

## Manifold alignment aware ants

Mohammadi, Mohammad; Tino, Peter; Bunte, Kerstin

DOI:

[10.1162/neco\\_a\\_01478](https://doi.org/10.1162/neco_a_01478)

License:

None: All rights reserved

*Document Version*

Peer reviewed version

*Citation for published version (Harvard):*

Mohammadi, M, Tino, P & Bunte, K 2022, 'Manifold alignment aware ants: a Markovian process for manifold extraction', *Neural Computation*, vol. 34, no. 3, pp. 595-641. [https://doi.org/10.1162/neco\\_a\\_01478](https://doi.org/10.1162/neco_a_01478)

[Link to publication on Research at Birmingham portal](#)

### **Publisher Rights Statement:**

This is an accepted manuscript version of an article first published in *Neural Computation* <https://direct.mit.edu/neco>. The final version of record is available at [https://doi.org/10.1162/neco\\_a\\_01478](https://doi.org/10.1162/neco_a_01478)

### **General rights**

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

### **Take down policy**

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact [UBIRA@lists.bham.ac.uk](mailto:UBIRA@lists.bham.ac.uk) providing details and we will remove access to the work immediately and investigate.

# Manifold Alignment Aware Ants: a Markovian process for manifold extraction

**Mohammad Mohammadi<sup>1</sup>, Peter Tino<sup>2</sup>, Kerstin Bunte<sup>1</sup>**

<sup>1</sup>Faculty of Science and Engineering, University of Groningen, the Netherlands.

<sup>2</sup>Department of Computer Science, University of Birmingham.

**Keywords:** manifold learning, swarm intelligence, dimensionality reduction, noise reduction, density estimation

## Abstract

The presence of manifolds is a common assumption in many applications including astronomy and computer vision. For instance, in astronomy low-dimensional stellar structures, such as streams, shells and globular clusters, can be found in the neighborhood of big galaxies such as the Milky Way. Since these structures are often buried in very large data sets, an algorithm which can not only recover the manifold, but also remove the background noise (or outliers) is highly desirable. While other works try to recover manifolds either by pushing all points towards manifolds or by down-sampling from dense regions, aiming to solve one of the problems, they generally fail to suppress the noise on manifolds and remove background noise simultaneously. Inspired by the collective behavior of biological ants in food-seeking process, we propose a new algorithm which employs several random walkers who are equipped with a local align-

ment measure to detect and denoise manifolds. During the walking process the agents release pheromone on data points which reinforces future movements. Over time the pheromone concentrates on the manifolds, while it fades in the background noise due to an evaporation procedure. We use the Markov Chain (MC) framework to provide a theoretical analysis of the convergence of the algorithm and its performance. Moreover, an empirical analysis, based on synthetic and real world data sets, is provided to demonstrate its applicability in different areas, such as: a) improving the performance of t-distributed stochastic neighbor embedding (t-SNE) and spectral clustering using the underlying MC formulas, b) recovering astronomical low-dimensional structures and c) improving the performance of the Fast Parzen Window density estimator.

## **1 Introduction**

New technological developments facilitate the collection of large amounts of high-dimensional data in different fields, such as astronomy, sensor networks, medical science, and computer vision. A typical challenge in dealing with high-dimensional data is the Curse of Dimensionality, where the data space is sparse, such that data points are far from their neighbors. However, in practice, the high dimensional data is often generated by a system governed by a small number of underlying components (Dixit, 2019). In other words, the high-dimensional data lies on a lower-dimensional topological structure called a manifold. Therefore, many dimensionality reduction methods, such as (Roweis and Saul, 2000; Belkin and Niyogi, 2003; Donoho and Grimes, 2003; Zhang and Zha, 2003; Coifman and Lafon, 2006), aim to identify this underlying manifold

from the data. However, such approaches might be tremendously impeded in their performance because of the presence of noise and outliers (Wang and Carreira-Perpinán, 2010). This led to the development of methods aiming to deal with noisy manifolds and/or outliers.

There are different strategies to decrease the impact of noise on manifold learning algorithms. The technique proposed in (Little et al., 2020) denoises data sets via down-sampling and picking samples in high dense regions, whereas others suppress the noise level via pushing all instances towards region with higher density. While the first category can deal with background noise (or outliers not belonging to any manifold), the basic assumption in the second category is that the noisy samples belong to a manifold, and their deviations from the manifold are caused by measurement noise. An example technique is Manifold Denoising (MD) (Hein and Maier, 2006), where the denoising process is modeled as a diffusion process on a neighborhood graph using the Laplacian to denoise the manifold. The aim of Klicpera et al. (2019) is to generalise the notion of diffusion operator on graphs so that larger node neighborhoods than one-hop in the neighborhood graph are considered. To this end, the authors first construct a transition graph (that can be re-normalized to define a MC) that encompasses a variety of (potentially unbounded) node neighborhoods, albeit strongly down-weighting increasingly large ones. The resulting dense transition graph then needs to be sparsified to consider only (potentially) higher order neighborhoods that “really matter”, i.e. have strong support in the transition structure. The Generalized Diffusion Convolution (GDC) is then formulated in such a new neighborhood structure. (Wang and Carreira-Perpinán, 2010) proposes Manifold Blurring Mean Shift (MBMS), where the mean shift directions are

only allowed to be parallel to manifold normals. Although increasing the number of iterations in the above algorithms may cause a better result, they eventually partition the manifold into local clusters. To overcome this problem, Locally Linear Denoising (LLD) (Gong et al., 2010) proposes a non-iterative method, assuming that a manifold can be explained as a set of overlapping linear patches. While it denoises each patch separately, it simultaneously uses the graph Laplacian to achieve a smooth manifold. However, such denoising methods may fail in applications, such as astronomy, where noisy manifolds are buried inside point clouds. Thus, a method capable of handling both types of noise, along with a manifold and background, is highly desirable.

Natural systems, such as ant colonies, have motivated many swarm algorithms in computer science. For instance, the cooperation among ants in the food-seeking process has inspired many solutions for optimization problems (Dorigo et al., 1991; Dorigo, 1992; Dorigo et al., 1996; Stützle and Hoos, 2000; Maniezzo, 1999; Blum et al., 2001) and clustering (Tsai et al., 2004; Chu et al., 2004; Runkler, 2005). In nature, when an ant leaves its nest to find food it deposits a chemical substance called pheromone on its path, which serves as information for other ants which are attracted to it. Each ants' decision is influenced by the pheromone: the higher the concentration is on a particular route, the more likely it is chosen as a path to follow. This seemingly uncoordinated local behavior in collection assists the hive to find the shortest path to the food source. Following this strategy, several algorithms were proposed and applied to optimization problems. Their simplicity and flexibility make these methods desirable for many optimization tasks that include graphs, such as vehicle/internet routing (Rizzoli et al., 2004) and water distribution systems (Gil et al., 2011).

Unfortunately, strategies explained above fail to denoise data comprising of noisy manifolds contained in point clouds, which are often faced in application domains involving particle simulations, such as astronomy. In this contribution, we propose a new algorithm to uncover noisy manifolds buried in high-dimensional background noise. Motivated by ant colony optimization strategies, our method employs multiple artificial ants who jump from point to point based on defined preferences and release pheromone. In previous work we proposed a strategy in which ants were allowed to transport data points towards the manifolds and hence denoising and “cleaning” them (Mohammadi and Bunte, 2020). In contrast, this contribution aims to detect noisy manifolds in potentially large amounts of background noise. Here, the position of data points is not changed, instead the ants prefer to move towards low-dimensional structures and deposit pheromone on the points they visit. The latter serves as positive feedback that accumulates and concentrates on the manifolds, leaving the noisy points with relatively less pheromone and hence highlighting the structures. Both our approach and GDC are built on a notion of a stochastic transition matrix reflecting the connection structure of a given graph. While in our case, the neighborhood graph represents local geometric structures of a point cloud, in GDC pairwise relationships are used (even though in applications the graphs do often represent point clouds). In general, our approach differs from graph diffusion (GD) in several important aspects:

1. In our approach the adjacency structure of the neighborhood graph by itself is not sufficient - we explicitly enforce a preference for neighbors that are aligned with low-dimensional structures in the point cloud the neighborhood graph represents;
2. Motivated by ant colony algorithms we allow for an additional mechanism of

positive reinforcement when discovering low-dimensional point structures in a noisy background, namely pheromone accumulation and evaporation;

3. We do not perform any graph convolution, instead we associate the “importance values” of each node with the amount of accumulated pheromone. Note that these values potentially reflect large scale geometric structures as the pheromone is deposited over long-range ant walks. In this sense, our approach shares the ambition of GDC for representing larger scale “important” neighborhoods. However, in GDC the larger scale structures are rapidly weighted down, whereas in our case, the pheromone deposits can survive over large scale structures, if frequently visited by the ants.

With empirical experiments we demonstrate the capability of our algorithm to highlight manifolds in synthetic and real-world applications. Furthermore, we show that the pheromone distribution can be used to improve the performance of the Fast Parzen Window density estimator. In addition to an empirical analysis we provide a thorough theoretical analysis to verify the ability of pheromone to encode information about distances of data points to the underlying manifold. First we only focus on a linear manifold and use the Markov Chain framework to study the behavior of the algorithm. In order to extend our analysis to non-linear manifolds we assume that a manifold can be approximated by locally linear patches, and we analyze how the pheromone distribution asymptotically behaves. Our analysis shows that the pheromone sorts data points according to their distances to the linear patch, i.e. the closer a data point is to the linear patch, the higher pheromone level it has. In summary, we propose an algorithm that provides valuable information, called pheromone, about how far points are from the un-

derlying manifolds. This information is beneficial not only to extract noisy manifolds buried in point clouds, but also to build more effective density estimators.

The organization of this contribution is as follows. In section 2, we provide background information to study random walks in general. Section 3 contains three alternative transition probabilities for random walks and our proposed method to extract manifolds from a point cloud is exhibited in section 4. A theoretical analysis of the role of the pheromone and its convergence is presented in section 5. Section 6 contains an empirical analysis on synthetic data sets and real-world application examples, demonstrating the performance in recovering manifolds and improving visualization and density estimation techniques.

## 2 Background and notation

Let  $X = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_N] \in \mathbb{R}^{D \times N}$  be the data matrix storing  $N$   $D$ -dimensional data points in its columns. Given a number  $r > 0$ , we define the neighborhood of any point  $\mathbf{x}_i$  as follows:

$$B_r(\mathbf{x}_i) = \{\mathbf{x}_j \mid \|\mathbf{x}_i - \mathbf{x}_j\|^2 < r, 1 \leq j \leq N\} \ , \quad (1)$$

where  $\|\cdot\|$  denotes the Euclidean norm. This definition of neighborhood has an advantage of being geometrically motivated (Belkin and Niyogi, 2003), which also inspires our theoretical analysis. Alternatives are methods such as k-nearest neighbor and mutual k-nearest neighbor (Von Luxburg, 2007). In the neighborhood graph nodes correspond to the data points from  $X$  and  $(\mathbf{x}_i, \mathbf{x}_k)$  represents the directed edge connecting  $\mathbf{x}_i$  to  $\mathbf{x}_k$ , with the associated weight  $w(\mathbf{x}_i, \mathbf{x}_k) \geq 0$ . The usual way to describe a Markov



Chain (MC) is to use a weighted graph where the nodes denote the states of the MC, and its weights represent the probability of transitions between the states. Thus, from the above neighborhood graph one can construct an MC by row-normalization, i.e.:

$$p_{ij} = \frac{w(\mathbf{x}_i, \mathbf{x}_j)}{\sum_{k \in B_r(\mathbf{x}_i)} w(\mathbf{x}_i, \mathbf{x}_k)} \quad (2)$$

where  $p_{ij}$  is the transition probability of jumping from  $\mathbf{x}_i$  to  $\mathbf{x}_j$ . A popular MC process is the random walk where a walker follows the transition probabilities to move on the graph. We will now review some of the key results and notions from the MC theory that will be needed in our study.

Assume  $S = \{1, 2, \dots, N\}$  represents the state space of a system. Let  $X_n$  denote the random variable of the state of the system in the  $n$ -th time step. In an MC, we assume that the next state of the system only depends on the current state, and the system's behavior is described through transition probabilities:

$$p_{ij} = P(X_1 = j | X_0 = i), \quad \forall i, j \in S \quad (3)$$

with  $X_0$  being the initial state of the system. In this section, we only consider homogeneous MCs where the transition probabilities are independent of the time  $n$ . Now, let  $P = [p_{ij}]$  be a matrix containing the transition probabilities. While  $P$  includes the transition probabilities for one time step, it can be extended to arbitrary number of steps:

**Theorem 1.** *For a MC with the transition probability matrix  $P$ , we have*

$$P(X_n = j | X_0 = i) = p_{ij}^{(n)}$$

where  $p_{ij}^{(n)}$  is  $(i, j)$ -th element of the  $n$ -th power of the matrix  $P$ , i.e.  $P^n = \underbrace{P \times \dots \times P}_{n \text{ times}}$  (Kulkarni, 1999).

This distribution is used to study an MC in a limited time interval. Let  $N_j(n)$  denote the number of visits state  $j$  accumulates over  $(n + 1)$  steps including the initial state. To study the behavior of a chain in this interval the occupancy time is defined as follows:

**Definition 2.1.** For any state  $j$  and initial state  $i$ , the occupancy time of state  $j$  up to  $n$  is the expected number of times the system spends in state  $j$  in a random walk of  $n + 1$  steps starting in  $i$ , i.e.:

$$m_{ij}(n) = E[N_j(n)|X_0 = i] ,$$

where  $E[\cdot]$  denotes the expectation operator.

Let  $M(n) = [m_{ij}(n)]$  denote the occupancy time matrix. The following theorem connects the occupancy time to the power of transition matrix (Kulkarni, 1999):

**Theorem 2.** Given an MC, its occupancy time matrix is

$$M(n) = \sum_{q=0}^n P^q . \tag{4}$$

While the above theorem helps to study an MC in a finite number of time steps, we would like to investigate the long-term behavior of the chain, which is encoded in the stationary distribution:

**Definition 2.2.** For an MC with the transition probability matrix  $P$  the distribution  $\pi = [\pi_1, \dots, \pi_N]$  is called its stationary distribution if the following balance equation holds:

$$\pi P = \pi . \tag{5}$$

Here, we only focus on a specific class of Markov Chains called *ergodic chains*.

**Definition 2.3.** An MC is called ergodic if the following condition holds:

$$\exists T \in \mathbb{N} : \forall n > T, \forall i, j \in S \quad P(X_n = j | X_0 = i) > 0 .$$

In other words, for sufficiently large time steps  $n$  the system can be found in any state.

The following theorem shows that the stationary distribution captures all information about the long-term behavior of an ergodic MC (Cinlar, 2013; Kulkarni, 1999).

**Theorem 3.** For an ergodic MC its stationary distribution  $\pi$  is unique with:

$$\forall i \in S; \lim_{n \rightarrow \infty} P(X_n = j | X_0 = i) = \pi_j , \quad (6)$$

$$\lim_{n \rightarrow \infty} P(X_n = j) = \pi_j , \quad (7)$$

$$\text{and the expected return time to state } j \text{ is } \mu_j = \frac{1}{\pi_j} . \quad (8)$$

### 3 Homogeneous Markov Chain

An important class of manifold learning techniques relies on weighted graph representations of data. Markov Chains constructed on graphs offer a powerful tool to represent non-linear manifolds and are used in several dimensionality reduction techniques. This is exemplified by a well-known family of dimensionality reduction techniques called "kernel eigenmap methods", which includes local linear embedding (Roweis and Saul, 2000); Laplacian Eigenmaps (Belkin and Niyogi, 2003); Hessian Eigenmaps (Donoho and Grimes, 2003); local tangent space alignment (Zhang and Zha, 2003) and diffusion map (Coifman and Lafon, 2006). The general idea is that the eigenvectors of Markov transition matrices are used to project high-dimensional data to a lower dimensional Euclidean space preserving the main structures of the data (Coifman and Lafon, 2006).

Therefore, applying MCs in the context of manifold learning has been employed in several techniques.

Although manifold learning techniques assume that the data lies on a lower-dimensional topological structure, in practice, it rarely clear due to the presence of noise. One way to improve their performance is to use "less noisy data" closely aligned with the underlying manifold. In this contribution, the goal is to perform a random walk that places emphasis on such manifold aligned sub-samples. In the following, we introduce three MCs that can recover a manifold by highlighting sample points closer to it.

