

---

# Least-Squares Log-Density Gradient Clustering for Riemannian Manifolds

---

Mina Ashizawa<sup>1</sup>

Hiroaki Sasaki<sup>2,3</sup>

Tomoya Sakai<sup>1</sup>

Masashi Sugiyama<sup>3,1</sup>

<sup>1</sup>The University of Tokyo, Japan    <sup>2</sup>Nara Institute of Science and Technology, Japan    <sup>3</sup>RIKEN, Japan

## Abstract

*Mean shift* is a mode-seeking clustering algorithm that has been successfully used in a wide range of applications such as image segmentation and object tracking. To further improve the clustering performance, mean shift has been extended to various directions, including generalization to handle data on Riemannian manifolds and extension to directly estimating the log-density gradient without density estimation. In this paper, we combine these ideas and propose a novel mode-seeking algorithm for Riemannian manifolds with direct log-density gradient estimation. Although the idea of combining the two extensions is rather straightforward, directly estimating the log-density gradient on Riemannian manifolds is mathematically challenging. We will provide a mathematically sound algorithm and demonstrate its usefulness through experiments.

## 1 Introduction

*Clustering* is one of the most important unsupervised learning tasks in machine learning and has been extensively studied for decades (Clarke et al., 2009; Murphy, 2012). Among various different types of clustering algorithms, *mode-seeking* is a well-studied and practically useful approach. *Mean shift* (Fukunaga and Hostetler, 1975; Cheng, 1995; Comaniciu and Meer, 2002; Carreira-Perpinán, 2015) is a seminal algorithm for mode-seeking clustering: Kernel density estimation is first performed on given data points and then the data points are updated along the gradient of the estimated density towards the modes. Finally, the data

points which converged to the same mode are given the same cluster label. A notable advantage of mean shift is that it does not require to specify the number of clusters in advance. Thanks to this useful property, mean shift has been successfully employed in a wide range of real-world applications such as image segmentation (Wang et al., 2004; Tao et al., 2007) and object tracking (Comaniciu et al., 2000; Collins, 2003).

The original mean shift algorithm considers data points in the Euclidean space. However, in practice, data points sometimes lie on a structured space such as the *Lie group* and *Grassmann manifold*. For data on such a structured space, kernel density estimation and gradient ascent with the Euclidean metric do not necessarily perform appropriately. To cope with this problem, mean shift has been extended to Riemannian manifolds (Boothby, 2003) and demonstrated to work well in experiments (Tuzel et al., 2005; Subbarao and Meer, 2006, 2009; Cetingul and Vidal, 2009).

Another important extension of mean shift is to avoid density estimation. The original mean shift algorithm uses kernel density estimation, which tends to perform poorly when the data dimension is high. Furthermore, a good density estimator does not necessarily mean a good density gradient estimator, and thus the two-step approach of first estimating the density and then computing its gradient does not always perform well. To cope with this problem, a method to directly estimate the log-density gradient without density estimation has been developed (Cox, 1985; Sasaki et al., 2014), and a mode-seeking clustering algorithm based on the direct log-density gradient estimator was experimentally shown to work well (Sasaki et al., 2014).

The purpose of this paper is to combine these two extensions and propose a novel clustering algorithm based on direct log-density gradient estimation and mode-seeking on Riemannian manifolds. Although the idea of combining the two extensions is rather straightforward, directly estimating the density gradient on Riemannian manifolds is mathematically challenging. We will provide a mathematically sound algorithm and demonstrate its usefulness through experiments.

## 2 Problem Formulation

In this section, we formulate the clustering problem by mode-seeking and review existing algorithms.

**Clustering by Mode-Seeking:** Suppose that we are given independent and identically distributed samples of size  $n$  on  $\mathbb{R}^d$ ,  $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^n$ , with unknown probability density  $p(\mathbf{x})$ . The goal of clustering is to split the set  $\mathcal{X}$  into  $c$  disjoint subsets  $\{\mathcal{X}_i\}_{i=1}^c$  so that samples in each subset share similar properties while samples in different subsets have different properties.

Various types of clustering algorithms have been explored so far (Clarke et al., 2009; Murphy, 2012). Among them, *mode-seeking* is one of the popular and well-studied approaches (Fukunaga and Hostetler, 1975; Cheng, 1995; Comaniciu and Meer, 2002; Carreira-Perpinán, 2015). In mode-seeking clustering, data samples  $\{\mathbf{x}_i\}_{i=1}^n$  are first gathered to the modes of data density  $p(\mathbf{x})$  by, e.g., gradient ascent  $\mathbf{x} \leftarrow \mathbf{x} + \varepsilon \nabla p(\mathbf{x})$ , where  $\varepsilon > 0$  is the step size and  $\nabla p(\mathbf{x})$  is the gradient of the density function  $p(\mathbf{x})$  with respect to  $\mathbf{x} = (x^{(1)}, \dots, x^{(d)})^\top$ . Then, data samples which converged to the same mode are given the same cluster label.

Below, we review representative mode-seeking clustering algorithms.

**Mean Shift Clustering:** *Mean shift* is a seminal algorithm of mode-seeking clustering (Fukunaga and Hostetler, 1975; Cheng, 1995; Comaniciu and Meer, 2002). In the mean shift algorithm, the probability density  $p(\mathbf{x})$  is first learned by kernel density estimation:

$$\hat{p}(\mathbf{x}) = \frac{c_{k,\sigma}}{n} \sum_{i=1}^n k\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{\sigma}\right\|^2\right), \quad (1)$$

where  $k(t)$  is a non-negative function,  $\sigma > 0$  is the bandwidth, and  $c_{k,\sigma}$  is the normalization constant such that the integration of  $\hat{p}(\mathbf{x})$  is equal to 1. As function  $k(t)$ , the exponential decaying function  $k(t) = \exp(-t/2)$  is often used in practice, which yields the Gaussian kernel:

$$k\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{\sigma}\right\|^2\right) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{2\sigma^2}\right).$$

The bandwidth  $\sigma$  can be systematically chosen by cross-validation with respect to the log-likelihood or squared error criteria.

