obtaining a result from our sensor fusion scheme that is comparable to the results attained by the best single sensor is still of value.

Fig. 6 further illustrates the poor embeddings and classification results achieved by the concatenation scheme. The same embeddings,

which are depicted in Fig. 5, are presented here, but this time with a different color – now according to the instantaneous frequency of the noise sensor (33). As can be observed, in contrast to the embeddings achieved by the proposed algorithm or by the multiplication scheme, the embeddings achieved by the concatenation scheme are well correlated with the instantaneous frequency in the noise sensor, indicating that the underlying structure is wrongly captured. This further illustrates the difference between the filtering effects of our algorithm and other methods, which are based to the fusion of data from all the sensors [9,10,12].

6. Conclusions

In this paper, we propose a new algorithm for fusing information measured by multiple, multimodal sensors. The primary focus is on a setting in which all sensors observe the same system, but each introduces different variables – some are related to various aspects of the system of interest, whereas others are sensor-specific and irrelevant. We present a nonlinear data fusion scheme for suppressing the sensor-specific variables while preserving the system variables measured by two or more sensors. Experimental results demonstrate the applicability of our method to artificial toy problem and to recorded multimodal data for the purpose of sleep stage assessment.

The core of the presented technique is an implementation of an abstract notion of intersection and union of multimodal data sets. While the intersection between two sets is well defined and theoretically explained [14], the union of two (or more) sets still calls for rigorous analysis. Future work will include such analysis and the development of a union scheme that respects uniqueness.

Acknowledgments

The authors would like to thank the anonymous reviewers for their helpful suggestions. The authors wish to thank Ofer Karp and Idan Amir for sharing their insights from prior contributions on the subject and for making their code available. Hau-tieng Wu acknowledges the support of Sloan Research Fellow FR-2015-65363. This research was partly supported by the European Union Seventh Framework Programme (FP7) under Marie Curie Grant 630657 and by the Israel Science Foundation (grant no. 1490/16).

References

- [1] B. Khaleghi, A. Khamis, F.O. Karray, S.N. Razavi, Multisensor data fusion: a review of the state-of-the-art, *Inf. Fusion* 14 (1) (2013) 28–44.
- [2] D. Lahat, T. Adali, C. Jutten, Multimodal data fusion: an overview of methods, challenges, and prospects, *Proc. IEEE* 103 (9) (2015) 1449–1477.
- [3] R. Gravina, P. Alinia, H. Ghasemzadeh, G. Fortino, Multi-sensor fusion in body sensor networks: state-of-the-art and research challenges, *Inf. Fusion* 35 (2017) 68–80.
- [4] J.B. Tenenbaum, V. de Silva, J.C. Langford, A global geometric framework for nonlinear dimensionality reduction, *Science* 260 (2000) 2319–2323.
- [5] S.T. Roweis, L.K. Saul, Nonlinear dimensionality reduction by locally linear embedding, *Science* 260 (2000) 2323–2326.
- [6] D.L. Donoho, C. Grimes, Hessian eigenmaps: new locally linear embedding techniques for high-dimensional data, *Proc. Nat. Acad. Sci.* 100 (2003) 5591–5596.
- [7] M. Belkin, P. Niyogi, Laplacian eigenmaps for dimensionality reduction and data representation, *Neural. Comput.* 15 (6) (2003) 1373–1396.
- [8] R.R. Coifman, S. Lafon, Diffusion maps, *Appl. Comput. Harmon. Anal.* 21 (1) (2006) 5–30.
- [9] M. Davenport, C. Hegde, M.F. Duarte, R.G. Baraniuk, Joint manifolds for data fusion, *IEEE Trans. Image Process.* 19 (10) (2010) 2580–2594.
- [10] Y. Keller, R. Coifman, S. Lafon, S.W. Zucker, Audio-visual group recognition using diffusion maps, *IEEE Trans. Signal Process.* 58 (1) (2010) 403–413, <http://dx.doi.org/10.1109/TSP.2009.2030861>.
- [11] O. Yair, R. Talmon, Local canonical correlation analysis for nonlinear common variables discovery, *IEEE Trans. Signal Process.* 65 (5) (2017) 1101–1115.
- [12] M. Salhov, O. Lindenbaum, A. Silberschatz, Y. Shkolnisky, A. Averbuch, Multi-view kernel consensus for data analysis and signal processing, arXiv: 1606.08819 (2016).
- [13] R.R. Lederman, R. Talmon, Learning the geometry of common latent variables using alternating-diffusion, *Appl. Comp. Harmon. Anal.* (2015), <http://dx.doi.org/10.1016/j.acha.2015.09.002>.