### 3.1 Weights Based on Kernels

A common way to construct an MC on a data set is to use kernel functions (Berry and Sauer, 2016), i.e. a map  $K : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$  that fulfills:

- $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^D \quad K(\mathbf{x}, \mathbf{y}) = K(\mathbf{y}, \mathbf{x})$
- $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^D \quad K(\mathbf{x}, \mathbf{y}) \geq 0$  ,

and quantifies a "similarity" between pairs of data points. The most popular kernel function is the Gaussian, which for two points  $\mathbf{x}_i$  and  $\mathbf{x}_j$  is defined as:

$$K_\sigma(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right) , \quad (9)$$

with the Euclidean distance  $\|\cdot\|$  and scale parameter  $\sigma > 0$ . Based on the Gaussian kernel the transition probability reads:

$$p_{ij} = \frac{K_\sigma(\mathbf{x}_i, \mathbf{x}_j)}{\sum_{k \in \mathcal{N}_i} K_\sigma(\mathbf{x}_i, \mathbf{x}_k)}$$

with  $\mathcal{N}_i$  being the set of point  $\mathbf{x}_i$ 's neighbors. This definition of transition probabilities reinforces random walkers to spend more time in denser regions.

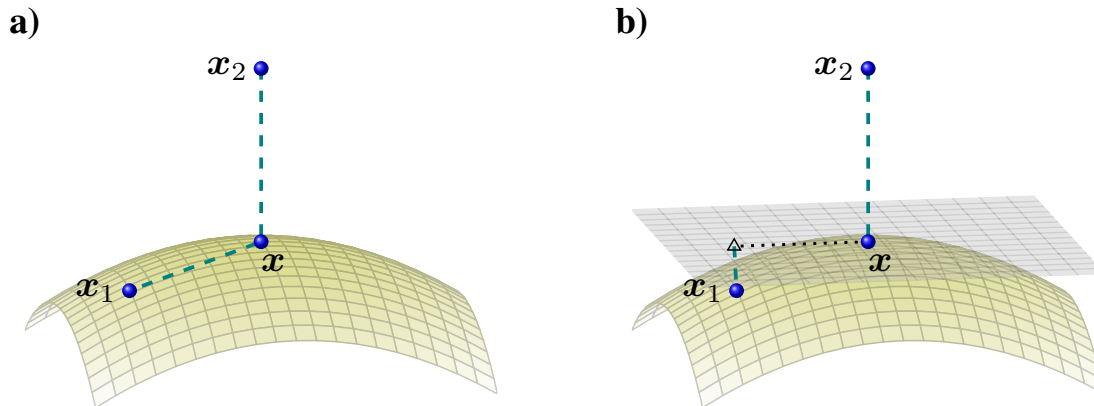


Figure 1: An illustration of a curved manifold. In contrast to  $x_2$  the points  $x$  and  $x_1$  lie on the curved manifold. The Euclidean distance (dashed) to both  $x_1$  and  $x_2$  is equal (see panel a). As shown in panel b the distance to the tangent (dashed) reflects the closeness to the manifold much better.

### 3.2 Weights Based on Tangent Spaces

Assuming the presence of a manifold the noise level of a data point can be related to its distance from the manifold. Thus, to form a MC that encourages random walkers to spend more time on the manifold the walkers need to know how far data points are from the manifold. As depicted in Fig. 1 a) the Euclidean distance fails to reveal this information. There, if a random walker resides on point  $x$  the Gaussian kernel gives the walker the same chance to stay on the manifold (point  $x_1$ ) or to leave it (jumping to point  $x_2$ )<sup>1</sup>. In order to overcome this drawback, we need to define a favorable measure that improves the approximate distances to the manifold.

A manifold can be approximated locally at a point  $x$  by its tangent space (Tu, 2011).

<sup>1</sup>A similar problem holds for the GDC, introduced in Klicpera et al. (2019). Since it enforces symmetric weights between every pairs of point, it uses a similar weight for jumping from  $x$  to  $x_2$  and vice versa.

Hence, we propose to use the tangent space to estimate the distance of  $\mathbf{x}$ 's neighbors to the manifold (see Figure 1b). Let  $\mathcal{N}$  represent the set of  $\mathbf{x}$ 's neighbors, for example enclosed in  $B_r(\mathbf{x})$ . The tangent space at  $\mathbf{x}$  is typically approximated by Principal Component Analysis (PCA) on  $\mathcal{N}$ <sup>2</sup>, yielding a set of eigenvalues and unit orthogonal eigenvectors of the local co-variance matrix,  $\{(\lambda_k, u_k)\}_{k=1}^D$ . Without loss of generality, we assume that the eigenvalues are in descending order with  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D \geq 0$  and normalized to  $\sum_k \lambda_k = 1$ . Then for a  $d$ -dimensional manifold the subspace spanned by the columns of  $U_1 = [u_1, u_2, \dots, u_d]$  provides an estimation for the tangent space and locally for the manifold.

Since we consider noisy data sets we only have access to a noisy version of the tangent space  $\hat{U}_1$ . Hence, depending on the noise level the difference between  $U_1$  and  $\hat{U}_1$  may be small or high. For more details, we refer the reader to (Kaslovsky and Meyer, 2014) and (Little et al., 2017). In the following we introduce two ways to define MCs highlighting samples close to manifold structures.

**A manifold with known intrinsic dimensionality:** Consider a  $d$ -dimensional manifold in  $\mathbf{x}$ 's neighborhood. For any  $\mathbf{x}_i \in \mathcal{N}$  one can estimate the distance to the tangent space by:

$$\delta_i^{\mathcal{M}} = \|(I - \hat{U}_1 \hat{U}_1^T)(\mathbf{x}_i - \mathbf{x}')\| , \quad (10)$$

where  $\mathbf{x}'$  determines the place where the tangent space touches the manifold in the neighborhood. Based on these quantities, we define the weight values to be used in the

---

<sup>2</sup>Note that methods, such as (Lerman et al., 2015), provide better estimation for the tangent spaces.

However, since in this contribution we consider big data sets we use the less costly PCA.

calculation of transition probabilities Eq. (2) as follows:

$$w_d(\mathbf{x}, \mathbf{x}_i) = \begin{cases} 1 - \frac{\delta_i^M}{\alpha} & \text{if } \alpha \geq \Delta_i \\ 0 & \text{if } \alpha < \Delta_i . \end{cases} \quad (11)$$

where  $\alpha$  is a factor to ensure the weight values are positive. Here, we determine  $\alpha$  such that only  $p$  percent of neighbors have non-zero weights. This definition constructs a random walk favoring jumps closer to the manifold. In other words, the random walker observes a flat region and it is more likely to move to a point close to it.

**A manifold with unknown intrinsic dimensionality** In most applications the dimensionality of a manifold  $d$  is unknown and eigenvalues are typically used to estimate it. As suggested in (Wang et al., 2008) the intrinsic dimensionality may be estimated by:

$$\hat{d} = \arg \max_d S_d ,$$

where  $S_d = d \cdot (\lambda_d - \lambda_{d+1})$ ,  $d \in \{1, \dots, D - 1\}$  and  $S_D = D \cdot \lambda_D$ . However, the performance of these types of criteria highly depend on the neighborhood size which is related to the noise level and the manifold curvature. In order to tackle this problem, we consider all possible values of  $d$  instead of picking a specific value. The basic idea is to compute the weight values for any  $d \in \{1, \dots, D\}$ , according to Eq. (11), followed by the calculation of their mean:

$$w(\mathbf{x}, \mathbf{x}_i) = \sum_{d=1}^D S_d \cdot w_d(\mathbf{x}, \mathbf{x}_i) , \quad (12)$$

where the eigengap  $S_d$  indicates the importance of each intrinsic dimensionality. Since  $S_d \geq 0$  and  $\sum_{d=1}^D S_d = 1$ , we may interpret  $S_d$  as the ‘‘probability’’ of the manifold being  $d$ -dimensional and the weight value  $w$  as the expected weight with respect to this distribution over manifold dimensionalities.

## 4 Ant Colony: Non-Homogeneous Markov Chain

The group behavior of decentralized natural systems, such as ant colony and bird flocking, has inspired many methods in computer science and is often summarized under the keyword Swarm Intelligence (SI). For instance, the biological behavior of ants in the food-seeking process has motivated several methods in combinatorial optimization and clustering. In the process two mechanisms, one behavioral and one environmental, help to find the shortest path to a food source:

- **Deposition:** when a bio-ant walks to (and from) a food source it releases a substance called pheromone on the ground. Ants are attracted by it when they choose a path to follow. The more ants use a route the more pheromone accumulates increasing the chance for the path to be selected by subsequent ants as well.
- **Evaporation:** the pheromone evaporates over time and hence less attractive trails with less pheromone eventually disappear.

Thus, although there is no central authority to control the ants' behavior, their indirect form of local interactions, via pheromone, helps them to find the shortest path to the closest food source.

In this contribution, we propose a new algorithm that aims to recover manifolds from noisy samples. The basic idea is that a set of ant-like agents are released in the data space to search for manifolds based on the following elements:

O1: an MC that highlights underlying manifolds in the data set, and

O2: deposition and evaporation mechanisms to update the pheromone.



The aim is that ants are more likely to visit points close to manifolds and deposit pheromone. Higher pheromone levels function as a positive feedback mechanism reinforcing further visitations. Eventually, by extracting points with more pheromone potential manifolds can be uncovered. In the following we formally define and illustrate the two elements outlined above.

#### 4.1 Manifold Alignment Aware Ants (M3A)

Inspired by swarm intelligence, M3A uses several ants who walk in the data space and search for manifolds. During the search process, the ants interact with one another through pheromone and with the environment. In section 3, we have already proposed possible formulations of transition probabilities for the random walk. However, they consider only the distribution of data points in the environment and do not provide any form of communication among ants. Thus, we modify the transition probabilities by adding a pheromone factor  $f$ . Given any data point  $\mathbf{x}_i$ , let  $f_i$  denote the amount of pheromone on it and  $\mathcal{N}_i$  represents the set of its neighbors. The new transition probability is defined as:

$$p_{ij} = \frac{(w_{ij})^{1-\gamma} (\hat{f}_j)^\gamma}{\sum_{\mathbf{x}_k \in \mathcal{N}_i} (w_{ik})^{1-\gamma} (\hat{f}_k)^\gamma}, \quad (13)$$

where  $w_{ij} = w(\mathbf{x}_i, \mathbf{x}_j)$ . To ensure that the weights  $w$  and the pheromone  $f$  are in the same scale, i.e.  $[0, 1]$ , the pheromone is normalized within the neighborhood by  $\hat{f}_j = \frac{f_j}{\sum_{\mathbf{x}_k \in \mathcal{N}_i} f_k}$ . The  $\gamma \in [0, 1]$  effectively controls how much ants are attracted to structures already found or are able to explore the data space for new regions. For instance, if we set  $\gamma$  to zero, then there would be no interaction among ants and, as a

result, every ant could independently explore the data space. Hence, the new transition probability takes into account both types of interactions O1 and O2.

Pheromone is an indirect form of communication among artificial ants who release it on the visited points. Therefore, the amount of pheromone on a sample varies and it should be updated over time. In order to make the algorithm suitable for parallelization the pheromone values are updated when the ants finish their walks of a predefined number of steps  $n$ , called *one round*<sup>3</sup>. Thus, the transition probability for the  $(t + 1)$ -th round can be re-written as:

$$p_{ij}^{(t+1)} = \frac{(w_{ij})^{1-\gamma} \left(\hat{f}_j^{(t)}\right)^\gamma}{\sum_{k \in \mathcal{N}_i} (w_{ik})^{1-\gamma} \left(\hat{f}_k^{(t)}\right)^\gamma}, \quad (14)$$

where  $f_j^{(t)}$  denotes the amount of pheromone on  $x_j$  after  $t$  rounds. Since the amount of pheromone on a data point depends on the number of visits and the extent of evaporation the updating rule consists of two parts:

- deposition based on the number of times a point has been visited, and
- evaporation depends on the environment and is controlled by hyper-parameter  $\rho$ .

Therefore, the pheromone level on any sample  $x_i$  can be updated as follows:

$$f_i^{(t+1)} = \frac{c}{M} \sum_{a=1}^M \frac{N_i^a(n)}{n+1} + (1 - \rho) \cdot f_i^{(t)}, \quad (15)$$

where  $M$  is the number of ants and  $N_i^a(n)$  is the number of times the  $a$ -th ant visits  $x_i$  over  $n + 1$  steps (including initial state) in the  $(t + 1)$ -th round. Note that the deposition

---

<sup>3</sup>In one round, we can distribute the agents among several processors and update the pheromone after  $n$  steps. Otherwise, communications between processors would be necessary to synchronize pheromone values in every step, which would cause unnecessary overhead.

---

**Algorithm 1: OneAnt**

---

```
1 Input: pheromone  $\mathbf{f} = [f_1, f_2, \dots, f_N]$ ;  
2  $\mathbf{N}_v = [0, 0, \dots, 0]$  ;  
3 Randomly select a node as the initial position (say  $\mathbf{x}_i$ ) ;  
4  $N_v[i] = N_v[i] + 1$ ;  
5 for  $s = 1$  to  $n$  do  
6   | randomly select its next destination following Eq. (13) (say  $\mathbf{x}_j$ );  
7   |  $N_v[j] = N_v[j] + 1$ ;  
8 end
```

**Result:** number of visits  $\mathbf{N}_v$ .

---

term is divided by  $M(n + 1)$  to prevent unlimited increase of the pheromone level (especially for big  $M$  or  $n$ ). The constant  $c > 0$  specifies the amount of pheromone a single ant deposits on a point in a single visit. The pseudo code of the M3A<sup>4</sup> is provided in Algorithm 2. In summary, we can highlight the following points:

- The proposed ant algorithm uses transition probabilities (Eq. 14) based on local tangent alignment and pheromone to highlight manifold structures, with the latter also reinforcing the agents to stay close.
- In contrast to previous works with the ant colony, the pheromone is deposited on the nodes, instead of edges. This vastly reduces the number of pheromone values that need to be stored. Moreover, the associated pheromone to each data point can be used to extract points close to manifolds.

---

<sup>4</sup>M3A implementation is provided at <https://github.com/mohammadimathstar/M3A>

---

**Algorithm 2:** Manifold Alignment Aware Ants (M3A)

---

```
1 Initialize the pheromone vector  $\mathbf{f} = [f_1^{(0)}, f_2^{(0)}, \dots, f_N^{(0)}]$ ;  
2 for  $t = 1$  to  $N_{iter}$  do  
3    $\mathbf{N}_v = [0, 0, \dots, 0]$ ;  
4   for  $a = 1$  to  $M$  do  
5      $\mathbf{N}_v^a = \text{OneAnt}(\mathbf{f})$ ;  
6   end  
7    $\mathbf{N}_v = \sum_{a=1}^M \mathbf{N}_v^a$ ;  
8   Update pheromone by Eq. (15).  
9 end
```

**Result:** The pheromone vector  $\mathbf{f} = [f_1, f_2, \dots, f_N]$ .

---

- The amount of pheromone is updated after a pre-defined number of steps  $n$ . In practice, this allows to distribute the computations among several processors and parallelize the random walks.

## 4.2 Complexity analysis

As a preprocessing step for the M3A algorithm, we implement two operations on each data point: a) finding its neighbors and b) performing PCA on its neighborhood. To perform the neighbor search we need to calculate the distances of all pairs of  $N$  samples, which leads to the complexity of  $\mathcal{O}(N^2)$ . However, there are approximate nearest neighbor search strategies to reduce the complexity. For instance, the k-d tree algorithm (Bentley, 1975) with the complexity  $\mathcal{O}(N \log N)$  can significantly decrease the

computational costs. Besides the neighbor search we perform local PCA using Singular Value Decomposition. In the worst case, its implementation for any point  $x$  scales cubic with the dimensionality of the data and quadratic with the size of its neighborhood  $|\mathcal{N}_x|$ . Nonetheless, approximate strategies decrease the computational complexity (Golub and Van Loan, 2012). In addition to speeding up the preprocessing step via approximate strategies the execution of the M3A can be accelerated via parallelization. We only need to distribute the ants (line 5 in algorithm 2) among several processors.

## 5 Theoretical Analysis

Many metaheuristic methods suffer from the lack of theoretical analysis. Therefore, their behavior and the effects of parameters on them are typically investigated only empirically. In this section, we provide an analysis of the convergence and the impact of hyper-parameters of the proposed algorithm in exemplary situations. More precisely, we consider a data set containing a noisy  $d$ -dimensional manifold and we use Eq. (11) to compute the weight values with  $\alpha$  determined such that  $p = 100\%$  of neighbors having non-zero weight values. First, we study the effect of noise on the performance of PCA and, as a result, on the M3A algorithm. Then, we examine the pheromone distribution and demonstrate its convergence to the stationary distribution and its capability in recovering a linear manifold. Finally, we study the performance of the algorithm on non-linear manifolds under the assumption that they can be approximated by local linear patches. We concentrate our analysis on a single patch and consider the other patches as conceptually grouped in a single state representing “the outside”.