Next, the gradient of the kernel density estimator  $\hat{p}(\mathbf{x})$  is computed:

$$\nabla \hat{p}(\mathbf{x}) = \frac{c_{k',\sigma}}{n} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{x}) k'\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{\sigma}\right\|^2\right) = \epsilon(\mathbf{x}) \mathbf{m}(\mathbf{x}),$$

where  $c_{k',\sigma} = -2c_{k,\sigma}/\sigma^2$ ,  $k'$  is the derivative of  $k$ ,

$$\epsilon(\mathbf{x}) := \frac{c_{k',\sigma}}{n} \sum_{i=1}^n k'\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{\sigma}\right\|^2\right) > 0,$$

and  $\mathbf{m}(\mathbf{x})$  is called the *mean shift vector* (Comaniciu and Meer, 2002):

$$\mathbf{m}(\mathbf{x}) := \frac{\sum_{i=1}^n \mathbf{x}_i k'\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{\sigma}\right\|^2\right)}{\sum_{i'=1}^n k'\left(\left\|\frac{\mathbf{x} - \mathbf{x}_{i'}\right\|^2\right)} - \mathbf{x}.$$

To obtain the modes of data density, the mean shift algorithm uses the *fixed-point iteration* for mode-seeking. More specifically, a necessary condition for local maximum  $\nabla \hat{p}(\mathbf{x}) = \mathbf{0}$  implies  $\mathbf{m}(\mathbf{x}) = \mathbf{0}$ , which yields  $\mathbf{x} \leftarrow \mathbf{x} + \mathbf{m}(\mathbf{x})$ . Since this fixed-point update rule can be rewritten as  $\mathbf{x} \leftarrow \mathbf{x} + \frac{1}{\epsilon(\mathbf{x})} \nabla \hat{p}(\mathbf{x})$ , it can be interpreted as gradient ascent with adaptive step size  $1/\epsilon(\mathbf{x}) > 0$ .

After repeating the fixed-point update for all samples  $\{\mathbf{x}_i\}_{i=1}^n$  until convergence, samples which converge to the same mode are given the same cluster label.

The mean shift algorithm has been successfully employed in various real-world applications such as image segmentation (Wang et al., 2004; Tao et al., 2007) and object tracking (Comaniciu et al., 2000; Collins, 2003).

**Mean Shift Clustering on Riemannian Manifolds:** The original mean shift algorithm considers data points in the Euclidean space. However, in practice, data points sometimes lie on a structured space such as the Lie group and Grassmann manifold. For data on such a structured space, kernel density estimation and gradient ascent with the Euclidean metric do not necessarily perform appropriately. To cope with this problem, the mean shift algorithm has been extended to Riemannian manifolds (Tuzel et al., 2005; Subbarao and Meer, 2006, 2009; Cetingul and Vidal, 2009). Here we briefly review such an extension. Let us consider data points  $\{\mathbf{X}_i\}_{i=1}^n$  on Riemannian manifold  $\mathcal{M}^m$  of dimension  $m \leq d$  embedded in  $\mathbb{R}^d$ .

As shown in Eq.(1), kernel density estimation in the original mean shift algorithm uses the squared Euclidean distance between  $\mathbf{x}$  and  $\mathbf{x}_i$ , i.e.,  $\|\mathbf{x} - \mathbf{x}_i\|^2$ . For data points on a Riemannian manifold, this may be replaced with the squared *geodesic distance* between  $\mathbf{X}$  and  $\mathbf{X}_i$ :

$$\hat{p}(\mathbf{X}) = \frac{c_{k,\sigma,\delta}}{n} \sum_{i=1}^n k\left(\frac{\delta(\mathbf{X}, \mathbf{X}_i)^2}{\sigma^2}\right), \quad (2)$$

where  $\delta(\mathbf{X}, \mathbf{X}')$  denotes the geodesic distance between  $\mathbf{X}$  and  $\mathbf{X}'$ , and  $c_{k,\sigma,\delta}$  is a positive constant.<sup>1</sup>

<sup>1</sup>Strictly speaking, Eq.(2) is inappropriate as a density

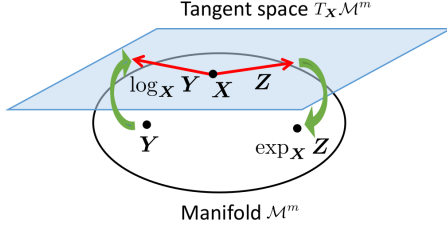


Figure 1: The exponential map and the logarithm map. The exponential map at the point  $\mathbf{X}$  transforms points on the tangent space at  $\mathbf{X}$  to the manifold, while the logarithm map at  $\mathbf{X}$  transforms points on the manifold to the tangent space at  $\mathbf{X}$ .

Then, the gradient of the estimated density is given by

$$\begin{aligned} \nabla \hat{p}(\mathbf{X}) &= \frac{-c_{k',\sigma,\delta}}{2n} \sum_{i=1}^n [\nabla \delta(\mathbf{X}, \mathbf{X}_i)^2] k' \left( \frac{\delta(\mathbf{X}, \mathbf{X}_i)^2}{\sigma^2} \right) \\ &= \frac{c_{k',\sigma,\delta}}{n} \sum_{i=1}^n [\log_{\mathbf{X}} \mathbf{X}_i] k' \left( \frac{\delta(\mathbf{X}, \mathbf{X}_i)^2}{\sigma^2} \right), \end{aligned} \quad (3)$$

where  $c_{k',\sigma,\delta} = -2c_{k,\sigma,\delta}/\sigma^2$ .  $\log_{\mathbf{X}} \mathbf{X}_i$  denotes the *logarithm map* of  $\mathbf{X}_i$  at  $\mathbf{X}$ , which satisfies the following relation (Boothby, 2003):

$$\log_{\mathbf{X}} \mathbf{X}_i = -\frac{1}{2} \nabla \delta(\mathbf{X}, \mathbf{X}_i)^2. \quad (4)$$

In the same way as the original mean shift algorithm, Eq.(3) can be expressed as  $\nabla \hat{p}(\mathbf{X}) = \epsilon(\mathbf{X})\mathbf{M}(\mathbf{X})$ , where  $\epsilon(\mathbf{X}) := \frac{c_{k',\sigma,\delta}}{n} \sum_{i=1}^n k' \left( \frac{\delta(\mathbf{X}, \mathbf{X}_i)^2}{\sigma^2} \right) > 0$  and  $\mathbf{M}(\mathbf{X})$  is the mean shift vector defined as