- [14] R. Talmon, H.-t. Wu, Latent common manifold learning with alternating diffusion: analysis and applications, arXiv: 1602.00078 (2016).
- [15] P.L. Lai, C. Fyfe, Kernel and nonlinear canonical correlation analysis, *Int. J. Neural Syst.* 10 (05) (2000) 365–377.
- [16] V.R. de Sa, Spectral clustering with two views, *ICML Workshop on Learning with Multiple Views*, (2005).
- [17] V.R. de Sa, P.W. Gallagher, J.M. Lewis, V.L. Malave, Multi-view kernel construction, *Mach. Learn.* 79 (1–2) (2010) 47–71, <http://dx.doi.org/10.1007/s10994-009-5157-z>.
- [18] B. Boots, G.J. Gordon, Two-manifold problems with applications to nonlinear system identification, *Proc. 29th Intl. Conf. on Machine Learning (ICML)*, (2012).
- [19] T. Michaeli, W. Wang, K. Livescu, Nonparametric canonical correlation analysis, arXiv: 1511.04839 (2015).
- [20] C.J. Dsilva, R. Talmon, C.W. Gear, R.R. Coifman, I.G. Kevrekidis, Data-driven reduction for multiscale stochastic dynamical systems, arXiv: 1501.05195 (2015).
- [21] A. Singer, R.R. Coifman, Non-linear independent component analysis with diffusion maps, *Appl. Comput. Harmon. Anal.* 25 (2) (2008) 226–239.
- [22] R. Talmon, R. Coifman, Empirical intrinsic geometry for nonlinear modeling and time series filtering, *Proc. Nat. Acad. Sci.* 110 (31) (2013) 12535–12540.
- [23] R. Talmon, R.R. Coifman, Intrinsic modeling of stochastic dynamical systems using empirical geometry, Tech. Report, YALEU/DCS/TR1467, 2014.
- [24] R. Talmon, S. Mallat, H. Zaveri, R.R. Coifman, Manifold learning for latent variable inference in dynamical systems, *IEEE Trans. Signal Process.* 63 (15) (2015) 3843–3856.
- [25] N.E. Karoui, H.-T. Wu, Graph connection laplacian methods can be made robust to noise, *Ann. Stat.* 44 (1) (2016) 346–372.
- [26] A. Singer, H.-T. Wu, Spectral convergence of the connection laplacian from random samples, *Inf. Inference* 6 (1) (2017) 58–123.
- [27] A. Singer, Lecture notes in “massive data analysis”, 2015, (University Lecture).
- [28] H.-t. Wu, R. Talmon, Y.-L. Lo, Assess sleep stage by modern signal processing techniques, *Biomed. Eng., IEEE Trans.* 62 (4) (2015) 1159–1168.
- [29] E.J. Candes, T. Tao, Near-optimal signal recovery from random projections: universal encoding strategies, *IEEE Trans. Inf. Theory* 52 (12) (2006) 5406–5425.
- [30] T. Lee-Chiong, *Sleep Medicine: Essentials and Review*, Oxford, 2008.
- [31] A. Rechtschaffen, A. Kales, *A Manual of Standardized Terminology, Techniques and Scoring System for Sleep Stages of Human Subjects*, Washington: Public Health Service, US Government Printing Office, 1968.
- [32] R. Berry, R. Budhiraja, D. Gottlieb, et al., Rules for scoring respiratory events in sleep: update of the 2007 AASM manual for the scoring of sleep and associated events, *J. Clin. Sleep Med.* 8 (5) (2012) 597–619.