## 5.1 Spectral Analysis

An important subject in perturbation theory is to study the effect of noise on the eigenvalues and eigenvectors of matrices (Kaslovsky and Meyer, 2014). Since we use principal directions of co-variance matrices to define the MCs, it is important to know the effect of noise on the eigenvectors of co-variance matrices. Let  $\mathcal{M}$  be a  $d$ -dimensional vector subspace embedded in a higher dimensional space<sup>5</sup>. In the case of noise-free samples it can be exactly recovered via PCA. However, in the presence of noise the recovered subspace  $\hat{\mathcal{M}}$  is perturbed. Here, the goal is to see the impact of noise on  $\hat{\mathcal{M}}$ .

Without loss of generality, let us consider the first  $d$  coordinates span the subspace  $\mathcal{M}$ . We assume the set  $\{\mathbf{l}_i\}_{i=1}^N$  contains  $N$  realizations of  $\mathcal{M}$ , such that their first  $d$  coordinates  $(l_i^{(1)}, \dots, l_i^{(d)})$  are uniformly distributed within  $B_r(\mathbf{0})$ . Thus, every point  $\mathbf{l}_i$  on the subspace  $\mathcal{M}$  has the form  $(l_i^{(1)}, \dots, l_i^{(d)}, 0, \dots, 0)$ . Now, suppose these realizations are disrupted via Gaussian noise with mean zero and standard deviation  $\sigma I_D$ . If  $\mathbf{x}_i$  denotes a noisy observation, then we have:

$$\mathbf{x}_i = \mathbf{l}_i + \mathbf{e}_i \text{ ,}$$

where  $\mathbf{e}_i$  is the noise vector. In this setting, the design matrix  $X$  can be described as:

$$X = L + E \text{ ,} \tag{16}$$

where the columns of  $L$  and  $E$  keep the noise-free realizations and noise vectors, re-

---

<sup>5</sup>The assumptions are made since we only focus on linear manifolds in the next subsection. For a more general case, i.e. non-linear manifold, we refer to (Kaslovsky and Meyer, 2014) where the effect of curvature on the covariance matrices is studied.

spectively. Let us denote the centered version of a matrix  $H = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N]$  as:

$$\tilde{H} = H - \mu_h \mathbf{1}_N ,$$

where  $\mu_h = \frac{1}{N} \sum_i \mathbf{h}_i$  and  $\mathbf{1}_N = [1, 1, \dots, 1]$  with  $N$  entries. Then the data co-variance matrix can be written as

$$\frac{1}{N} \tilde{X} \tilde{X}^T = \frac{1}{N} \tilde{L} \tilde{L}^T + \Delta$$

with the matrix  $\Delta$  representing the perturbation caused by the noise:

$$\Delta = \frac{1}{N} (\tilde{L} \tilde{E}^T + \tilde{E} \tilde{L}^T + \tilde{E} \tilde{E}^T) .$$

From eigendecomposition, we obtain:

$$\frac{1}{N} \tilde{L} \tilde{L}^T = U \Lambda U^T = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} \Lambda_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} U_1 & U_2 \end{bmatrix}^T ,$$

where  $\Lambda_1$  is a  $d \times d$  diagonal matrix containing non-zero eigenvalues in descending order, and  $U_1$  includes their corresponding eigenvectors. Note that  $U_2$  can be any orthogonal basis for the last  $D - d$  coordinates. Similarly, we can write:

$$\frac{1}{N} \tilde{X} \tilde{X}^T = \hat{U} \hat{\Lambda} \hat{U}^T = \begin{bmatrix} \hat{U}_1 & \hat{U}_2 \end{bmatrix} \begin{bmatrix} \hat{\Lambda}_1 & 0 \\ 0 & \hat{\Lambda}_2 \end{bmatrix} \begin{bmatrix} \hat{U}_1 & \hat{U}_2 \end{bmatrix}^T .$$

The subspace  $\mathcal{M}$  is spanned by  $U_1$  and the recovered subspace  $\hat{\mathcal{M}}$  is spanned by  $\hat{U}_1$ .

The orthogonal projectors onto  $\mathcal{M}$  and  $\hat{\mathcal{M}}$  are derived as follow:

$$Q = U_1 U_1^T \quad \hat{Q} = \hat{U}_1 \hat{U}_1^T .$$

Here, the Frobenius distance  $\|Q - \hat{Q}\|_F$  is used to compare two subspaces because it corresponds to the sum of the squared sines of the principal angles between  $\mathcal{M}$  and  $\hat{\mathcal{M}}$

(Kaslovsky and Meyer, 2014). If the number of samples  $N$  and the probability constants  $\varepsilon$  and  $\varepsilon_\lambda$  satisfy the following inequalities:

$$N > 4 \left( \max(\sqrt{d}, \sqrt{D-d}) + \varepsilon \right), \quad \varepsilon < 0.7\sqrt{d(D-d)}, \quad \varepsilon_\lambda < \frac{3}{\sqrt{d+2}}\sqrt{N} \quad (17)$$

then the following theorem offers a bound on the angle between  $\mathcal{M}$  and  $\hat{\mathcal{M}}$ .

**Theorem 4.** (Kaslovsky and Meyer, 2014) *Let*

$$\delta = \frac{r^2}{d+2} \left( 1 - \frac{1}{\sqrt{N}}\zeta_1(\varepsilon_\lambda) \right) - \sigma \frac{1}{\sqrt{N}}\zeta_2(\varepsilon_\lambda) - \sigma^2 \left( \sqrt{d} + \sqrt{D-d} + \frac{1}{\sqrt{N}}\zeta_3(\varepsilon) \right)$$

and

$$\beta = \frac{1}{\sqrt{N}} \left[ \sigma \sqrt{d(D-d)}\eta(\varepsilon, \varepsilon_\lambda) + \frac{1}{\sqrt{N}}\zeta_{\text{numer}}(\varepsilon) \right].$$

Additionally, if the following conditions hold:

- (Condition 1)  $\delta > 0$ ,
- (Condition 2)  $\beta < \frac{1}{2}\delta$

then

$$\|Q - \hat{Q}\| \leq \frac{2\sqrt{2}\beta}{\delta} \quad (18)$$

with probability greater than

$$1 - 2de^{-\varepsilon_\lambda^2} - 9e^{-\varepsilon^2} \quad (19)$$

over the joint random selection of the sample points and random realization of the noise,



where the following definitions have been made to ease the presentation:

$$\begin{aligned}
\zeta_1(\varepsilon_\lambda) &= \frac{2}{\sqrt{N}} - \frac{1}{N^{\frac{3}{2}}} + \left(1 - \frac{1}{N}\right)\varepsilon_\lambda\sqrt{8(d+2)} \ , \\
\zeta_2(\varepsilon, \varepsilon_\lambda) &= \frac{2rd}{\sqrt{d+2}} \left(1 + \varepsilon_\lambda \frac{5\sqrt{d+2}}{\sqrt{N}}\right) \left(1 + \frac{6\varepsilon}{5d}\right) \ , \\
\zeta_3(\varepsilon) &= \frac{5}{2} \left(\sqrt{d} + \varepsilon\sqrt{2}\right) \left(\sqrt{D-d} + \varepsilon\sqrt{2}\right) \ , \\
\eta(\varepsilon, \varepsilon_\lambda) &= \left(1 + \frac{6}{5} \frac{\varepsilon}{\sqrt{d(D-d)}}\right) \left[\sigma + \frac{r}{\sqrt{d+2}} \left(1 + \varepsilon_\lambda \frac{5\sqrt{d+2}}{\sqrt{N}}\right)\right] \ , \\
\zeta_{\text{numer}}(\varepsilon) &= \sigma^2\sqrt{d(D-2)} \left(1 + \frac{6}{5} \frac{\varepsilon}{\sqrt{d(D-d)}}\right) \left(\sqrt{D-d} + \varepsilon\sqrt{2}\right) \ .
\end{aligned}$$

This theorem provides an upper bound for the difference between the true manifold  $\mathcal{M}$  and the recovered subspace  $\hat{\mathcal{M}}$ . Therefore, it is safe to say that the definition of weights, in Eq. (11), is reasonable in keeping a random walker close to the  $\mathcal{M}$  (via encouraging to stay close to  $\hat{\mathcal{M}}$ ). Moreover, the above theorem can help us to explain the relation between the quality of the recovered subspace  $\hat{\mathcal{M}}$  and the hyper-parameters  $N$ ,  $\sigma$  and  $r$ . An empirical example demonstrates the relationship in practice. We generate samples from a noisy one-dimensional linear manifold embedded in  $\mathbb{R}^{10}$  for various (hyper-) parameter settings. Their influence is demonstrated in Figure 2. Panel (a) shows that the approximation performance increases with a growing number of samples  $N$ . On the other hand, there is a direct connection between the noise level  $\sigma$  and the approximation error as shown in panel (b). Moreover, Figure 2 (c) demonstrates that a bigger neighborhood radius  $r$  recovers an increasingly accurate subspace, which is not the case for non-linear manifolds, where this also depends on the curvature.

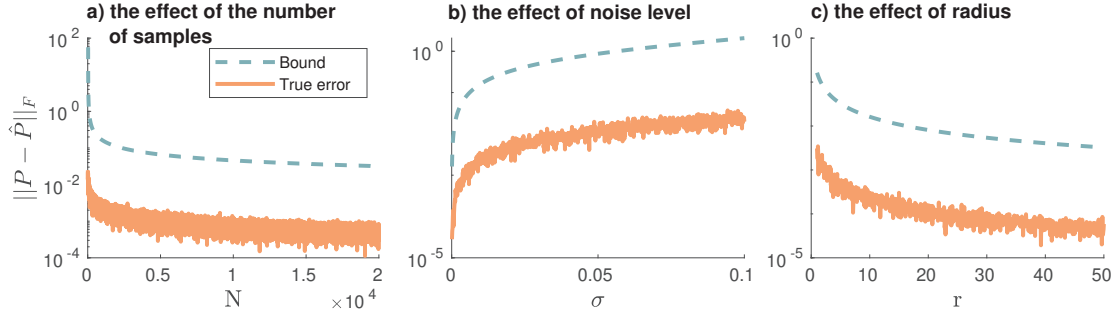


Figure 2: The effect of (hyper-)parameters on the quality of the recovered subspaces: a) the effect of the number of samples  $N$  within the neighborhood ( $\sigma = 0.01$ ,  $r = 1$ ), b) the role of the noise level  $\sigma$  ( $N = 1000$ ,  $r = 1$ ), and c) the impact of neighborhood radius  $r$  ( $N = 1000$ ,  $\sigma = 0.01$ ).

## 5.2 Pheromone distribution for a linear manifold

The proposed M3A algorithm associates a pheromone value with each sample. In this section, we show that the pheromone values encode distances of the samples to the reconstructed subspace  $\hat{\mathcal{M}}$ . Let the data set fulfill the conditions explained in the previous subsection, i.e. the noisy data points are uniformly distributed along a linear manifold inside  $B_r(\mathbf{0})$ . In order to simplify our analysis we make the following assumptions in constructing the MC:

- The radius  $r$  is big enough such that all pairs of samples are neighbors and any sample belongs to its neighborhood. Therefore, instead of using local PCA for each point we apply PCA on the whole data set.
- In Eq. (10) we set  $\mathbf{x}'$  to the origin  $\mathbf{0}$ . Thus, for any point  $\mathbf{x}_i$  its distance to  $\hat{\mathcal{M}}$  is:

$$\delta_i^{\mathcal{M}} = \|(I - \hat{U}_1 \hat{U}_1^T) \mathbf{x}_i\| ,$$

and with Eq. (11) the weight values are computed as:

$$w_{ki} = 1 - \frac{\delta_i^M}{\alpha} = w_i \quad 1 \leq k \leq N . \quad (20)$$

We define  $w_i = w_{ki}$  since  $w_{ki}$  does not depend on  $k$ .

- Since the presence of powers in the transition probabilities (i.e.  $\gamma$  and  $1 - \gamma$ ) makes our analysis intractable, we analyse the special case with removed powers, such that our arguments can be presented in a tractable manner:

$$p_{ji}^{(t+1)} = \frac{w_i f_i^{(t)}}{\sum_k w_k f_k^{(t)}} \quad (21)$$

- The pheromone values are uniformly initialized by:

$$f_i^{(0)} = \frac{1}{N} \quad \forall i . \quad (22)$$

- In the updating rule of pheromone Eq. (15) we set the constant  $c$  to the evaporation rate  $\rho$ :

$$f_i^{(t+1)} = \frac{\rho}{M} \sum_{a=1}^M \frac{N_i^a(n)}{n+1} + (1 - \rho) f_i^{(t)} . \quad (23)$$

Without loss of generality we re-label the data points by  $\{1, \dots, N\}$  such that

$$w_1 \geq w_2 \geq \dots \geq w_N . \quad (24)$$

Thus, the new labels sort samples according to their closeness to the subspace  $\hat{\mathcal{M}}$ .

**Corollary 4.1.** *The pheromone values across the samples form a probability distribution called pheromone distribution denoted by  $\mathbf{f}^{(t)} = [f_1^{(t)}, f_2^{(t)}, \dots, f_N^{(t)}]$ .*

*Proof.* It can be proven by induction over  $t$ . From (22) we have  $\sum_{i=1}^N f_i^{(0)} = 1$  and if we assume that the statement is valid for  $t$ , i.e.  $\sum_{i=1}^N f_i^{(t)} = 1$ , we can re-write (23):

$$\sum_{i=1}^N f_i^{(t+1)} = \frac{\rho}{M} \sum_{a=1}^M \left( \frac{1}{n+1} \overbrace{\sum_{i=1}^N N_i^a(n)}^{=1} \right) + 1 - \rho = 1 .$$

Note that the expression inside the parentheses is equal to 1 since  $\sum_{i=1}^N N_i^a(n)$  equals to the number of steps (i.e.  $n$ ) plus 1 (for the initial point).  $\square$

In this setting we show that:

- the pheromone values on sample points are sorted by their distances to the subspace  $\hat{\mathcal{M}}$ , and
- the pheromone distribution  $\mathbf{f}^{(t)}$  converges as  $t \rightarrow \infty$ .

Our analysis is based on the Markov Chain framework and we consider two cases. First, we assume there is only one ant who walks for an unbounded number of steps in every round (i.e.  $n \rightarrow \infty$ ). Second, we assume an unbounded number of ants (i.e.  $M \rightarrow \infty$ ) who walk for  $n$  steps each in every round. In other words, we study the long-term and the short-term behavior of random walks, respectively.

### 5.2.1 Single ant, unbounded path length

We assume a single ant (i.e.  $M = 1$ ) performing  $n \rightarrow \infty$  steps in each round. From theorem 3, the long-term fraction of time spending in a point  $x_i$  is equal to its stationary distribution value  $\pi_i$ . Therefore, for any point  $x_i$ , its pheromone value is updated

according to:

$$\begin{aligned}
f_i^{(t+1)} &= \frac{\rho}{M} \sum_{a=1}^M \pi_i^{(t+1)} + (1 - \rho) f_i^{(t)} \\
&= (1 - \rho) f_i^{(t)} + \rho \cdot \pi_i^{(t+1)} ,
\end{aligned} \tag{25}$$

where  $\pi^{(t+1)}$  is the stationary distribution associated to the transition probability matrix in the  $(t + 1)$ -th round. The stationary distribution can be computed via the balance equation:

$$\begin{aligned}
\pi_i^{(t+1)} &= \sum_j \pi_j^{(t+1)} p_{ji}^{(t+1)} \\
&= \sum_j \pi_j^{(t+1)} \frac{w_i f_i^{(t)}}{\sum_k w_k f_k^{(t)}} \\
&= \frac{w_i f_i^{(t)}}{\sum_k w_k f_k^{(t)}} \\
&= p_{ji}^{(t+1)} \quad \forall j .
\end{aligned} \tag{26}$$

From (22) and (24), it can be shown (by induction) for any  $t > 0$ :

$$\pi_1^{(t)} \geq \pi_2^{(t)} \geq \dots \geq \pi_N^{(t)} \tag{27}$$

and by (25):

$$f_1^{(t)} \geq f_2^{(t)} \geq \dots \geq f_N^{(t)} . \tag{28}$$

Since the weight values encode the distance of samples to the subspace  $\hat{\mathcal{M}}$ , see Eq. (11), the following corollary holds.