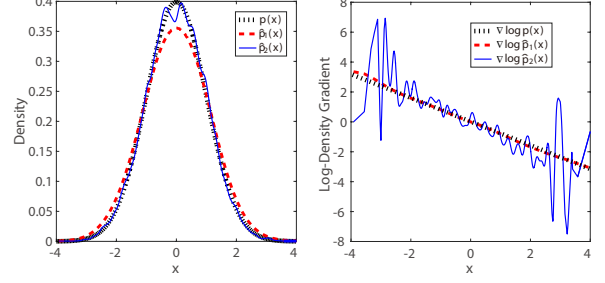
$$\mathbf{M}(\mathbf{X}) := \frac{\sum_{i=1}^n [\log_{\mathbf{X}} \mathbf{X}_i] k' \left( \frac{\delta(\mathbf{X}, \mathbf{X}_i)^2}{\sigma^2} \right)}{\sum_{i'=1}^n k' \left( \frac{\delta(\mathbf{X}, \mathbf{X}_{i'})^2}{\sigma^2} \right)}.$$

This shows that mean shift vector  $\mathbf{M}(\mathbf{X})$  lies on the *tangent space* at  $\mathbf{X}$ , which is denoted by  $T_{\mathbf{X}}\mathcal{M}^m$  (see Figure 1).

Thus, when a point  $\mathbf{X}$  on the manifold is updated along with this vector, it no longer lies on the manifold. To cope with this problem, the updated point is projected back to the manifold by the *exponential map*  $\mathbf{X} \leftarrow \exp_{\mathbf{X}} \mathbf{M}(\mathbf{X})$  (see Figure 1 again) (Boothby, 2003).

After repeating this update for all samples  $\{\mathbf{X}_i\}_{i=1}^n$  until convergence, samples which converge to the same mode are given the same cluster label.

estimator on Riemannian manifolds: To ensure that the integration of  $\hat{p}(\mathbf{X})$  is equal to 1, the normalizing constant of each  $k \left( \frac{\delta(\mathbf{X}, \mathbf{X}_i)^2}{\sigma^2} \right)$  has to depend on  $\mathbf{X}_i$  (Pelletier, 2005; Arvanitidis et al., 2016). Nonetheless, Eq.(2) has been employed in the mean shift algorithm because of its computational efficiency and simple form (Subbarao and Meer, 2009).



(a) Data densities (b) Log-density gradients

Figure 2: Density estimation and log-density gradient estimation.  $\hat{p}_2(x)$  is a better estimate of true density  $p(x)$  than  $\hat{p}_1(x)$ , but  $\nabla \log \hat{p}_2(x)$  is a worse estimate of true log-density gradient  $\nabla \log p(x)$  than  $\nabla \log \hat{p}_1(x)$ . Thus, a good density estimator is not necessarily a good log-density gradient estimator.

The mean shift algorithm on Riemannian manifolds has been shown to work well experimentally (Tuzel et al., 2005; Subbarao and Meer, 2006, 2009; Cetingu and Vidal, 2009).

### Least-Squares Log-Density Gradient Clustering:

Another important extension of the original mean shift algorithm is to avoid density estimation (Sasaki et al., 2014). Kernel density estimation used in the original mean shift algorithm tends to perform poorly when the data dimension  $d$  is high. Furthermore, a good density estimator does not necessarily mean a good density gradient estimator (see Figure 2), and thus the two-step approach of first estimating the density and then computing its gradient does not always perform well. To cope with these problems, a method to directly estimate the density gradient without density estimation has been developed (Cox, 1985; Sasaki et al., 2014). Here, we briefly review the direct log-density gradient estimator called the *least-squares log-density gradient* (LSLDG) and the mode-seeking clustering algorithm based on it called *LSLDG clustering* (Sasaki et al., 2014).

Let  $\mathbf{g}(\mathbf{x}) := (g^{(1)}(\mathbf{x}), \dots, g^{(d)}(\mathbf{x}))^\top$  be the gradient of the log-density function, where  $g^{(j)}(\mathbf{x}) := \partial_j \log p(\mathbf{x})$  and  $\partial_j = \frac{\partial}{\partial x^{(j)}}$ . The key idea of LSLDG is to directly fit a model  $\tilde{g}^{(j)}(\mathbf{x})$  to the true log-density gradient  $g^{(j)}(\mathbf{x})$  under the squared loss:

$$\begin{aligned} J^{(j)}(\tilde{g}^{(j)}(\mathbf{x})) &:= \int \left( \tilde{g}^{(j)}(\mathbf{x}) - g^{(j)}(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x} - C \\ &= \int \tilde{g}^{(j)}(\mathbf{x})^2 p(\mathbf{x}) d\mathbf{x} - 2 \int \tilde{g}^{(j)}(\mathbf{x}) \partial_j p(\mathbf{x}) d\mathbf{x} \\ &= \int \tilde{g}^{(j)}(\mathbf{x})^2 p(\mathbf{x}) d\mathbf{x} + 2 \int \partial_j \tilde{g}^{(j)}(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}, \end{aligned}$$

where  $C = \int g^{(j)}(\mathbf{x})^2 p(\mathbf{x}) d\mathbf{x}$  and the last equality follows from *integration by parts* under

$\lim_{x^{(j)} \rightarrow \pm\infty} \tilde{g}^{(j)}(\mathbf{x})p(\mathbf{x}) = 0$ . Then the empirical approximation of  $J^{(j)}$  is given as

$$J^{(j)}(\tilde{g}^{(j)}(\mathbf{x})) \approx \frac{1}{n} \sum_{i=1}^n \tilde{g}^{(j)}(\mathbf{x}_i)^2 + \frac{2}{n} \sum_{i=1}^n \partial_j \tilde{g}^{(j)}(\mathbf{x}_i). \quad (5)$$