- [33] A. Karni, D. Tanne, B.S. Rubenstein, J.J. Askenasy, D. Sagi, Dependence on REM sleep of overnight improvement of a perceptual skill, *Science* 265 (5172) (1994) 679–682.
- [34] J.-e. Kang, M.M. Lim, R.J. Bateman, J.J. Lee, L.P. Smyth, J.R. Cirrito, N. Fujiki, S. Nishino, D.M. Holtzman, Amyloid- β dynamics are regulated by orexin and the sleep-wake cycle, *Science* 326 (Nov 13) (2009) 1005–1007.
- [35] F. Roche Campo, X. Drouot, A.W. Thille, F. Galia, B. Cabello, M.-P. D’Ortho, L. Brochard, Poor sleep quality is associated with late noninvasive ventilation failure in patients with acute hypercapnic respiratory failure. *Crit. Care Med.* 38 (2) (2010) 477–485.
- [36] H.R. Colten, B.M. Altevogt, Functional and economic impact of sleep loss and sleep-related disorders, in: H.R. Colten, B.M. Altevogt (Eds.), *Sleep Disorders and Sleep Deprivation: An Unmet Public Health Problem*, The National Academies Press, 2006.
- [37] V. Bajaj, R.B. Pachori, Automatic classification of sleep stages based on the time-frequency image of EEG signals, *Comput. Methods Programs Biomed.* 112 (3) (2013) 320–328.
- [38] N. Kannathal, M. Choo, U. Acharya, P. Sadasivan, Entropies for detection of epilepsy in EEG, *Comput. Methods Programs Biomed.* 80 (2005) 187–194.
- [39] S. Blanco, R. Quiroga, O. Rosso, S. Kochen, Time-frequency analysis of electroencephalogram series, *Phys. Rev. E* 51 (3) (1995) 2624–2631.
- [40] S. Geng, W. Zhou, Q. Yuan, D. Cai, Y. Zeng, EEG non-linear feature extraction using correlation dimension and hurst exponent, *Neurol. Res.* 33 (9) (2011) 908–912.
- [41] G.S. Chung, B.H. Choi, K.K. Kim, Y.G. Lim, J.W. Choi, D.-U. Jeong, K.-S. Park, REM sleep classification with respiration rates, *Information Technology Applications in Biomedicine, 2007. ITAB 2007. 6th International Special Topic Conference on*, (2007), pp. 194–197.
- [42] G. Guerrero-Mora, P. Elvia, A. Bianchi, J. Kortelainen, M. Tenhunen, S. Himanen, M. Mendez, E. Arce-Santana, O. Gutierrez-Navarro, Sleep-wake detection based on respiratory signal acquired through a pressure bed sensor, *Engineering in Medicine and Biology Society (EMBC), 2012 Annual International Conference of the IEEE*, (2012), pp. 3452–3455.
- [43] J. Sloboda, M. Das, A simple sleep stage identification technique for incorporation in inexpensive electronic sleep screening devices, *Aerospace and Electronics Conference (NAECON), Proceedings of the 2011 IEEE National*, (2011), pp. 21–24.
- [44] R.R. Lederman, R. Talmon, H.-t. Wu, Y.-L. Lo, R.R. Coifman, Alternating diffusion for common manifold learning with application to sleep stage assessment, *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE*, 2015, pp. 5758–5762.
- [45] S. Mallat, Group invariant scattering, *Pure Appl. Math.* 10 (65) (2012) 1331–1398.