**Corollary 4.2.** *The pheromone values are sorted according to their closeness to the linear manifold  $\hat{\mathcal{M}}$ .*

In the next step, we would like to prove the convergence of the pheromone distribution, but we need to first establish the following lemma:

*Lemma 5.* The sequence  $\{a^{(t)}\}_{t=0}^{\infty}$  with

$$a^{(t)} = \sum_k w_k f_k^{(t)} \quad (29)$$

- is monotonically increasing (i.e.  $a^{(t+1)} \geq a^{(t)}$ ) and
- is convergent (i.e.  $a^{(t)} \rightarrow a$  for some  $a \geq 0$ ).

*Proof.* a) To demonstrate that the sequence  $\{a^{(t)}\}_{t=0}^{\infty}$  is monotonically increasing, we need to show:

$$I^{(t)} = a^{(t+1)} - a^{(t)} = \sum_k w_k \cdot \Delta^{(t)} f_k \geq 0 ,$$

where  $\Delta^{(t)} f_i = f_i^{(t+1)} - f_i^{(t)}$ . We can write:

$$\begin{aligned} \Delta^{(t)} f_i &\stackrel{\text{Eq.(25)}}{=} \rho(\pi_i^{(t+1)} - f_i^{(t)}) \\ &\stackrel{\text{Eq.(26)}}{=} \rho \left( \frac{w_i}{\sum_k w_k f_k^{(t)}} - 1 \right) f_i^{(t)} \\ &= \frac{\rho}{a^{(t)}} (w_i - a^{(t)}) f_i^{(t)} . \end{aligned} \quad (30)$$

For every fixed  $t$  we have:  $\{f_i^{(t)}\}_{i=1}^N$  and  $\{w_i - a^{(t)}\}_{i=1}^N$ <sup>6</sup> that are monotonically decreasing sequences. Thus, the sequence  $\{\Delta^{(t)} f_i\}_{i=1}^N$  will be decreasing as well. Therefore, we have:

$$I^{(t)} = \sum_{k \leq k^*} w_k \cdot \overbrace{\Delta^{(t)} f_k}^{\geq 0} + \sum_{k > k^*} w_k \cdot \overbrace{\Delta^{(t)} f_k}^{< 0}$$

where  $k^* = \max\{k : \Delta^{(t)} f_k \geq 0\}$  and we define the following values:

$$w = \min_{k \leq k^*} w_k \quad , \quad w' = \max_{k > k^*} w_k .$$

Since  $\{w_i\}_{i=1}^N$  is decreasing we have  $w > w'$ :

---

<sup>6</sup>From inequality (24) we know that  $\{w_i\}_{i=1}^N$  is a decreasing sequence. Therefore, the sequence  $\{w_i - a^{(t)}\}_{i=1}^N$  is also decreasing, since  $a^{(t)}$  is a constant number for a fixed  $t$ .

$$\begin{aligned}
I^{(t)} &\geq w \cdot \sum_{k \leq k^*} \overbrace{\Delta^{(t)} f_k}^{\geq 0} + w' \cdot \sum_{k > k^*} \overbrace{\Delta^{(t)} f_k}^{< 0} \\
&\stackrel{w > w'}{>} w \cdot \sum_{k \leq k^*} \Delta^{(t)} f_k + w \cdot \sum_{k > k^*} \Delta^{(t)} f_k \\
&= w \cdot \sum_k \Delta^{(t)} f_k = w \cdot \sum_k (f_k^{(t+1)} - f_k^{(t)}) \\
&= w \cdot \left( \sum_k f_k^{(t+1)} - \sum_k f_k^{(t)} \right) = 0 .
\end{aligned}$$

b) We know that:

$$0 \leq w_i \leq 1, \quad 0 \leq f_i^{(t)} \leq 1$$

and therefore

$$0 \leq a^{(t)} \leq N .$$

Since  $a^{(t)}$  is an increasing and bounded sequence it is convergent.  $\square$

The following theorem states that the pheromone values are convergent over time, and, as a result, the algorithm is convergent to a unique distribution.

**Theorem 6.** *For any point  $x_i$ , its pheromone value  $f_i^{(t)}$  is convergent as  $t \rightarrow \infty$ .*

*Proof.* Since we know that the pheromone values  $f_i^{(t)}$  are bounded from below and above for all  $t$ , it is enough to show that  $\{f_i^{(t)}\}_{t=1}^{\infty}$  is a monotonic sequence. From Eq. (30) we just need to show that  $w_i - a^{(t)}$  becomes only positive or only negative after some round  $T$ . Let us define the sequence  $\{b^{(t)}\}$  as  $b^{(t)} = w_i - a^{(t)}$ . From lemma 5 we know sequence  $a^{(t)}$  is increasing and  $a^{(t)} \rightarrow a$  (as  $t \rightarrow \infty$ ). Thus,  $b^{(t)}$  will be a decreasing sequence and convergent to  $w_i - a$ :

- If  $w_i \geq a$ , then  $b^{(t)} \geq 0$  for any  $t$ . Thus, from Eq. (30) we get  $\Delta^{(t)} f_i \geq 0$  and then  $f_i^{(t)}$  is a monotonically increasing sequence.

- If  $w_i < a$  then there exists  $T \in \mathbb{N}$  such that

$$\forall t > T : a^{(t)} > w_i \implies b^{(t)} < 0 ,$$

so  $f_i^{(t)}$  is a decreasing sequence for  $t > T$ .

□

**Corollary 6.1.** *From the convergence of pheromone values and the updating rule in Eq.*

*(25) it is clear that  $\pi_i^{(t)}$  is convergent (as  $t \rightarrow \infty$ ) for any  $i \in \{1, 2, \dots, N\}$ .*

After proving the convergence of the pheromone distribution and stationary distribution, we show in the following that the pheromone distribution (coming from the algorithm) converges to the stationary distribution.

**Theorem 7.** *The pheromone distribution  $\mathbf{f}^{(t)}$  converges to the stationary distribution  $\pi^{(t)}$  as  $t \rightarrow \infty$ , i.e. for any  $\epsilon > 0$ , there exists  $T \in \mathbb{N}$  such that*

$$\forall i \in \{1, \dots, N\}, \forall t > T \rightarrow |\pi_i^{(t)} - f_i^{(t)}| < \epsilon .$$

*Proof.* Theorem 6 states that the pheromone distribution converges (i.e.  $\forall i f_i^{(t)} \rightarrow f_i$  as  $t \rightarrow \infty$ ). Therefore, the update rule can be written as follows:

$$f_i = (1 - \rho)f_i + \rho \cdot \lim_{t \rightarrow \infty} \pi_i^{(t)} \implies \lim_{t \rightarrow \infty} \pi_i^{(t)} = f_i .$$

Thus, the pheromone and stationary distribution converge towards the same distribution.

□

Theorems 6 and 7 indicate two important points: first, that the pheromone distribution is convergent and second, that  $\mathbf{f}^{(t)}$  converges to a well-known distribution called the stationary distribution.



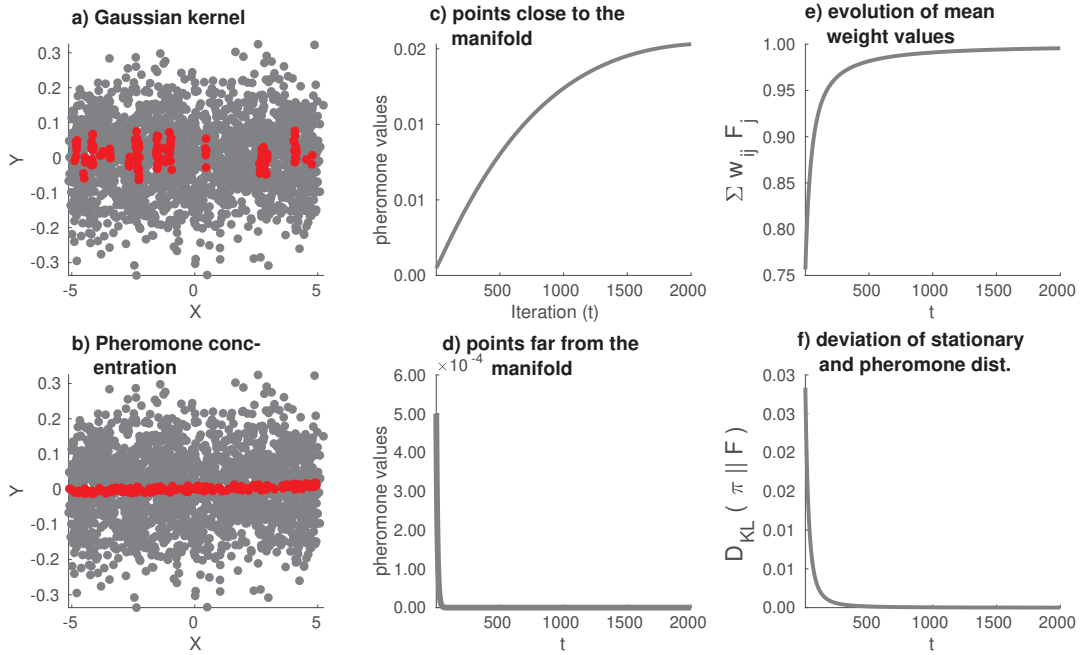


Figure 3: 1D manifold consisting of 2000 data points disrupted by Gaussian noise ( $N(0, 0.1)$ ) depicted as gray points in panel a and b. The 120 points with highest density (using a Gaussian kernel) and pheromone values (resulting from Eq. (20)) are highlighted in red in a) and b) respectively. The development of the pheromone values over iterations  $t$  for points close to the manifold is steeply increasing (as shown in panel c), while it is abruptly dropping to zero for points far away (see panel d). Panel e) shows the evolution of the mean value of weights  $a^{(t)}$  (see Eq. 29) and f) the KL- divergence from the stationary to the pheromone distribution over the iterations  $t$ .

In order to demonstrate the above results with an example we generate a 1D manifold such that 2000 data points are uniformly distributed on a line and disrupted by

Gaussian noise  $N(0, 0.1)$ <sup>7</sup> (gray points as shown in Figure 3a and b). In panel (a) the Gaussian kernel with  $\sigma = 1$  and  $r = 0.1$ <sup>8</sup> is used to detect 120 samples with the highest density highlighted in red. Panel (b) shows the samples with the highest pheromone values using the weights defined in Eq. (20). With the Gaussian kernel failing we clearly see the superiority of the new algorithm in recovering the linear manifold nearly perfectly, since it uses PCA in the creation of the MC. Panels (c) and (d) depict the evolution of pheromone values of samples close to and far from the underlying 1D manifold, respectively. Closer samples to the manifold exhibiting lower noise levels receive more and more visits in smaller time intervals (from theorem 3) and thus accumulate more and more pheromone over time. Panel (e) displays that the mean of weight values  $a^{(t)} = \sum_k w_k f_k$  is monotonic increasing and convergent, as explained in Lemma 5. Finally, panel (f) uses the Kullback Leibler (KL) divergence to show that the pheromone distribution converges towards the stationary distribution, confirming theorem 7.

### 5.2.2 Unbounded number of ants, fixed bounded path length

Since our algorithm allows employing multiple ants we study its performance in the presence of a swarm of ants, which walk for a limited number of steps  $n$  in every round.

---

<sup>7</sup>Other noise models, such as uniform noise (“tubes” around the manifold) or the Laplace distribution (strongly concentrating points along the manifold) can be alternative scenarios dependent on the process with which the manifolds are sampled. However, here we concentrate on Gaussian noise that is suitable for our applications to keep the paper more concise.

<sup>8</sup>Note that we select  $r$  such that the red points are more distributed along the manifold. For bigger value for  $r$ , the red points are more concentrated around  $x = 0$  and fail to recover the 1D manifold.

Let  $M$  denote the number of ants and  $N_j(n)$  represents the number of times a random walker visits point  $\mathbf{x}_j$  during  $n$  steps on  $(t + 1)$ -th round. For a specific case  $M = 1$  we have (according to Definition 2.1):

$$E[N_j(n)] = E_{P(X_0=\mathbf{x}_i)}[m_{ij}(n)] \quad (31)$$

and assuming uniform distribution over the initial states (sample points),

$$E[N_j(n)] = \frac{1}{N} \sum_{i=1}^N m_{ij}(n) . \quad (32)$$

For  $M > 1$  we still have the above equality for each ant since pheromone values are kept fixed over one round. Hence, in each round an ant walks independently and others do not have any impact on it. Initializing each ant on one of the sample points with probability  $P(X_0 = \mathbf{x}_i)$  we obtain by the Law of Large Numbers:

$$\frac{\sum_{a=1}^M N_j^a(n)}{M} \rightarrow E[N_j(n)], \quad \text{as } M \rightarrow \infty . \quad (33)$$

As the number of ants  $M$  goes to infinity and from Eq. (32) and (33) the updating rule, in Eq. (23), can be re-written:

$$f_j^{(t+1)} = (1 - \rho)f_j^{(t)} + \frac{\rho}{N} \sum_{i=1}^N \frac{m_{ij}(n)}{n + 1} . \quad (34)$$

The following corollary says that, like the weight values, the pheromone values form a decreasing sequence.

**Corollary 7.1.** *For each  $t$  the sequence of pheromone values  $\{f_i^{(t)}\}_{i=1}^N$  is monotonically decreasing.*

*Proof.* We use induction on  $t$  to show that the following statement holds for any arbitrary pair of points  $\mathbf{x}_i$  and  $\mathbf{x}_j$  with  $i < j$ :

$$f_i^{(t)} \geq f_j^{(t)} . \quad (35)$$

For simplicity we drop the exponent  $t$  in the following computations. In the first round  $t = 1$  we have:

$$f_i^{(0)} = f_j^{(0)} = \frac{1}{N} \xrightarrow{w_i \geq w_j} p_{ki} \geq p_{kj} \quad \forall k ,$$

and

$$p_{ki}^{(2)} = \sum_l p_{kl} p_{li} \geq \sum_l p_{kl} p_{lj} = p_{kj}^{(2)} ,$$

where  $p_{ki}^{(2)}$  is the  $(k, i)$ -the element of the second power of the transition matrix  $P$ .

Similarly, it can be generalized for any  $q$ :

$$p_{ki}^{(q)} \geq p_{kj}^{(q)} \quad \forall k .$$

As a result of theorem 2, we derive:

$$m_{ki}(n) \geq m_{kj}(n) \quad \forall k .$$

Thus, the updating rule (see Eq. (34)) shows that the inequality (35) holds in the first round. Now, let the statement hold up to  $t$ -th round (i.e.  $f_i^{(t)} \geq f_j^{(t)}$ ). In the  $(t + 1)$ -th round we have:

$$p_{ki} = \frac{w_i f_i}{\sum w_k f_k} \xrightarrow{w_i \geq w_j} \frac{w_j f_j}{\sum w_k f_k} = p_{kj} .$$

Similar to the first round, we obtain for any  $q$ :

$$\forall k \quad p_{ki}^{(q)} \geq p_{kj}^{(q)} \implies m_{ki}(n) \geq m_{kj}(n) \implies f_i^{(t+1)} \geq f_j^{(t+1)}$$

□

Therefore, the pheromone values, similar to weights, encode information about how close the individual sample points are to the subspace  $\hat{\mathcal{M}}$ . Analogously to Theorems

6 and 7 the following theorem guarantees the convergence of the pheromone and the stationary distributions.