As the log-density gradient model  $\tilde{g}^{(j)}(\mathbf{x})$ , a linear-in-parameter model is used:

$$\tilde{g}^{(j)}(\mathbf{x}) = \boldsymbol{\theta}^{(j)\top} \boldsymbol{\psi}^{(j)}(\mathbf{x}) = \sum_{l=1}^b \theta_l^{(j)} \psi_l^{(j)}(\mathbf{x}),$$

where  $b$  denotes the number of parameters,  $\boldsymbol{\theta}^{(j)} \in \mathbb{R}^b$  is the parameter vector, and  $\boldsymbol{\psi}^{(j)}(\mathbf{x}) \in \mathbb{R}^b$  is the vector of basis functions. By plugging this linear model into Eq.(5) and adding the  $\ell_2$ -regularizer to avoid overfitting, the following optimization problem is obtained:

$$\hat{\boldsymbol{\theta}}^{(j)} = \underset{\boldsymbol{\theta}^{(j)} \in \mathbb{R}^b}{\operatorname{argmin}} \left[ \boldsymbol{\theta}^{(j)\top} \widehat{\mathbf{G}}^{(j)} \boldsymbol{\theta}^{(j)} + 2\boldsymbol{\theta}^{(j)\top} \widehat{\mathbf{h}}^{(j)} + \lambda^{(j)} \boldsymbol{\theta}^{(j)\top} \boldsymbol{\theta}^{(j)} \right],$$

where  $\lambda^{(j)} \geq 0$  is the regularization parameter, and  $\widehat{\mathbf{G}}^{(j)}$  and  $\widehat{\mathbf{h}}^{(j)}$  are defined as

$$\widehat{\mathbf{G}}^{(j)} := \frac{1}{n} \sum_{i=1}^n \boldsymbol{\psi}^{(j)}(\mathbf{x}_i) \boldsymbol{\psi}^{(j)}(\mathbf{x}_i)^\top, \quad \widehat{\mathbf{h}}^{(j)} := \frac{1}{n} \sum_{i=1}^n \partial_j \boldsymbol{\psi}^{(j)}(\mathbf{x}_i)$$

The optimal solution  $\hat{\boldsymbol{\theta}}^{(j)}$  is obtained analytically by

$$\hat{\boldsymbol{\theta}}^{(j)} = -(\widehat{\mathbf{G}}^{(j)} + \lambda^{(j)} \mathbf{I}_b)^{-1} \widehat{\mathbf{h}}^{(j)},$$

where  $\mathbf{I}_b$  denotes the  $b \times b$  identity matrix. All hyper-parameters such as the regularization parameter  $\lambda^{(j)}$  and basis parameters in  $\boldsymbol{\psi}^{(j)}(\mathbf{x})$  can be systematically chosen by cross-validation with respect to the squared error criterion  $J^{(j)}$ . This direct log-density gradient estimator is called LSLDG.

To derive a mean shift like fixed-point algorithm from LSLDG, for the Gaussian kernel,  $\phi_l^{(j)}(\mathbf{x}) := \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_l\|^2}{2\sigma^{(j)^2}}\right)$ , where the centers  $\{\mathbf{c}_l\}_{l=1}^b$  are chosen randomly from  $\{\mathbf{x}_i\}_{i=1}^n$  without overlap, its partial derivative was proposed to be used as basis functions (Sasaki et al., 2014):

$$\psi_l^{(j)}(\mathbf{x}) := \partial_j \phi_l^{(j)}(\mathbf{x}) = \frac{1}{\sigma^{(j)^2}} (c_l^{(j)} - x^{(j)}) \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_l\|^2}{2\sigma^{(j)^2}}\right).$$

Suppose  $\sum_{l=1}^b \hat{\theta}_l^{(j)} \phi_l^{(j)}(\mathbf{x}) \neq 0$ . Then the LSLDG solution can be expressed as

$$\begin{aligned} \hat{g}^{(j)}(\mathbf{x}) &= \sum_{l=1}^b \hat{\theta}_l^{(j)} \psi_l^{(j)}(\mathbf{x}) \\ &= \frac{1}{\sigma^{(j)^2}} \sum_{l=1}^b \hat{\theta}_l^{(j)} (c_l^{(j)} - x^{(j)}) \phi_l^{(j)}(\mathbf{x}) = \hat{\epsilon}^{(j)}(\mathbf{x}) \hat{m}^{(j)}(\mathbf{x}), \end{aligned}$$

where  $\hat{\epsilon}^{(j)}(\mathbf{x}) = \frac{1}{\sigma^{(j)^2}} \sum_{l=1}^b \hat{\theta}_l^{(j)} \phi_l^{(j)}(\mathbf{x})$  and  $\hat{m}^{(j)}(\mathbf{x})$  is the  $j$ -th element of the mean shift vector defined as

$$\hat{m}^{(j)}(\mathbf{x}) := \frac{\sum_{l=1}^b \hat{\theta}_l^{(j)} c_l^{(j)} \phi_l^{(j)}(\mathbf{x})}{\sum_{l'=1}^b \hat{\theta}_{l'}^{(j)} \phi_{l'}^{(j)}(\mathbf{x})} - x^{(j)}.$$

Then, a necessary condition for local maximum  $\partial_j \hat{p}(\mathbf{x}) = 0$  implies  $\hat{m}^{(j)}(\mathbf{x}) = 0$ , which yields  $x^{(j)} \leftarrow x^{(j)} + \hat{m}^{(j)}(\mathbf{x})$ . This update formula can be regarded as a weighted variant of the original mean shift algorithm (Cheng, 1995), and it is reduced to the original mean shift algorithm when  $b = n$  and  $\theta_l^{(j)} = 1/n$ .

The LSLDG clustering algorithm was demonstrated to work well in experiments (Sasaki et al., 2014).

### 3 Proposed Method

In this section, we propose to combine the Riemannian extension of the mean shift algorithm and the LSLDG algorithm, and develop a novel clustering algorithm for Riemannian manifolds. Below, we consider data points  $\{\mathbf{X}_i\}_{i=1}^n$  on a Riemannian manifold  $(\mathcal{M}^m, H)$  of dimension  $m(\leq d)$  embedded in  $\mathbb{R}^d$  with Riemannian metric  $H$ .

**Direct Log-Density-Gradient Estimation on Riemannian Manifolds:** If LSLDG is naively applied to data on Riemannian manifolds, the estimated gradient vector does not necessarily lie on the tangent space. To prevent this problem, we propose to use the *common* parameter vector for all dimensions with basis functions confined in the tangent space. More specifically, let the true log-density gradient vector be  $\mathbf{g}(\mathbf{X}) := \nabla \log p(\mathbf{X}) \in T_{\mathbf{X}} \mathcal{M}^m$ , where  $\nabla$  denotes the Riemannian gradient. We model  $\mathbf{g}(\mathbf{X})$  by  $\tilde{\mathbf{g}}(\mathbf{X}) = \sum_{l=1}^b \theta_l \boldsymbol{\psi}_l(\mathbf{X}) \in T_{\mathbf{X}} \mathcal{M}^m$ , where  $\theta_l$  is the common parameter,  $b$  is the number of parameters, and  $\boldsymbol{\psi}_l(\mathbf{X})$  is the vector of basis functions which we assume to be on the tangent space  $T_{\mathbf{X}} \mathcal{M}^m$ .