**Theorem 8.** *The following statements hold for the pheromone distribution:*

(a) *For any point  $x_i$  its pheromone values  $f_i^{(t)}$  converge in time  $t$  (i.e.  $f_i^{(t)} \rightarrow f_i$  as  $t \rightarrow \infty$ ).*

(b) *The pheromone distribution  $f^{(t)}$  converges to the stationary distribution  $\pi^{(t)}$  as  $t \rightarrow \infty$ , i.e. for any  $\epsilon > 0$  there exists  $T \in \mathbb{N}$  such that:*

$$\forall i \in \{1, \dots, N\}, \forall t > T \rightarrow |\pi_i^{(t)} - f_i^{(t)}| < \epsilon .$$

*Proof.* The proof of (a) follows that of theorem 6.

b) From the convergence of pheromone values and the updating rule (34) we can write:

$$f_i = \frac{1}{N(n+1)} \sum_{i=1}^N \lim_{t \rightarrow \infty} m_{ji}^{(t)}(n) .$$

From (26) we know :

$$\forall i, j; \quad p_{ji}^{(t)} = \pi_i^{(t)} \implies p_{ji}^{(t)(q)} = \pi_i^{(t)} \quad \forall q , \quad (36)$$

where  $p_{ji}^{(t)(q)}$  is the  $(j, i)$ -th element of  $(P^{(t)})^q$  in the  $t$ -th round. Therefore, we have

$$\begin{aligned} f_i &= \lim_{t \rightarrow \infty} \left( \frac{1}{N(n+1)} \sum_{i=1}^N m_{ji}^{(t)}(n) \right) \stackrel{\text{theo. (2)}}{=} \lim_{t \rightarrow \infty} \left( \frac{1}{N(n+1)} \sum_{i=1}^N \sum_{q=0}^n p_{ji}^{(t)(q)} \right) \\ &\stackrel{\text{Eq. (36)}}{=} \lim_{t \rightarrow \infty} \frac{1}{N(n+1)} \sum_{i=1}^N \sum_{q=0}^n \pi_i^{(t)} = \lim_{t \rightarrow \infty} \pi_i^{(t)} . \end{aligned}$$

□

Hence, it can be said that both scenarios ( $n \rightarrow \infty$  or  $M \rightarrow \infty$ ) result in the same pheromone structure and asymptotic:

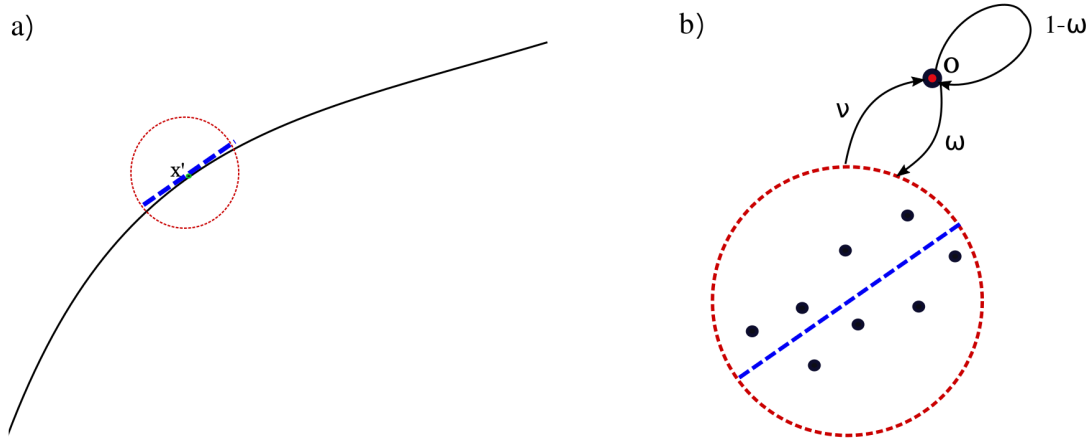


Figure 4: a) approximation of a piece of non-linear manifolds with a linear patch. b) the linear patch with an external node ( $o$ ), simulating points outside the neighborhood.

- pheromone values are sorted according to their closeness to the recovered linear manifold, and finally
- they are convergent to the stationary distribution of the Markov Chain.

### 5.3 Pheromone distribution for a non-linear manifold

For non-linear manifolds we consider a simplified scenario for our analysis. Although it is slightly different from the algorithm's definition, we use it as an example to explain the success of M3A in recovering non-linear manifolds. Motivated by the assumption that a non-linear manifold can locally be approximated by linear patches (see Fig. 4 a), we focus on a small part of the manifold (red circle), and we model all points outside via a single state called  $o$ . Therefore, we assume that the random walker can jump from any point inside the red circle to  $o$  with a probability proportional to the fixed-number  $\nu$  and vice versa with a probability proportional to the fixed-number  $\omega$  (see Fig. 4 b). Consequently, the transition probabilities are defined as:

$$p_{\mathbf{o}j}^{(t+1)} = \begin{cases} \frac{\omega}{N}, & \text{if } j \neq \mathbf{o} \\ 1 - \omega & \text{if } j = \mathbf{o} \end{cases} \quad \text{and} \quad p_{i(\neq \mathbf{o})j}^{(t+1)} = \begin{cases} \frac{\nu}{\sum_{k \neq \mathbf{o}} w_k f_k^{(t)} + \nu}, & \text{if } j = \mathbf{o} \\ \frac{w_j f_j^{(t)}}{\sum_{k \neq \mathbf{o}} w_k f_k^{(t)} + \nu}, & \text{if } j \neq \mathbf{o} \end{cases},$$

where  $N$  is the number of data points inside the red circle. Note that in order to differentiate the transition probabilities and the stationary distribution in this section from the previous one, we use the prime symbol ( $'$ ).

From the balance equation we can compute the stationary distribution. Thus, for any state  $j$  ( $\neq \mathbf{o}$ ) we have:

$$\begin{aligned} \pi_j^{(t+1)} &= \sum_{i \neq \mathbf{o}} \pi_i^{(t+1)} p_{ij}^{(t+1)} + \pi_{\mathbf{o}}^{(t+1)} p_{\mathbf{o}j}^{(t+1)} \\ &= \sum_{i \neq \mathbf{o}} \pi_i^{(t+1)} \frac{w_j f_j^{(t)}}{\sum_k w_k f_k^{(t)} + \nu} + \pi_{\mathbf{o}}^{(t+1)} \frac{\omega}{N} \\ &= \frac{w_j f_j^{(t)}}{\sum_k w_k f_k^{(t)} + \nu} \sum_{i \neq \mathbf{o}} \pi_i^{(t+1)} + \pi_{\mathbf{o}}^{(t+1)} \frac{\omega}{N} \\ &= \left( \frac{w_j f_j^{(t)}}{\sum_k w_k f_k^{(t)}} \right) \left( \frac{\sum_k w_k f_k^{(t)}}{\sum_k w_k f_k^{(t)} + \nu} \right) (1 - \pi_{\mathbf{o}}^{(t+1)}) + \pi_{\mathbf{o}}^{(t+1)} \frac{\omega}{N}. \quad (37) \end{aligned}$$

In the previous section we found  $p_{ij}^{(t+1)} = \frac{w_j f_j^{(t)}}{\sum_k w_k f_k^{(t)}}$  and hence we can re-write the terms above as follows:

$$\begin{aligned} \pi_j^{(t+1)} &= p_{ij}^{(t+1)} \left( \frac{\sum_k w_k f_k^{(t)}}{\sum_k w_k f_k^{(t)} + \nu} \right) (1 - \pi_{\mathbf{o}}^{(t+1)}) + \pi_{\mathbf{o}}^{(t+1)} \frac{\omega}{N} \\ &= a p_{ij}^{(t+1)} + b, \quad \text{where} \quad (38) \\ a &= \left( \frac{\sum_k w_k f_k^{(t)}}{\sum_k w_k f_k^{(t)} + \nu} \right) (1 - \pi_{\mathbf{o}}^{(t+1)}) > 0 \\ \text{and } b &= \pi_{\mathbf{o}}^{(t+1)} \frac{\omega}{N} > 0. \end{aligned}$$

From Eq. (26) we obtain:

$$\pi_j^{(t+1)} = a \pi_j^{(t+1)} + b, \quad (39)$$

where  $\pi_j^{(t+1)}$  denotes the stationary distribution of the random walk in the previous section. Thus, this equation connects the stationary distributions of the new MC  $\pi'$  to  $\pi$  in the previous section. In other words, the new distribution is monotonically increasing with respect to the previous one.

The theorem 3 connects the long-term behavior of a random walk to its stationary distribution. From (27) and (39) we see that the stationary distribution (and similarly the pheromone distribution) is sorted according to the distances of data points (inside the neighborhood) to the linear patch constructed by PCA. From theorem 2, we know that the occupancy time (i.e. the short-term behavior) is directly connected to the transition probabilities of the MC. For any  $i$  and  $j$  ( $\neq \mathbf{o}$ ), we have:

$$p_{ij}^{(t+1)} = \frac{w_j f_j}{\sum_k w_k f_k + \nu} = \left( \frac{w_j f_j}{\sum_k w_k f_k} \right) \left( \frac{\sum_k w_k f_k}{\sum_k w_k f_k + \nu} \right) = ap_{ij}^{(t+1)}$$

Thus, the new transition probabilities are monotonically increasing function of the previous transition probabilities. Therefore, similar results as in section 5.2.2 hold.

By simplifying the scenario, we show in this section that the algorithm sorts data points according to their distances to the linear patches (as estimators for a manifold). From (Little et al., 2017), we can conclude that the performance of the new algorithm depends on several factors including selecting an appropriate neighborhood size  $r$ . Note that  $r$  should be big enough to include enough number of samples for recovering the tangent space, and small enough to prevent to cover high-curvature. Therefore, in the presence of an appropriate  $r$  value, our analysis shows that the pheromone distribution contains valuable information in order to highlight non-linear manifolds.



## 6 Experiments

In this section we use synthetic and real world data sets to investigate the performance of M3A in different scenarios. At first we discuss the influence and specification of its hyper-parameters. Furthermore, complementary to our analysis in section 5.1 we empirically examine the impact of the new MC formulations for denoising and visualization of non-linear manifolds. Moreover, we show how the proposed algorithm helps to discover manifolds and to build better probabilistic models in the sense of sparseness, descriptiveness, and preservation of structural details.

### 6.1 Strategy for hyper-parameter selection

Our analysis in section 5 shows that the neighborhood size  $r$ , the number of ants  $M$ , and the number of steps  $n$  in one round play a major role for the performance of M3A. Thus, an automatic strategy to find appropriate values for these quantities is highly desirable. In the following we present such strategies for practical application.

In subsection 5.1, we demonstrate that a proper value for  $r$  depends on the noise level  $\sigma$  of the manifold. We observe that a bigger  $r$  helps to reduce the effect of the noise and achieve a better approximation for the linear manifold. However, as it is shown in (Kaslovsky and Meyer, 2014) and (Little et al., 2017) for non-linear manifolds, the presence of curvature encourages to use a smaller radius  $r$ , such that the manifold looks almost linear within the local neighborhood. Hence, the selection of  $r$  is a trade-off between the noise level and curvature of the manifolds. Here, to specify a suitable value for  $r$  we use a mixture of Gaussian models to assess quantitatively the recovery of the manifold. First, we place the centers of Gaussian distributions on points with higher

pheromone values, and then their co-variance matrices are computed based on the local neighbors. Finally, we use the average log-likelihood (ALL) function to evaluate how good models (for each  $r$ ) fit the data set and to pick the best  $r$  (see Figure 5c).

In subsection 5.2 we investigate two strategies: a) using a few ants walking many steps in every round, and b) using many ants walking a limited number of steps. As a rule of thumb to determine a suitable number of steps and ants ( $n$  and  $M$ ) we assume they fulfill the following inequality:

$$n \cdot M \geq z \cdot N \quad (40)$$

for a constant  $z$ . Thus, one may use a small number of workers  $M$  (or processors/cores) with a big number of steps  $n$  or use many workers with a small number of steps in each round. In our experimentation, we set the number of steps and the number of ants to  $n = N$  and  $M = z = 50$ , respectively.

The pheromone determines how much the ants are reinforced to prefer points frequently visited before. And the evaporation rate  $\rho$  and parameter  $\gamma$  control the impact of the pheromone on the random walk process, and thus their determination is connected to  $M$  and  $n$ . If  $M$  and  $n$  are both small, it means the algorithm has neither enough ants nor enough time to highlight manifolds; thus, our analysis in subsection 5.2 is not valid any more. In this case, we suggest using smaller  $\rho$  and  $\gamma$  to prevent highlighting background noise. In this contribution, we set them to  $\rho = \gamma = 0.1$  for all the experiments.

Furthermore, the maximum amount of pheromone  $c$  added on a sample per round is fixed to 2, and the factor  $\alpha$  for the weight  $w_d$ , in Eq. (11), is determined such that only  $p = 50\%$  of neighbors have non-zero values. Finally, for any sample  $x_i$ , we

approximate the point  $\boldsymbol{x}'$  that touches the underlying manifold by  $\boldsymbol{x}' = \frac{1}{|\mathcal{N}_i|} \sum_{\boldsymbol{x}_k \in \mathcal{N}_i} \boldsymbol{x}_k$ , (see Eq. (10)).

## 6.2 Homogeneous Markov Chain

In this section we extend our theoretical analysis with an empirical analysis on non-linear manifolds. First, we demonstrate the denoising capability of the MCs defined in Eq. (11) and (12) on highlighting the underlying manifold. Second, we investigate their impact on the performance of t-distributed stochastic neighbor embedding (t-SNE) (Van der Maaten and Hinton, 2008), a widely used dimension reduction and visualization technique.

**Denoising the manifold:** To check the success of the weighting methods (section 3) in keeping ants close to the manifold, we need to study the short and long-term behavior of the random walk defined by Eq. (2). According to theorems 2 and 3 the short and long-term behaviors of MCs are directly connected to the power series of the transition probability and to the stationary distribution, respectively. Here, we use a synthetic data set to investigate the constructed MCs.

In Figure 5, panels (a,b) show 6000 samples generated from the s-curve manifold disrupted by Gaussian noise  $N(0, 0.2)$ . In order to select  $r$  we follow the strategy introduced before (in section 6.1). We apply M3A with different  $r$  using 20% of the samples (1200 points) with the highest pheromone values as means for Gaussian distributions and their neighbors to determine the co-variance matrices. The quality of the models is computed by the average log-likelihood (ALL) function and the results are shown in

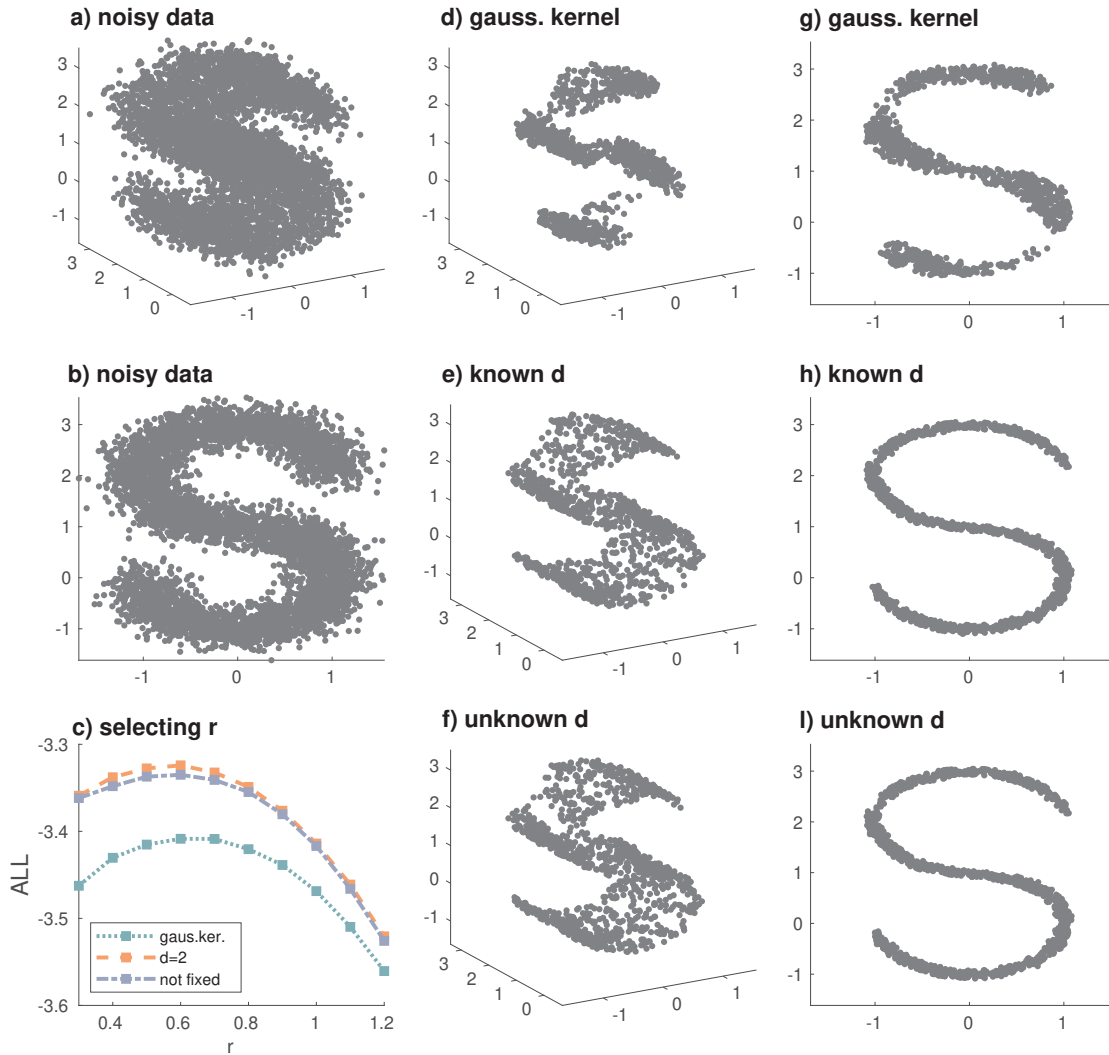


Figure 5: 2 views of the noisy s-curve data set (a,b), radius selection via average log likelihood for different  $r$  (c), and 1200 points with highest stationary distributions values for the MCs defined by Eq. (11) (e,h) and Eq. (12) (f,l), using  $r = 0.6$ .

Figure 5c. Panels (d-l) show 1200 points with the highest stationary distribution values for different MCs. As it can be seen the MC based on the Gaussian kernel function fails to recover the manifold (see panels d and g). However, the other methods are successful in keeping random walkers close to the manifold (see panels e, f, h and l). Note that panels (e, f) also show that, since we use our knowledge about the intrinsic dimension-

ality of the manifold in Eq. (11), it gives a slightly better result than the MC defined by Eq. (12). A similar result is achieved for the short-term behavior of the MCs where the two formulation of Eq. (11) and (12) outperforms the Gaussian kernel function.

**Manifold aligned similarities in t-SNE:** A specific case of dimensionality reduction is the visualization where high-dimensional data is embedded in two or three dimensions. A well-known non-linear tool to visualize high-dimensional data is t-distributed stochastic neighbor embedding (t-SNE) (Van der Maaten and Hinton, 2008). It transforms pairwise similarities of data points to probability distributions of the high-dimensional data as well as the low-dimensional embedding and then minimizes the Kullback-Leibler divergence between them. While the original t-SNE is based on the Gaussian kernel function Eq. (9), we replace it with the tangent space dissimilarity measure defined in Eq. (11) and (12). Since high-dimensional data spaces are typically sparse we use a k-nearest neighbors instead of radius neighborhood, and we set  $p = 100\%$  (i.e., we normalize the distances in Eq. (11) via dividing them by their maximum values). Furthermore, we compare the results to the original t-SNE using the Euclidean distance to compute the neighborhood probabilities, and a t-SNE version using the Mahalanobis distance instead. The latter is based on the co-variance of the data and can therefore align with global directions of major variance. Therefore, the locally aligned formulation, in Eq. (11) and (12), is compared to “no alignment” and “global alignment”.

We demonstrate the strategy using two real non-linear dimensionality reduction benchmark data sets: 1) COIL20: consists of images (with  $32 \times 32 = 1024$  pixels) of 20 objects rotated 72 times (5 degrees per image) forming one dimensional manifolds, and

2) USPS: contains 9298 gray-scale images (with  $16 \times 16$  pixels) of handwritten digits. As a preprocessing step we apply PCA to decrease the dimensionality to 20. The different t-SNE outputs are compared and evaluated using the quality of group compactness (QGC) (Gorban and Zinovyev, 2010) as a measure of how close samples from the same class remain in the embedding space. Let  $c(i; k)$  denote the number of points in the  $k$ -neighborhood of  $\mathbf{x}_i$  with the same label, then the compactness of a class is defined as:

$$QGC_k(l) = \frac{1}{k \cdot N(l)} \sum_{y_i=l} c(i; k)$$

where  $N(l)$  is the number of samples with the label  $l$ . If the QGC for a class is close to one, it means the class is compact and well-separated from others. Since t-SNE is non-convex it generates different results for each run and hence we repeat the experiment 10 times and report the average quality.

Figure 6 shows example t-SNE embeddings for COIL20 and USPS (panel a and c), accompanied by results using Eq. (11) with  $d = 1$  and  $k = 10$  and Eq. (12) with  $k = 20$  (panels b and d). Panels e and g depict the average QGC curves for the different similarity measures, respectively. It can be seen that the new similarity measures outperform others, especially for COIL20 since we use our knowledge about the intrinsic dimensionality  $d = 1$  of the manifolds. For the USPS, it is clear that it does a better job for more global structures, i.e. clusters visible for  $k > 200$ . Since the new weights only consider the distance to the estimated manifold, it can recover more closed loops with less distortion (see panels a and b). Moreover, the four objects corresponding to four types of toy cars cannot be separated using the original t-SNE because it maps cars with the same orientations close to each other. As seen in panel g, the new weights do

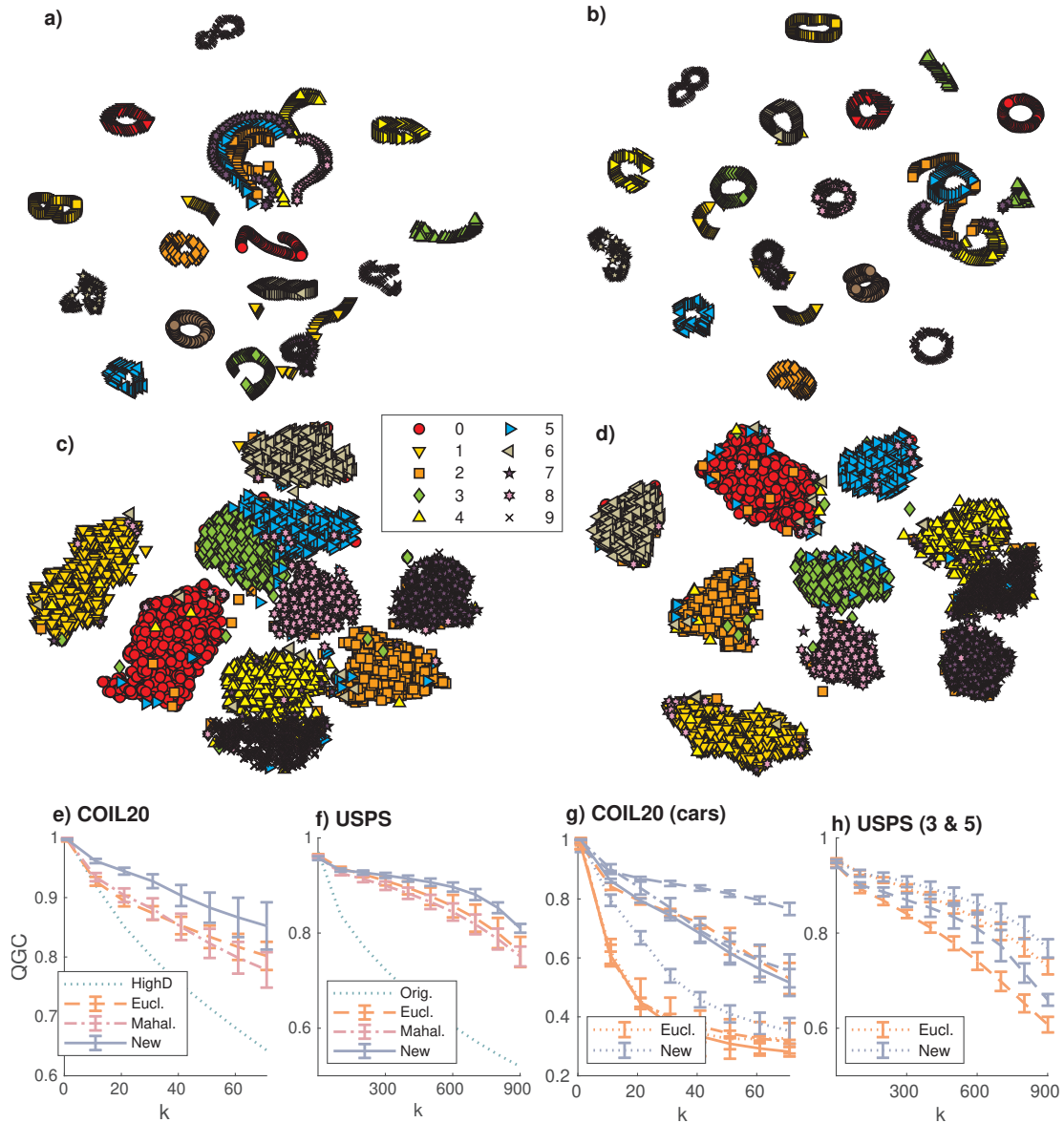


Figure 6: The impact of the new similarity measures on the t-SNE: Top row COIL20: a) using original t-SNE (perpl. 20) and b) tSNE+ using Eq. (11) ( $d = 1, k = 10$ ). Middle row USPS: c) using original t-SNE (perpl. = 30) and d) tSNE+ using Eq. (12) ( $k = 20$ ). Bottom row: average QGC measures for COIL20 and USPS (panel e and f). g) QGC curves for four classes of toy cars in COIL20 and h) of digits 3 and 5 in USPS.

a better job in embedding them. Not knowing the intrinsic dimensionality of manifolds in the handwritten digit data set, we use Eq. (12) to define a similarity measure for t-SNE. Although t-SNE with the Gaussian kernel function can recover most classes, a strong overlap between samples from digit 3 and 5 is visible. The new weights however, successfully separate them and increase the QGC measure, as seen in panels d and h.

### 6.3 Manifold Alignment Aware Ants (M3A)

In this section we analyze the capability of M3A in recovering structures surrounded by background noise. First, we compare it to related state-of-the-art techniques for denoising low-dimensional manifolds. And second, we demonstrate the performance of our algorithm on a real-world astronomical data set, namely the GAIA DR2 catalog (Collaboration et al., 2018), to extract stellar structures.

**Extracting low-dimensional manifolds:** Typically, manifold learning assumes that data points are lying on a low-dimensional manifold embedded in a higher-dimensional space. However, in practice, manifolds are disrupted by high-dimensional noise, for example stemming from equipment or the presence of samples that do not belong to the manifolds. In those cases, the goal is to extract the manifolds by removing the background noise and simultaneously suppressing the noise level on the manifolds. To compare the performance of denoising techniques we create a synthetic 1D manifold forming a circle ( $r = 6$ ), which misses a part of its arc. Furthermore, we randomly select 3000 points from the manifold and add Gaussian noise  $N(0, 0.3)$ . To model background noise, 3000 samples are generated following a uniform distribution in the



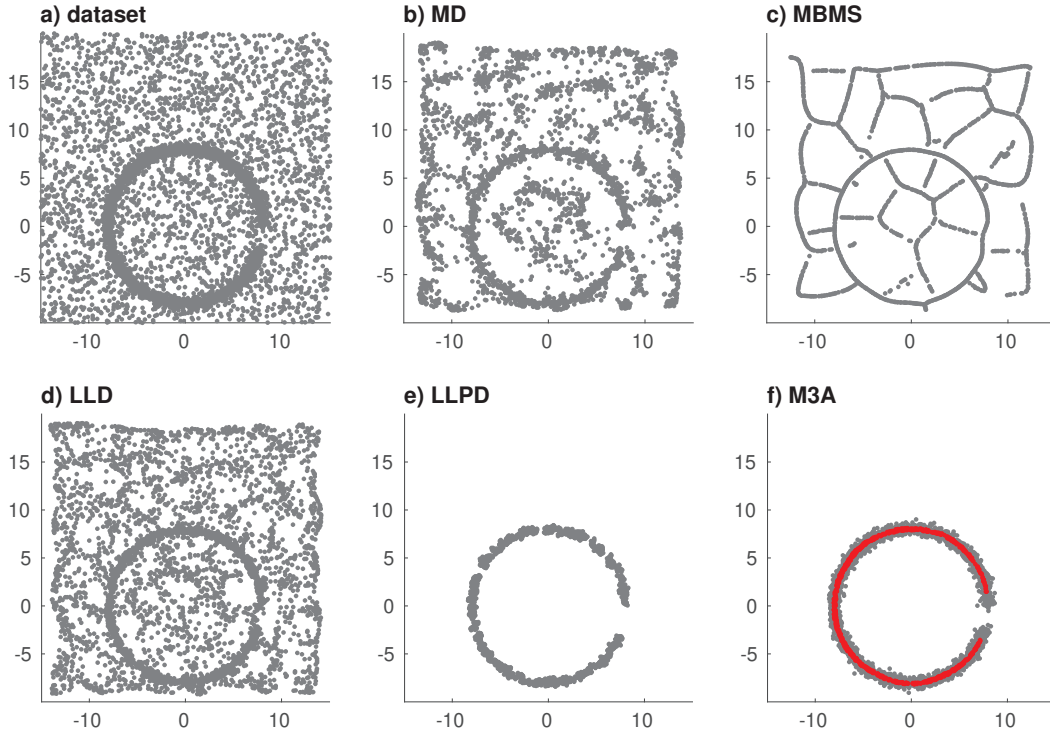


Figure 7: Panel a): 1D circular manifold embedded in uniform background noise. Best results of: b) MD, c) MBMS, d) LLD, and e) LLPD (containing 2000 samples, i.e. 35%). Panel f): denoised M3A result where gray and red points show 2700 (45%) and 780 (%13) samples with  $r = 3$  and  $r = 2$ .

square  $[-15, 15] \times [-10, 20]$  (see Figure 7a).

In the following we compare four denoising techniques to the proposed M3A. We use: a) MD with hyper-parameters settings  $\delta t = 0.1$ ,  $N_{\text{iter}} = 20$  and  $k = 40$ , b) MBSM with  $k = 40$  and  $N_{\text{iter}} = 40$ , and c) LLD with  $\lambda = 0.01$ ,  $k = 40$  and  $\sigma = 1$ . Lastly, LLPD is deployed using a threshold of 0.15. Moreover, in order to select  $r$  in M3A we follow a two step strategy. First, we apply the algorithm using  $r = 3$  to remove the background noise. Note that since the density of the background is lower than the manifold it can be done using a wide range of  $r$ . Then, we use the strategy explained in

section 6.1 to find the best  $r$  for the remaining points.

Some techniques, such as MD, MBMS and LLD, are not designed to deal with the presence of the background noise and are therefore not directly suitable for the given situation, as shown in Figure 7. Since the Manifold Denoising (MD) algorithm pushes samples towards dense regions it has two disadvantages: first, it creates some new dense structures in data space, and second, the real manifold becomes discontinuous (see panel b). Since Manifold Blurring Mean Shift (MBMS) limits the movement of samples to be parallel to the manifold normals, it can prevent discontinuity on the manifold. However, it creates many artificial 1D manifolds and fails to recover the missing arc in the circle structure of interest (see panel c). Similar to MBMS, the presence of background noise highly influences the result of the LLD algorithm, which also fails to recover the manifold (see panel d). In contrast to previous methods, LLPD considers the presence of the background noise and extracts the samples in dense regions recovering the noisy circle. Although it does not create new structures, it can not denoise the manifold. It also fragments the structure into smaller clusters (see panel e). On the other hand, our method encourages random walkers to stay close to the manifold and overcomes all three problems: 1) it recovers the circle from the background noise, 2) it preserves the missing part in the arc without creating artificial structures, and 3) it reduces the manifold noise by subsampling based on the highest pheromone level (see panel f).

**Detecting dense structures in real world data:** An important task in astronomy is to extract stellar structures, such as galaxies and globular clusters. Here, we apply M3A

to a part of the GAIA DR2 catalog (Collaboration et al., 2018) containing 1,071,714 light sources. The data set contains seven features, including five denoting positional information (right ascension (ra), declination (dec), motion along ra and dec, and parallax), as well as photometric information, that is G-band magnitude and (B-R) color. As a preprocessing step, we normalize the data set and select the  $k$  nearest neighbor (with  $k = 20$ ) to speed up the computation. First, we use the positional information to detect stellar structures. Figure 8a shows that M3A detects five of the known globular clusters (GC): 1) NGC4147 2) NGC5024 3) NGC5053 4) NGC5272 5) NGC5466. Second, we would like to cluster the members of a GC. To extract stars of a GC, we use objects which are not further than 0.05 degree from the center of the GC. Then, we use their photometric information since it reveals the stage of life a star is in (for more details see (Mohammadi et al., 2019)). Panels b and d display the clusters detected by spectral clustering using a Gaussian kernel (Von Luxburg, 2007). However, the resulting clusters do not correlate to any astrophysical meaning. If M3A is applied on the extracted light sources, we obtain the transition probability matrix  $P$ . After symmetrization  $P' = \frac{1}{2}(P + P^T)$ , the spectral clustering can be implemented based on  $P'$ . From panels c and e, we see that the resulting clusters coincide with the three known groups called: Main Sequence (MS), Horizontal Branch (HB) and Red Giant Branch (RGB) (see (Mohammadi et al., 2019)). It can be seen that MS is separable since it is a 2D manifold while HB and RGB are 1D manifolds. Moreover, since HB is almost a line perpendicular to RGB, the weight values between these two clusters are small, and as a result, they are separable. In summary, M3A provides: 1) a pheromone distribution that can be used to find stellar structures, and 2) a transition probability matrix that can

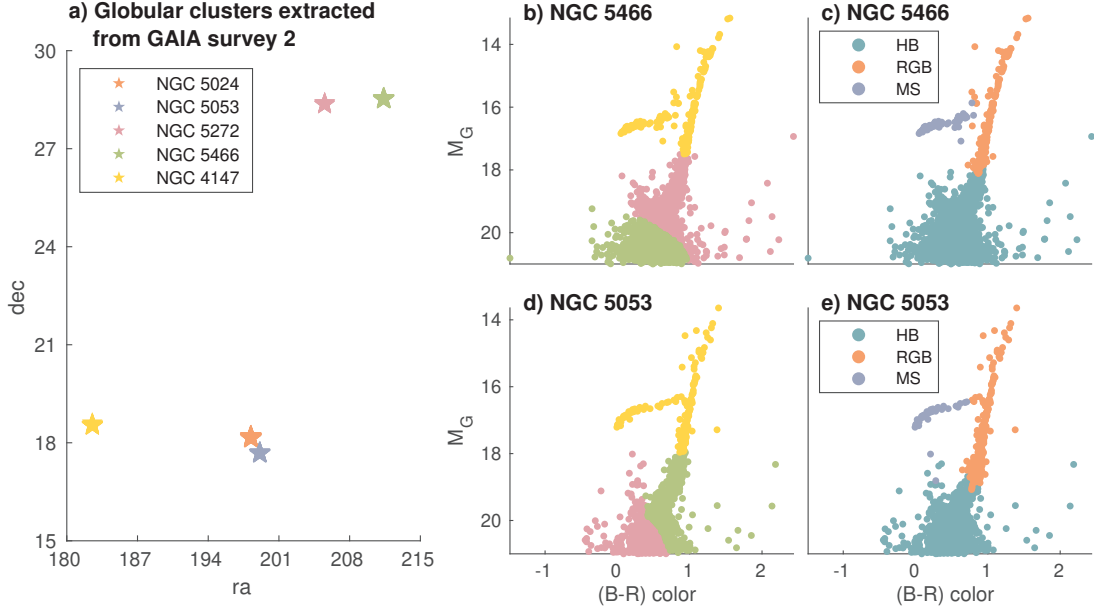


Figure 8: 5 clusters formed by 0.02% of samples with highest pheromone values resulting from M3A ( $r = 0.1$ ,  $N_{\text{steps}} = 5 \times 10^6$ ) applied to a window in the sky between right ascension ( $ra$ )  $\in ]180, 220[$  and declination ( $dec$ )  $]15, 30[$ . b-e) the application of spectral clustering on the photometric data of two globular clusters: b,d) using Gaussian kernel as a similarity measure. c,e) using the transition probability in Eq. (14) ( $r = 0.5$ ).

be used to group stars according to the stage of their life.