This common-parameter model is fitted to the true log-density gradient  $\mathbf{g}(\mathbf{X})$  under the squared loss *on a manifold*:

$$\begin{aligned} J(\tilde{\mathbf{g}}(\mathbf{X})) &:= \int_{\mathcal{M}^m} \|\tilde{\mathbf{g}}(\mathbf{X}) - \mathbf{g}(\mathbf{X})\|_H^2 p(\mathbf{X}) d\operatorname{vol}_{\mathbf{X}} - \bar{C} \\ &= \int_{\mathcal{M}^m} \|\tilde{\mathbf{g}}(\mathbf{X})\|_H^2 p(\mathbf{X}) d\operatorname{vol}_{\mathbf{X}} \\ &\quad - 2 \int_{\mathcal{M}^m} \langle \tilde{\mathbf{g}}(\mathbf{X}), \mathbf{g}(\mathbf{X}) \rangle_H p(\mathbf{X}) d\operatorname{vol}_{\mathbf{X}} \\ &= \int_{\mathcal{M}^m} \|\tilde{\mathbf{g}}(\mathbf{X})\|_H^2 p(\mathbf{X}) d\operatorname{vol}_{\mathbf{X}} \\ &\quad - 2 \int_{\mathcal{M}^m} \langle \tilde{\mathbf{g}}(\mathbf{X}), \nabla p(\mathbf{X}) \rangle_H d\operatorname{vol}_{\mathbf{X}}, \quad (6) \end{aligned}$$

where  $\bar{C} = \int_{\mathcal{M}^m} \|\mathbf{g}(\mathbf{X})\|_H^2 p(\mathbf{X}) d\text{vol}_{\mathbf{X}}$ ,  $\|\cdot\|_H^2 = \langle \cdot, \cdot \rangle_H$  denotes the inner product operator,  $d\text{vol}_{\mathbf{X}}$  denotes a *volume element* of a Riemannian manifold (Petersen, 2006), and the last equality follows from the relation  $\mathbf{g}(\mathbf{X}) = \nabla \log p(\mathbf{X}) = \nabla p(\mathbf{X})/p(\mathbf{X})$ . Applying the “*integration by parts*” formula for manifolds without boundary (Lee, 2012) to the second term in Eq.(6), we obtain

$$\int_{\mathcal{M}^m} \langle \tilde{\mathbf{g}}(\mathbf{X}), \nabla p(\mathbf{X}) \rangle_H d\text{vol}_{\mathbf{X}} = - \int_{\mathcal{M}^m} p(\mathbf{X}) \text{div} \tilde{\mathbf{g}}(\mathbf{X}) d\text{vol}_{\mathbf{X}},$$

where “div” denotes the *divergence* (Petersen, 2006). Approximating the expectation over  $p(\mathbf{X})$  by the average of samples  $\{\mathbf{X}_i\}_{i=1}^n$  and adding the  $\ell_2$ -regularization term, the following optimization problem is obtained:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \mathbb{R}^b}{\text{argmin}} \left[ \boldsymbol{\theta}^\top \hat{\mathbf{G}} \boldsymbol{\theta} + 2\boldsymbol{\theta}^\top \hat{\mathbf{h}} + \lambda \boldsymbol{\theta}^\top \boldsymbol{\theta} \right],$$

where, for  $l, l' = 1, \dots, b$ ,

$$\hat{G}_{l,l'} = \frac{1}{n} \sum_{i=1}^n \langle \psi_l(\mathbf{X}_i), \psi_{l'}(\mathbf{X}_i) \rangle_H, \quad (7)$$

$$\hat{h}_l = \frac{1}{n} \sum_{i=1}^n \text{div}(\psi_l(\mathbf{X}_i)), \quad (8)$$

and  $\lambda \geq 0$  is the regularization parameter. The optimal solution  $\boldsymbol{\theta}$  can be obtained analytically as

$$\hat{\boldsymbol{\theta}} = -(\hat{\mathbf{G}} + \lambda \mathbf{I}_b)^{-1} \hat{\mathbf{h}}.$$

All hyper-parameters such as the regularization parameter  $\lambda$  and basis parameters in  $\psi_l(\mathbf{X})$  can be systematically chosen by cross-validation with respect to the squared error criterion  $J$ . Finally, the log-density gradient estimator is given by  $\hat{\mathbf{g}}(\mathbf{X}) = \sum_{l=1}^b \hat{\theta}_l \psi_l(\mathbf{X})$ .

We call this method *Riemannian LSLDG* (R-LSLDG).

**Mode-Seeking on Riemannian Manifolds:** Based on the estimated log-density gradient  $\hat{\mathbf{g}}(\mathbf{X})$ , we propose a mode-seeking algorithm on Riemannian manifolds.

Let  $\phi_l(\mathbf{X}) = \exp\left(-\frac{\delta(\mathbf{X}, \mathbf{C}_l)^2}{2\sigma^2}\right)$ , where Gaussian centers  $\{\mathbf{C}_l\}_{l=1}^b$  are randomly chosen from data samples  $\{\mathbf{X}_i\}_{i=1}^n$  without overlap. In the same way as the original LSLDG clustering, we use its gradient as basis function vector  $\psi_l(\mathbf{X})$ :

$$\psi_l(\mathbf{X}) = -\frac{[\nabla \delta(\mathbf{X}, \mathbf{C}_l)^2] \phi_l(\mathbf{X})}{2\sigma^2} = \frac{[\log_{\mathbf{X}} \mathbf{C}_l] \phi_l(\mathbf{X})}{\sigma^2},$$

where we used Eq.(4).

Under the assumption that  $\sum_{l=1}^b \hat{\theta}_l \phi_l(\mathbf{X}) \neq 0$  analogously to the original mean shift algorithm, the mean

shift vector is given as

$$\widehat{\mathbf{M}}(\mathbf{X}) = \frac{\sum_{l=1}^b \hat{\theta}_l [\log_{\mathbf{X}} \mathbf{C}_l] \phi_l(\mathbf{X})}{\sum_{l'=1}^b \hat{\theta}_{l'} \phi_{l'}(\mathbf{X})},$$

which always belongs to the tangent space  $T_{\mathbf{X}}\mathcal{M}^m$ . Then it is projected back to the manifold by the exponential map  $\mathbf{X} \leftarrow \exp_{\mathbf{X}} \widehat{\mathbf{M}}(\mathbf{X})$ . After repeating this update for all samples  $\{\mathbf{X}_i\}_{i=1}^n$  until convergence, samples which converge to the same mode are given the same cluster label. We call this clustering method *R-LSLDG clustering* (R-LSLDGC).

Note that the common parameter formulation of LSLDG can be regarded as the limit of the multi-task LSLDG method (Yamane et al., 2016).

For further improvement, we can use a different bandwidth for each Gaussian center  $\mathbf{C}_l$  similarly to Comaniciu et al. (2001). Specifically, we use a basis function

$$\bar{\psi}_l(\mathbf{X}) = \frac{1}{\sigma_l^2} [\log_{\mathbf{X}} \mathbf{C}_l] \bar{\phi}_l(\mathbf{X}),$$

where  $\bar{\phi}_l(\mathbf{X}) = \exp(-\delta(\mathbf{X}, \mathbf{C}_l)^2/(2\sigma_l^2))$ . The mean shift vector is given by

$$\widehat{\mathbf{M}}(\mathbf{X}) = \frac{\sum_{l=1}^b \frac{\hat{\theta}_l}{\sigma_l^2} [\log_{\mathbf{X}} \mathbf{C}_l] \bar{\phi}_l(\mathbf{X})}{\sum_{l'=1}^b \frac{\hat{\theta}_{l'}}{\sigma_{l'}^2} \bar{\phi}_{l'}(\mathbf{X})},$$

where  $\hat{\theta}_l$  is learned by a basis function  $\bar{\psi}_l(\mathbf{X})$ .

**Grassmann Manifold:** In the experiments in the next section, we use the *Grassmann manifold* (Edelman et al., 1998) as an example of Riemannian manifolds. The Grassmann manifold  $\mathcal{G}_{d_1, d_2}$  is the set of  $d_2$ -dimensional linear subspace in  $\mathbb{R}^{d_1}$  for  $d_2 \leq d_1$ :

$$\mathcal{G}_{d_1, d_2} := \{\text{span}(\mathbf{X}) \mid \mathbf{X}^\top \mathbf{X} = \mathbf{I}_{d_2}, \mathbf{X} \in \mathbb{R}^{d_1 \times d_2}\},$$

where  $\text{span}(\mathbf{X})$  denotes the subspace spanned by the columns of  $\mathbf{X}$ . Denoting by  $T_{\mathbf{X}}\mathcal{G}_{d_1, d_2}$  the tangent space on the Grassmann manifold  $\mathcal{G}_{d_1, d_2}$  at location  $\mathbf{X} \in \mathcal{G}_{d_1, d_2}$ , the canonical metric  $\langle \cdot, \cdot \rangle_H$  for the Grassmann manifold is equivalent to the Euclidean metric  $\langle \cdot, \cdot \rangle$  (Edelman et al., 1998). Thus,  $\langle \mathbf{W}, \mathbf{Z} \rangle_H = \text{tr}(\mathbf{W}^\top \mathbf{Z})$  holds for any  $\mathbf{W}, \mathbf{Z} \in T_{\mathbf{X}}\mathcal{G}_{d_1, d_2}$ , where  $\text{tr}(\mathbf{A}) = \sum_{i=1}^d \mathbf{A}_{i,i}$  for a square matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$ . The exponential map for  $\mathbf{Z} \in T_{\mathbf{X}}\mathcal{G}_{d_1, d_2}$  is given by

$$\exp_{\mathbf{X}} \mathbf{Z} = (\mathbf{X} \mathbf{V} \cos \boldsymbol{\Sigma} + \mathbf{U} \sin \boldsymbol{\Sigma}) \mathbf{V}^\top,$$

where  $\mathbf{U}, \boldsymbol{\Sigma}$ , and  $\mathbf{V}$  come from the compact singular value decomposition of  $\mathbf{Z}$ , i.e.,  $\mathbf{Z} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top$ , and “cos” and “sin” for matrix  $\boldsymbol{\Sigma}$  act element-by-element along the diagonal of  $\boldsymbol{\Sigma}$ . The logarithm map for  $\mathbf{Y} \in \mathcal{G}_{d_1, d_2}$  is given by

$$\log_{\mathbf{X}} \mathbf{Y} = (\mathbf{I}_{d_1} - \mathbf{X} \mathbf{X}^\top) \mathbf{Y} \mathbf{Y}^\top \mathbf{X}.$$

The square geodesic distance between two points on the Grassmann manifold is given by

$$\delta(\mathbf{X}, \mathbf{Y})^2 = d_2 - \text{tr}(\mathbf{Y}^\top \mathbf{X} \mathbf{X}^\top \mathbf{Y}).$$

Thus, we have

$$\nabla \delta(\mathbf{X}, \mathbf{Y})^2 = -2(\mathbf{I}_{d_1} - \mathbf{X} \mathbf{X}^\top) \mathbf{Y} \mathbf{Y}^\top \mathbf{X},$$

where we used the transformation from the partial derivative  $\frac{\partial}{\partial \mathbf{X}} f(\mathbf{X})$  of a function  $f(\mathbf{X})$  to the manifold gradient  $\nabla f(\mathbf{X})$ :  $\nabla f(\mathbf{X}) = (\mathbf{I}_{d_1} - \mathbf{X} \mathbf{X}^\top) \frac{\partial}{\partial \mathbf{X}} f(\mathbf{X})$ . Then,  $\widehat{G}_{l,l'}$  and  $\widehat{h}_l$  defined by Eqs.(7) and (8) can be computed as