## 6.4 Improved density estimation

A common task in machine learning is to estimate the underlying probability density function (pdf) given a limited number of data points. A typical way to model it is to use a finite mixture model

$$\hat{f}(\mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^K p(k) \hat{f}(\mathbf{x}; \theta_k) ,$$

where  $p(k)$  and  $\theta_k$  are the mixture weight and the parameter values of the  $k$ -th component, respectively, while  $\boldsymbol{\theta}$  denotes the set of all parameters. If the number of models

( $K$ ) is much smaller than the number of samples ( $N$ ) it provides a sparse representation for the data set. A well-known algorithm is the Gaussian Mixture Model (GMM) (Bishop, 2006), which uses a mixture of full-rank Gaussian distributions to model the pdf. Its parameters are learned via optimizing the likelihood function through an iterative process, called Expectation-Maximization (EM). It is known that EM is sensitive to the initialization step, and it may get stuck in local optima. Moreover, the time needed for EM is dependent on the number of components and increases for large data sets with many samples and dimensions. One way to avoid this problem is to use a non-parametric method, such as Parzen Window (PW), which is a special case of the finite mixture model with  $K = N$ , and it only has one parameter  $\sigma$ :

$$\hat{f}(\mathbf{x}; \sigma) = \frac{1}{N} \sum_{k=1}^N \mathcal{K}(\mathbf{x}; \mathbf{x}_k, \sigma)$$

where  $\mathcal{K}$  is the Gaussian kernel function. However, if the data points are distributed along a low-dimensional manifold the spherical Gaussian kernel is often not an optimal choice (Vincent and Bengio, 2003). Therefore, (Vincent and Bengio, 2003) proposed the Manifold Parzen Window (MPW) which allows the Gaussian distributions to have elliptical shapes instead of spherical. Whereas PW and MPW may provide a reliable estimation, they both suffer from high computational costs in the test phase, due to the large number of Gaussian models. To prevent this problem several methods, such as Simplifying mixture models (SMM) (Zhang and Kwok, 2010) and Hierarchical clustering of a mixture model (HCMM) (Goldberger and Roweis, 2005) were proposed. They start with a large mixture model and then construct a simpler model, such that the distance between the original model and the simplified one is minimized. In contrast to these approaches, (Wang et al., 2009) proposed a new method called Fast Parzen

---

**Algorithm 3:** Partitioning the data space

---

**Result:** set  $\mathcal{S}$ , containing the center of balls ;

1 **Input:** The pheromone vector  $\mathbf{f} = [f_1, f_2, \dots, f_N]$  and the balls' radius  $r_{\text{ball}}$ ;

2 Set  $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ ,  $\mathcal{S} = \phi$ ;

3 **while**  $\mathcal{D} \neq \phi$  **do**

4      $l = \arg \max_{\mathbf{x}_i \in \mathcal{D}} f_i$  ;

5     find  $l$ 's neighbors  $\mathcal{N}_l = B_{r_{\text{ball}}}(\mathbf{x}_l)$  ;

6      $\mathcal{S} \leftarrow \mathcal{S} \cup \{l\}$  ;

7      $\mathcal{D} \leftarrow \mathcal{D} - \{l\} - \mathcal{N}_l$  ;

8 **end**

---

Window (FPW), which partitions the data space via hyper-balls positioned randomly with fixed radii  $r_{\text{ball}}$ . Then, it fits a full-rank Gaussian distribution for each ball only. Extending the latter idea we use the pheromone values obtained by M3A to find the best position for the center of the hyper-balls. We start with the whole data set  $\mathcal{D}$  and pick the point  $\mathbf{x}_l$  with the highest pheromone as the center of a ball. Then we remove all points within the ball from  $\mathcal{D}$ . Then, we continue picking the next point in  $\mathcal{D}$  with the highest pheromone as the new center and remove its neighbor from  $\mathcal{D}$ . Following this strategy, we obtain a more compact list of centers. Its pseudo-code is summarized in Algorithm 3.

To demonstrate the performance of the new strategy, we use three data sets:

- 1) a synthetic spiral shape manifold with Gaussian noise  $N(0, 0.04)$  for which 1000 samples are generated as a training set (see Figure 9a), and another 20000 samples

as a test set.

- 2) two intersected circular manifolds with radius 2: a) in the above circle the density varies and the noise level is 0.2; b) at the bottom the noise level varies between 0 to 0.3 (see Fig. 10a). We generated 3350 and 16750 points for training and testing, respectively. And
- 3) an astronomical simulation of a jellyfish galaxy that contains 3D position information of 58531 particles.

To evaluate the performance, we use the average log-likelihood (ALL) to measure how well the adapted model can describe the data set. We furthermore use 10-fold cross-validation on the training set to find the best hyper-parameter values.

Table 1 reports the results of the above algorithms on the spiral shape manifold. From cross-validation, hyper-parameters are found where  $m$ ,  $k$ , and  $\sigma$  are the number of components and neighbors, and the scale of Gaussian kernel, respectively, and  $r_{rw}$  is the neighbor size for M3A. Note that  $\text{GMM}^k$  is the Gaussian mixture model when k-means is used to initialize components. We repeat the experiment for M3A 10 times. While ALL for the M3A+FPW (i.e., algorithm 3) is comparable with PW, it provides a much sparser representation for the data set and speeds up the computation in the testing phase. Moreover, Figure 9 shows the contour curves of the pdfs. We observe that most methods represent the center of the spiral fairly well, with the exception of PW and SMM, which provide the most noisy model. Strikingly, PW, SMM, and HCMM discontinue the manifold in the less dense tail on the right side of the spiral, dividing it into several small clusters. Our strategy M3A+FPW, on the other hand, not only

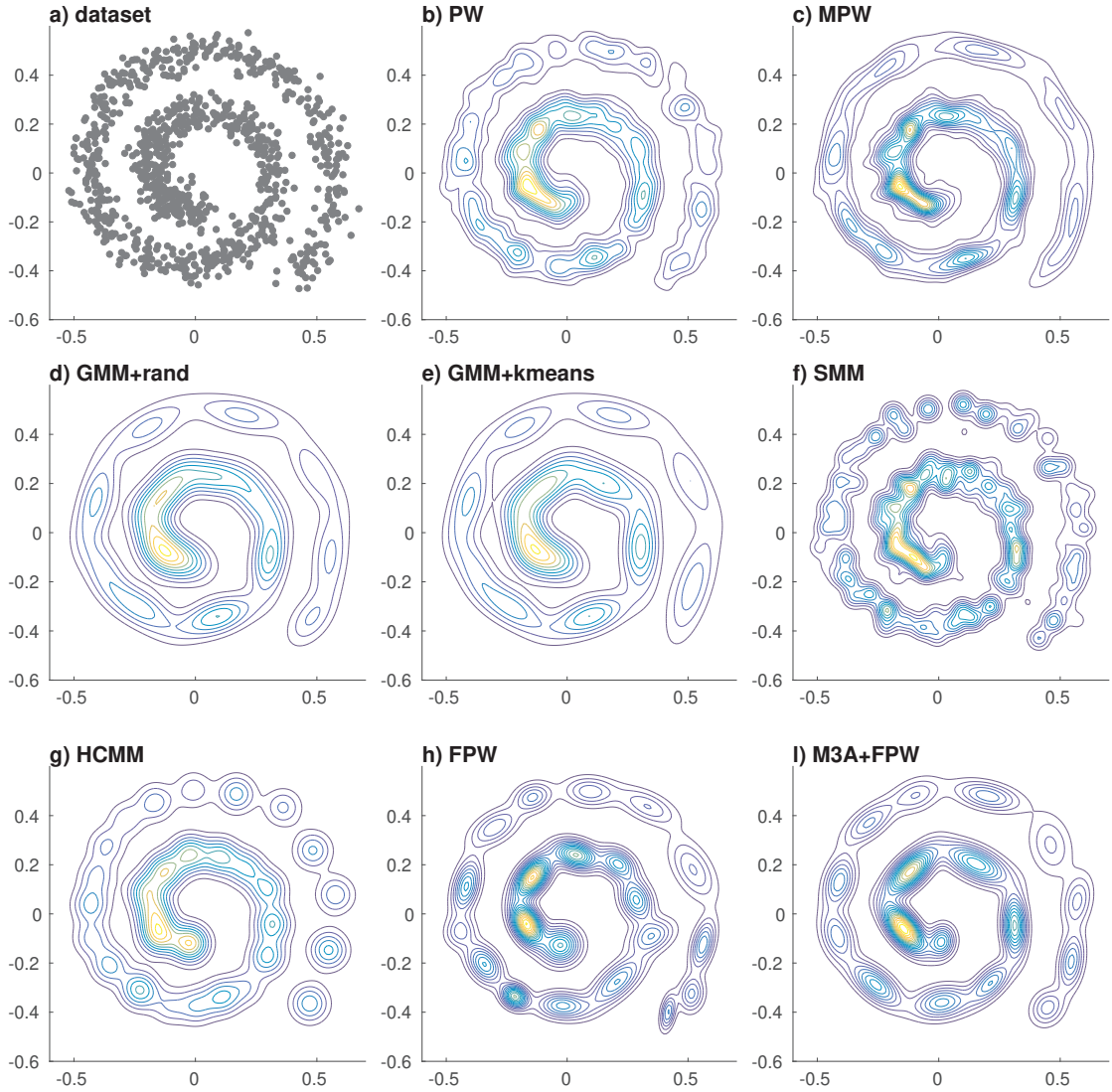


Figure 9: a) the noisy spiral data set, b-h) the contour curves for the density estimators: PW, MPW, GMM+rand, GMM+kmeans, SMM, HCMM, FPW, and i) M3A+FPW: neighborhood radius  $r_{rw} = 0.13$  and balls' size  $r_{ball} = 0.21$ .

provides a compact model but also successfully tracks the underlying structure and stays very close to the original manifold as compared to others.

In Table 2 we present the outputs of the studied algorithms on the two circular manifolds. as in the previous examples, cross-validation is used to tune hyper-parameters, and then the performance of the algorithms on the test set is reported. While the



Table 1: Comparison of density estimators for the spiral data set.

Algorithm	Parameters	ALL (std)	
PW	$\sigma = 0.03$	0.3130	(0.0)
MPW	$k = 45$	0.3003	(0.0)
FPW	$r_{\text{ball}} = 0.18$ ( $m = 20.8$ )	0.2851	(0.0092)
SMM	$m = 95$	0.2762	(0.0116)
HCMM	$m = 35, \sigma = 0.06$	0.2832	(0.0049)
GMM	$m = 20$	0.2944	(0.0028)
GMM <sup>k</sup>	$m = 15$	0.2802	(0.0049)
M3A+FPW	$r_{rw} = 0.13, r_{\text{ball}} = 0.21$ ( $m = 17.1 \pm 1$ )	0.3101	(0.0052)

FPW+M3A yields a sparse representation for the data, it also outperforms (in terms of ALL) others, except for the GMMs methods. However, since the M3A can be parallelized, it is more suitable for big datasets, in comparison to GMM techniques. Moreover, we display the contour curves of pdfs in Fig. 10 where the true manifold (without noise) is shown in red. It can be seen that the M3A helps FPW to be more compact and more successful in tracking the manifolds, as indicated by the black diamond in the plot. In comparison to GMMs, FPW+M3A uniformly distributes Gaussian models on the manifolds, even in less dense regions. It is especially important in some applications such as astronomy where low-density streams have high importance in studying the evolution of astronomical structures. Therefore, similarly to previous experiments, we conclude that FPW+M3A builds a compact model, which is more successful in following the manifolds, regardless of their densities.

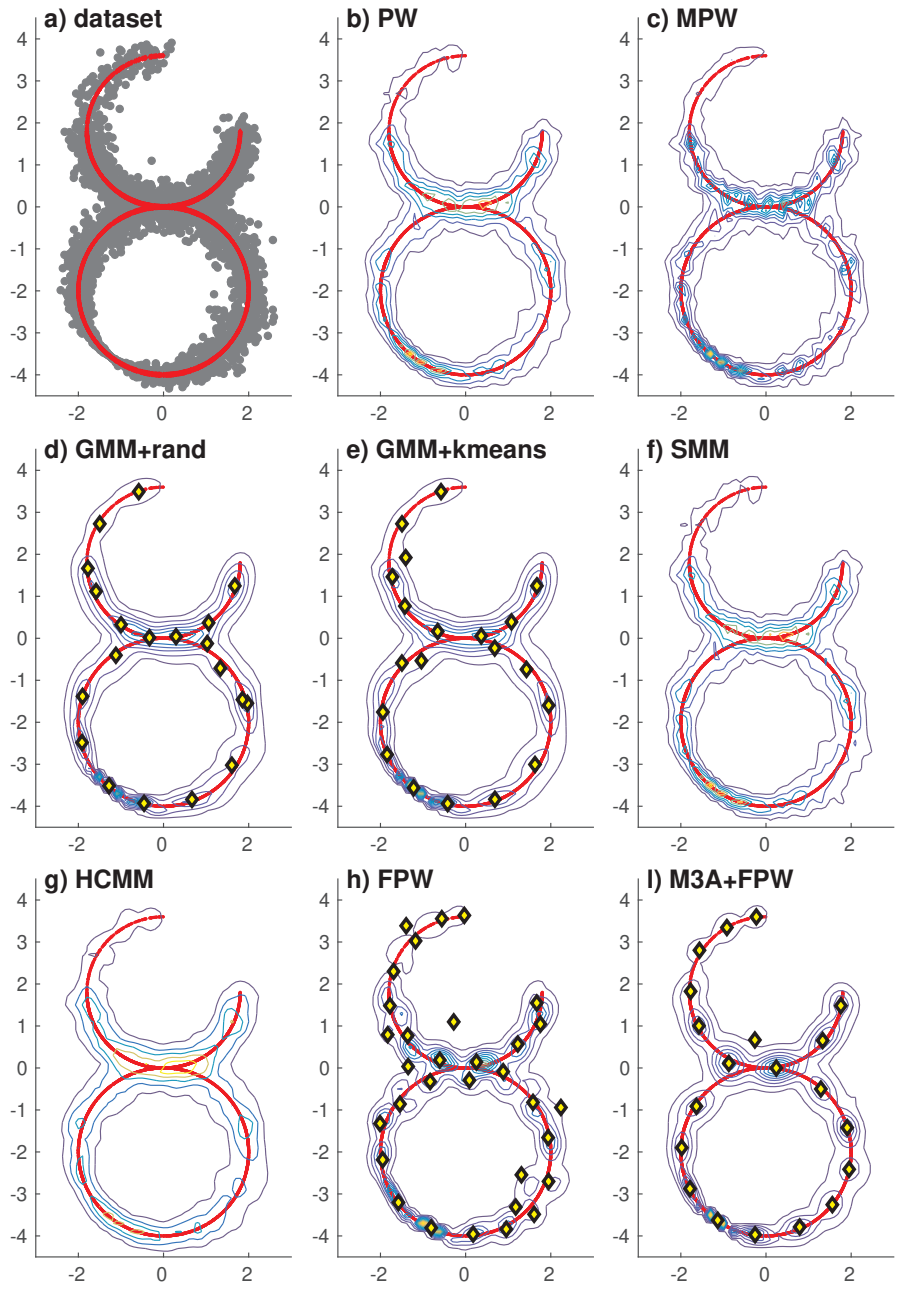


Figure 10: a) two circular manifolds with variation in density and noise level, b-l) the contour curves for the compared eight density estimators. Note that black diamonds represent the means of the Gaussian distributions.