$$\begin{aligned} \widehat{G}_{l,l'} &= \frac{1}{n\sigma_l^2\sigma_{l'}^2} \sum_{i=1}^n \sum_{j=1}^d [\mathbf{F}(\mathbf{X}_i, \mathbf{C}_l)]^{(j)} [\mathbf{F}(\mathbf{X}_i, \mathbf{C}_{l'})]^{(j)} \\ &\quad \times \phi_l(\mathbf{X}_i) \phi_{l'}(\mathbf{X}_i), \\ \widehat{h}_l &= \frac{1}{n\sigma_l^2} \sum_{i=1}^n \sum_{j=1}^d \left[ (\mathbf{I}_{d_1} - \mathbf{X}_i \mathbf{X}_i^\top) \frac{\partial \mathbf{F}(\mathbf{X}_i, \mathbf{C}_l)}{\partial \mathbf{X}^{(j)}} \right]^{(j)} \phi_l(\mathbf{X}_i) \\ &\quad + \frac{1}{n\sigma_l^4} \sum_{i=1}^n \sum_{j=1}^d \left[ [\mathbf{F}(\mathbf{X}_i, \mathbf{C}_l)]^{(j)} \right]^2 \phi_l(\mathbf{X}_i), \end{aligned} \quad (9)$$

where  $[\cdot]^{(j)}$  denotes the  $j$ -th element of the vectorization of the matrix,  $d = d_1 d_2$ ,  $\mathbf{F}(\mathbf{X}, \mathbf{C}) = (\mathbf{I}_{d_1} - \mathbf{X} \mathbf{X}^\top) \mathbf{C} \mathbf{C}^\top \mathbf{X}$ . Details for the derivation of Eq.(9) are deferred to the supplementary material.

## 4 Experiments

In this section, we experimentally compare the performance of the proposed R-LSLDG clustering (R-LSLDGC) algorithm with the original mean shift (MS), the Riemannian mean shift (R-MS), the LSLDG clustering (LSLDGC), and *spectral clustering* (SC) (Ng et al., 2002) in terms of the *adjusted Rand index* (ARI) (Hubert and Arabie, 1985), which takes the maximum value 1 when the obtained clustering solution perfectly matches the ground truth.

SC requires the number of clusters to be fixed in advance. For this reason, we provide the true number of clusters only to SC. The similarity between samples used in SC is defined as  $\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / (2\tau^2))$ , where  $\tau$  is the median of  $\{\|\mathbf{x}_i - \mathbf{x}_j\|\}_{i,j=1}^n$ .

The number of basis functions in LSLDG and R-LSLDGC is set at  $b = \min(100, n)$ . The bandwidth in all algorithms and the regularization parameter in LSLDG and R-LSLDGC are chosen by 5-fold cross-validation from the 8 candidates values  $\{10^{-3}, \dots, 10\}$  at the regular interval in logarithmic scale.

All experiments were carried out using a PC equipped with two 2.60GHz Intel<sup>®</sup> Xeon<sup>®</sup> E5-2640 v3 CPUs.

**Toy Data:** Let  $\{\mathbf{X}_i \in \mathbb{R}^{d_1 \times d_2}\}_{i=1}^n$  be samples containing 3 clusters on a Grassmann manifold. For  $d_1 = 3, \dots, 7$ ,  $d_2 = 2$ , and  $n = 150$ , each sample  $\mathbf{X}_i$  is generated as

$$\mathbf{X}_i = \begin{pmatrix} \cos \tau_i & -\sin \tau_i & \mathbf{O}_{2, d_1-2} \\ \sin \tau_i & \cos \tau_i & \mathbf{O}_{d_1-2, 2} \\ \mathbf{O}_{d_1-2, 2} & \mathbf{I}_{d_1-2} \end{pmatrix} \mathbf{S} \begin{pmatrix} \cos \eta_i & -\sin \eta_i \\ \sin \eta_i & \cos \eta_i \end{pmatrix}, \quad (10)$$

where  $\mathbf{S}$  is a randomly generated element on a Grassmann manifold  $\mathcal{G}_{d_1, d_2}$  and  $\mathbf{O}_{d, d'}$  is the  $d \times d'$  null matrix. For  $N(\mu, \sigma^2)$  being the normal distribution with mean  $\mu$  and variance  $\sigma^2$ ,  $\tau_i$  and  $\eta_i$  are generated as

$$\begin{aligned} \tau_i &\sim \begin{cases} N\left(0, \frac{\pi^2}{15^2}\right) & \text{for } i = 1, \dots, \frac{n}{3}, \\ N\left(\frac{2\pi}{3}, \frac{\pi^2}{15^2}\right) & \text{for } i = \frac{n}{3} + 1, \dots, \frac{2n}{3}, \\ N\left(\frac{4\pi}{3}, \frac{\pi^2}{15^2}\right) & \text{for } i = \frac{2n}{3} + 1, \dots, n, \end{cases} \\ \eta_i &\sim N(0, \gamma^2), \end{aligned}$$

where  $\gamma = 0, \pi/8, \pi/4, \pi/2$  controls the variance of the angle of right rotation matrix. Note that when  $\gamma$  is zero, the right rotation matrix is reduced to the identity matrix. In Eq.(10), the left rotation matrix is rotation of the subspace for generating 3 clusters, while the right rotation matrix is rotation within the subspace. As plotted in Figure 3, the larger  $\gamma$  collapses the cluster structure on the Euclidean space but not on the Grassmann manifold.

The ARI values are summarized in Table 1, showing that R-LSLDGC significantly outperforms other methods when  $\gamma = \pi/4, \pi/2$ , while SC tends to perform very well when  $\gamma = 0, \pi/8$ . However, we should note that this comparison is not completely fair since SC is provided with the true number of clusters, while other methods estimate the number of clusters from data.

Compared with the plain LSLDG, R-LSLDGC performs better for large  $\gamma$ , thanks to the geometry-aware formulation. On the other hand, the plain LSLDG tends to outperform R-LSLDGC when  $\gamma = 0$ . This is caused by the difference in modeling: R-LSLDGC adopts the common-parameter formulation and thus only a single model is used for jointly estimating the gradient of all dimensions. In contrast, the plain LSLDG adopts the coordinate-wise formulation, i.e., the gradient along each dimension is estimated separately. Due to this high flexibility, when  $\gamma = 0$  (i.e., no manifold distortion is introduced), the plain LSLDG sometimes performs better than R-LSLDGC.