Astronomical simulations are often used to study the evolution of astronomical objects through time. Since a simulation often includes many Millions of particles, it is

Table 2: Comparison of density estimators for the two circular manifolds.

Algorithm	Parameters	ALL (std)	
PW	$\sigma = 0.09$	-2.668	(0.0)
MPW	$k = 140$	-2.689	(0.0)
FPW	$r_{\text{ball}} = 0.75$ ( $m = 31$ )	-2.682	(0.0116)
SMM	$m = 800$	-2.671	(0.0007)
HCMM	$m = 130, \sigma = 0.15$	-2.700	(0.0014)
GMM	$m = 20$	-2.620	(0.0036)
GMM <sup>k</sup>	$m = 20$	-2.621	(0.0034)
M3A+FPW	$r_{rw} = 0.6, r_{\text{ball}} = 0.95$ ( $m = 20.6 \pm 1$ )	-2.650	(0.0027)

not possible to track each particle individually. Therefore, it is desirable to encode the distribution of particles via a sparse pdf and follow its evolution instead. Since our goal in this section is to compare different density estimators, we only use one time snapshot of a simulated jellyfish galaxy, as shown in Figure 11 a and b. This data set is challenging since the density changes significantly in the data space. While the head of the jellyfish is very dense the outer parts depict much lower density. Those parts contain some low-dimensional structures called streams, as highlighted by boxes in panels a and b. During the evolution process these structures are more affected by other astronomical objects and change over time. Thus, in addition to the dense regions a suitable density estimator should capture these lower density regions well.

We compare our M3A+FPW algorithm with the result of GMM initialized by k-means visually and additionally report the quality of all density estimators as before.

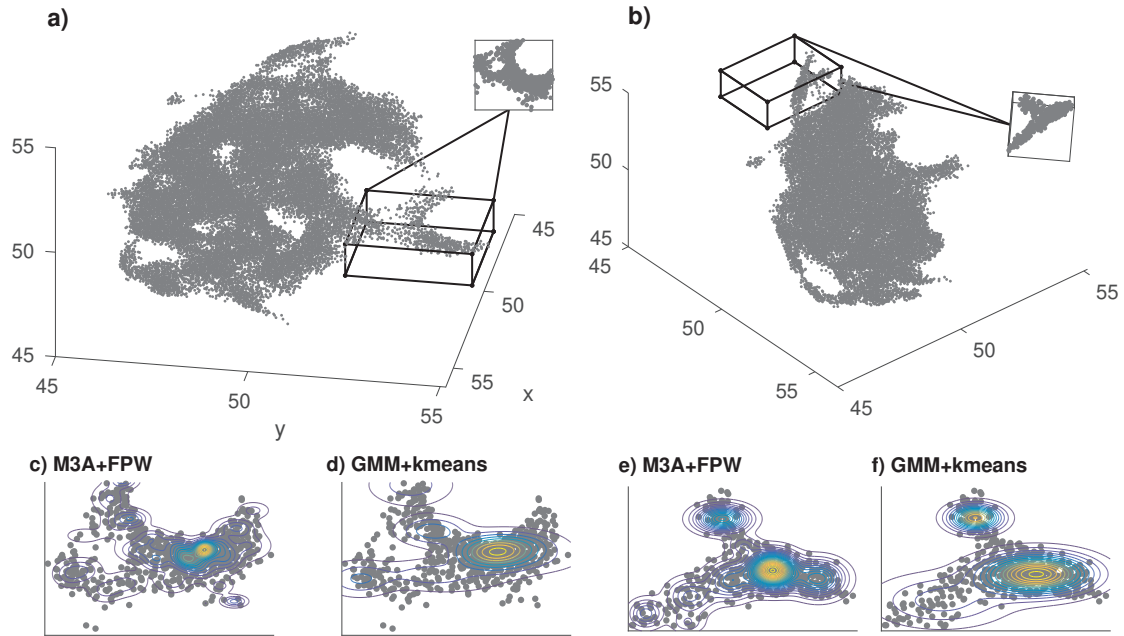


Figure 11: The simulated Jellyfish galaxy (viewpoint a and b). c-f) contour curves of density estimators projected by the first two principal directions of Gaussian co-variance matrices for zoomed regions a) (left) and b) (right) using Algorithm 3 (c,e) and GMM initialized by k-means (d,f).

From Table 3 it can be seen that the GMM outperforms the other methods in terms of ALL, which is expected since it aims to optimize the likelihood function. However, its optimization process is time-consuming for big data sets, which makes it challenging to be applied in big simulations, as indicated by large training times reported in the last column of Table 3<sup>9</sup>. PW as before provides an unnecessary complex model, which is not very desirable for this application. Alternatively, the M3A+FPW strategy presents a sparse model while preserving a comparable ALL.

<sup>9</sup>Here, we distribute the M3A ants among 10 processors (i.e. 5 ants per CPU) and restrict the number of step in each round to  $n = 10000$ .

Table 3: Comparison of density estimators for the simulated jellyfish galaxy.

Algorithm	Parameters	ALL (std)	$d_H$ (std)	time in s (std)
PW	$\sigma = 0.08$	-3.775 (0.018)	0.264 (0.040)	0.2 (0.0)
FPW	$r = 0.4, m = 949$	-3.855 (0.023)	0.633 (0.071)	9.2 (0.4)
SMM	$m = 4100$	-3.815 (0.021)	0.699 (0.115)	144.8 (10.5)
SMM	$m = 600$	-4.249 (0.032)	0.393 (0.053)	63.0 (5.2)
HCMM	$\sigma = 0.17, m = 600$	-3.967 (0.014)	0.697 (0.091)	12.4 (0.5)
GMM	$m = 300$	-3.772 (0.025)	2.916 (0.600)	1244.4 (257.4)
GMM <sup>k</sup>	$m = 300$	-3.767 (0.024)	3.060 (0.654)	909.3 (63.1)
M3A+FPW	$r_{rw} = 0.4,$ $r_{ball} = 0.5, m = 566$	-3.838 (0.025)	0.737 (0.109)	165.5 (2.1)
M3A+FPW	$r_{rw} = 0.4,$ $r_{ball} = 0.65, m = 309$	-3.882 (0.024)	0.823 (0.072)	151.2 (11.9)

In addition to ALL we also report the Hausdorff distance

$$d_H(A, B) = \max \left\{ \max_{a \in A} \min_{b \in B} d(a, b), \max_{b \in B} \min_{a \in A} d(a, b) \right\}$$

where  $d$  denotes the Euclidean distance, and  $A$  and  $B$  are the set of particles and the set of centers of Gaussian models, respectively.

It is informative since it shows how much the model captures the original low-density structures. To compare our method to GMM we increase the balls' radius to  $r_{ball} = 0.65$ . From the table it is clear that the Hausdorff distance of the proposed method is much lower than the one achieved by GMM. To visually see its success we added the contour plots (see panels c-f). It can be seen that the new method is more

successful in fitting the pdf along with the low-dimensional structure. In summary, if selecting a model is a trade-off between simplicity (sparsity) and capturing structural details the M3A+FPW strategy, outlined in Algorithm 3, provides a sparse model with a very competitive ALL value.

## Conclusion

In some applications, such as astronomy, it is common to have low-dimensional structures buried inside big data sets. Therefore, it is desirable to have a method that extracts these structures of varying density while their continuity is kept unchanged. Although there are several methods, they often fail to either extract manifolds or keep their continuity. Inspired by the ant colony algorithm we propose a new method where a value, called pheromone, is assigned to each sample. Later these quantities can be used to reveal manifolds without the undesirable effects. To study the behavior of the algorithm and the effect of its hyper-parameters, we provide a theoretical analysis using the Markov Chain framework, where we consider a noisy manifold and apply the algorithm examining two cases: a) a random walker with an unbounded number of steps, and b) an unbounded number of walkers with a fixed number of steps. We show for both cases that the pheromone distribution captures information about the distances of data points to the underlying manifold. In addition to the theory, we empirically analyze the algorithm using synthetic and real data sets, demonstrating three different scenarios: a) denoising and clustering manifolds, b) visualizing data, and c) density estimation. In all investigated scenarios and application examples strategies based using M3A exhibit

comparable or superior results in suppressing the noise, capturing the manifolds and providing sparser models.

## Acknowledgments

This work was supported by the European H2020-MSCA-ITN *SURvey Network for Deep Imaging Analysis and Learning* (SUNDIAL), project ID 721463. We thank the Center for Information Technology of the University of Groningen for their support and for providing access to the Peregrine high performance computing cluster. Furthermore, this work has made use of data from the European Space Agency (ESA) mission *Gaia* (<https://www.cosmos.esa.int/gaia>), processed by the *Gaia* Data Processing and Analysis Consortium (DPAC, <https://www.cosmos.esa.int/web/gaia/dpac/consortium>). Funding for the DPAC has been provided by national institutions, in particular the institutions participating in the *Gaia* Multilateral Agreement.

## References

- Belkin, M. and Niyogi, P. (2003). Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396.
- Bentley, J. L. (1975). Multidimensional binary search trees used for associative searching. *Communications of the ACM*, 18(9):509–517.
- Berry, T. and Sauer, T. (2016). Local kernels and the geometric structure of data. *Applied and Computational Harmonic Analysis*, 40(3):439–469.

- Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.
- Blum, C., Roli, A., and Dorigo, M. (2001). Hc-aco: The hyper-cube framework for ant colony optimization. In *Proceedings of MIC*, volume 2, pages 399–403.
- Chu, S.-C., Roddick, J. F., Su, C.-J., and Pan, J.-S. (2004). Constrained ant colony optimization for data clustering. In *Pacific Rim International Conference on Artificial Intelligence*, pages 534–543. Springer.
- Cinlar, E. (2013). *Introduction to stochastic processes*. Courier Corporation.
- Coifman, R. R. and Lafon, S. (2006). Diffusion maps. *Applied and computational harmonic analysis*, 21(1):5–30.
- Collaboration, G., Brown, A., Vallenari, A., Prusti, T., de Bruijne, J., Babusiaux, C., Bailer-Jones, C., Biermann, M., Evans, D., Eyer, L., et al. (2018). Gaia data release 2. summary of the contents and survey properties. *Astronomy and Astrophysics*, 616:1–22.
- Dixit, P. D. (2019). Introducing user-prescribed constraints in markov chains for non-linear dimensionality reduction. *Neural computation*, 31(5):980–997.
- Donoho, D. L. and Grimes, C. (2003). Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data. *Proceedings of the National Academy of Sciences*, 100(10):5591–5596.
- Dorigo, M. (1992). Optimization, learning and natural algorithms. *Ph. D. Thesis*, Politecnico di Milano.



- Dorigo, M., Maniezzo, V., and Colorni, A. (1991). Positive feedback as a search strategy. Technical Report 91-016, Politecnico di Milano, Italy.
- Dorigo, M., Maniezzo, V., and Colorni, A. (1996). Ant system: optimization by a colony of cooperating agents. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 26(1):29–41.
- Gil, C., Baños, R., Ortega, J., Márquez, A. L., Fernández, A., and Montoya, M. (2011). Ant colony optimization for water distribution network design: a comparative study. In *International Work-Conference on Artificial Neural Networks*, pages 300–307. Springer.
- Goldberger, J. and Roweis, S. T. (2005). Hierarchical clustering of a mixture model. In *Advances in Neural Information Processing Systems*, pages 505–512. Citeseer.
- Golub, G. H. and Van Loan, C. F. (2012). *Matrix computations*, volume 3. JHU press.
- Gong, D., Sha, F., and Medioni, G. (2010). Locally linear denoising on image manifolds. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 265–272. JMLR Workshop and Conference Proceedings.
- Gorban, A. N. and Zinovyev, A. (2010). Principal manifolds and graphs in practice: from molecular biology to dynamical systems. *International journal of neural systems*, 20(03):219–232.
- Hein, M. and Maier, M. (2006). Manifold denoising. In *NIPS*, volume 19, pages 561–568.

- Kaslovsky, D. N. and Meyer, F. G. (2014). Non-asymptotic analysis of tangent space perturbation. *Information and Inference: a Journal of the IMA*, 3(2):134–187.
- Klicpera, J., Weißenberger, S., and Günnemann, S. (2019). Diffusion improves graph learning. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E. B., and Garnett, R., editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 13333–13345.
- Kulkarni, V. G. (1999). *Modeling, Analysis, Design, and Control of Stochastic Systems*. Springer.
- Lerman, G., McCoy, M. B., Tropp, J. A., and Zhang, T. (2015). Robust computation of linear models by convex relaxation. *Foundations of Computational Mathematics*, 15(2):363–410.
- Little, A. V., Maggioni, M., and Murphy, J. M. (2020). Path-based spectral clustering: Guarantees, robustness to outliers, and fast algorithms. *Journal of machine learning research*, 21:1–66.
- Little, A. V., Maggioni, M., and Rosasco, L. (2017). Multiscale geometric methods for data sets i: Multiscale svd, noise and curvature. *Applied and Computational Harmonic Analysis*, 43(3):504–567.
- Maniezzo, V. (1999). Exact and approximate nondeterministic tree-search procedures for the quadratic assignment problem. *INFORMS journal on computing*, 11(4):358–369.

- Mohammadi, M. and Bunte, K. (2020). Multi-agent based manifold denoising. In *International Conference on Intelligent Data Engineering and Automated Learning*, pages 12–24. Springer.
- Mohammadi, M., Petkov, N., Bunte, K., Peletier, R. F., and Schleif, F.-M. (2019). Global cluster detection in the gaia survey. *Neurocomputing*, 342:164–171.
- Rizzoli, A. E., Oliverio, F., Montemanni, R., and Gambardella, L. M. (2004). Ant colony optimisation for vehicle routing problems: from theory to applications. *Galleria Rassegna Bimestrale Di Cultura*, 9(1):1–50.
- Roweis, S. T. and Saul, L. K. (2000). Nonlinear dimensionality reduction by locally linear embedding. *science*, 290(5500):2323–2326.
- Runkler, T. A. (2005). Ant colony optimization of clustering models. *International Journal of Intelligent Systems*, 20(12):1233–1251.
- Stützle, T. and Hoos, H. H. (2000). Max–min ant system. *Future generation computer systems*, 16(8):889–914.
- Tsai, C.-F., Tsai, C.-W., Wu, H.-C., and Yang, T. (2004). Acodf: a novel data clustering approach for data mining in large databases. *Journal of Systems and Software*, 73(1):133–145.
- Tu, L. W. (2011). *An introduction to manifolds*. Springer.
- Van der Maaten, L. and Hinton, G. (2008). Visualizing data using t-sne. *Journal of machine learning research*, 9(11):2579–2605.

- Vincent, P. and Bengio, Y. (2003). Manifold parzen windows. *Advances in neural information processing systems*, pages 849–856.
- Von Luxburg, U. (2007). A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416.
- Wang, W. and Carreira-Perpinán, M. A. (2010). Manifold blurring mean shift algorithms for manifold denoising. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 1759–1766. IEEE.
- Wang, X., Tino, P., and Fardal, M. A. (2008). Multiple manifolds learning framework based on hierarchical mixture density model. In *ECML PKDD*, pages 566–581.
- Wang, X., Tino, P., Fardal, M. A., Raychaudhury, S., and Babul, A. (2009). Fast parzen window density estimator. In *2009 International Joint Conference on Neural Networks*, pages 3267–3274. IEEE.
- Zhang, K. and Kwok, J. T. (2010). Simplifying mixture models through function approximation. *IEEE Transactions on Neural Networks*, 21(4):644–658.
- Zhang, Z. and Zha, H. (2003). Nonlinear dimension reduction via local tangent space alignment. In *International Conference on Intelligent Data Engineering and Automated Learning*, pages 477–481. Springer.