R-MS tends to outperform the plain MS and plain LSLDG when  $\gamma = \pi/2$ , thanks to the geometry-aware formulation. However, its performance degrades as the dimension  $d_1$  increases. In contrast, R-LSLDGC performs reliably even for large  $d_1$ , which substantiates

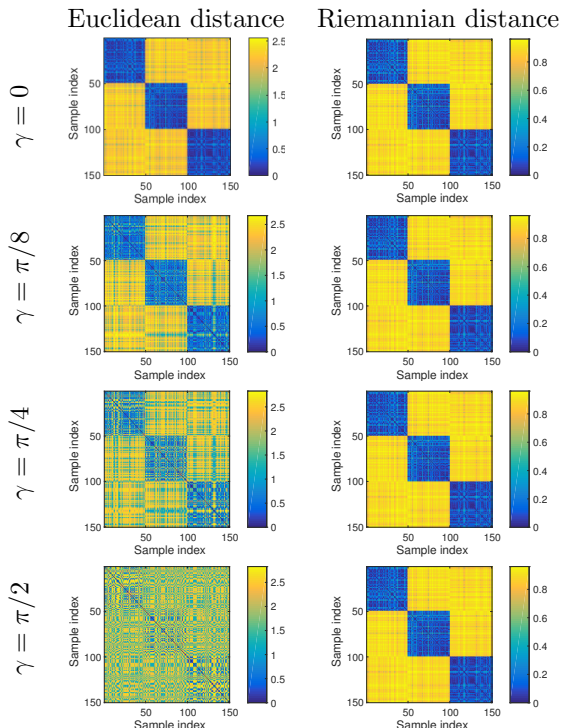


Figure 3: Distance matrices with the Euclidean and Riemannian distances when  $d_1 = 3$ ,  $d_2 = 2$ , and  $n = 150$ . A larger  $\gamma$  collapses the cluster structure on the Euclidean space, while  $\gamma$  does not influence the cluster structure on the Grassmann manifold.

the usefulness of direct gradient estimation without going through kernel density estimation on Riemannian manifolds.

Figure 4 plots the computation time of each method for different  $d_1$  averaged over  $\gamma = 0, \pi/8, \pi/4, \pi/2$ . The graph shows that SC, MS, and LSLDGC are quite fast, taking less than a few seconds. On the other hand, R-MS and R-LSLDGC take around 10 seconds, due to relatively heavy calculation of the logarithm map. This is the price we have to pay for better accuracy.

**Image Clustering:** The MNIST data set (Lecun et al., 1998) contains  $0, \dots, 9$  handwritten digits images. The images are down-sampled to  $7 \times 7$  pixels from its original size  $28 \times 28$ . Following the experimental setup in Wang et al. (2014) partially, 20 images are drawn from one digit class, and these images are vectorized and concatenated to form a matrix of size  $7^2 \times 20 = 980$ . Singular value decomposition (SVD)  $\mathbf{Z} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  is then applied to the matrix, and the top 3 left singular vectors are used as a sample on Grassmann manifold  $\mathcal{G}_{49,3}$ . Three digit classes are systematically selected to make the 3-clusters data sets of size  $n = 150$ , and we draw 50 samples from each digit class.

The ARI values are reported in Table 2, showing that the proposed R-LSLDGC outperforms other methods.

Table 1: The average and standard error of clustering accuracy measured by ARI (larger is better) over 20 runs for toy data with  $d_2 = 2$ . Bold face denotes the best and comparable methods in terms of the mean ARI according to the t-test at the significance level 5%.

$\gamma$	$d_1$	MS	LSLDGC	SC	R-MS	R-LSLDGC
0	3	.12 (.01)	.82 (.03)	<b>1.00 (.00)</b>	.51 (.06)	.81 (.04)
	4	.12 (.01)	<b>.88 (.03)</b>	<b>.94 (.04)</b>	.28 (.05)	.78 (.04)
	5	.12 (.01)	.89 (.02)	<b>1.00 (.00)</b>	.17 (.01)	.77 (.05)
	6	.12 (.01)	.90 (.03)	<b>1.00 (.00)</b>	.22 (.03)	.81 (.03)
	7	.09 (.01)	<b>.92 (.02)</b>	<b>.97 (.03)</b>	.18 (.01)	.85 (.03)
$\pi/8$	3	.10 (.01)	.62 (.03)	<b>.96 (.01)</b>	.51 (.06)	.81 (.04)
	4	.08 (.00)	.69 (.04)	<b>.95 (.03)</b>	.28 (.05)	.78 (.04)
	5	.08 (.01)	.67 (.05)	<b>.85 (.06)</b>	.17 (.01)	<b>.77 (.05)</b>
	6	.09 (.01)	.61 (.05)	<b>.92 (.05)</b>	.22 (.03)	.81 (.03)
$\pi/4$	3	.19 (.04)	.34 (.03)	.50 (.04)	.51 (.06)	<b>.81 (.04)</b>
	4	.13 (.03)	.36 (.02)	.56 (.04)	.28 (.05)	<b>.78 (.04)</b>
	5	.05 (.02)	.37 (.03)	.49 (.05)	.17 (.01)	<b>.77 (.05)</b>
	6	.04 (.00)	.31 (.03)	.35 (.04)	.22 (.03)	<b>.81 (.03)</b>
$\pi/2$	3	.13 (.02)	.21 (.02)	.09 (.01)	.51 (.06)	<b>.81 (.04)</b>
	4	.10 (.02)	.21 (.02)	.08 (.01)	.28 (.05)	<b>.78 (.04)</b>
	5	.02 (.00)	.15 (.02)	.09 (.01)	.17 (.01)	<b>.77 (.05)</b>
	6	.02 (.00)	.14 (.03)	.11 (.01)	.22 (.03)	<b>.81 (.03)</b>
7	.02 (.00)	.17 (.02)	.09 (.01)	.18 (.01)	<b>.85 (.03)</b>	

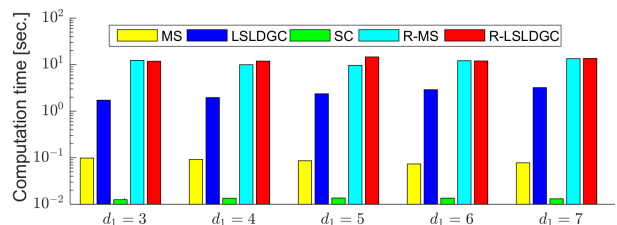


Figure 4: The average computation time of each method over all  $\gamma$  on toy data.

Figure 5 plots the average computation time over 20 runs. R-LSLDGC takes more than 100 seconds, but its computation time is still comparable to R-MS.

**Motion Segmentation:** The Hopkins 155 data set (Tron and Vidal, 2007) contains feature vectors automatically extracted from motions sequences of frame length  $F$  (see examples of the sequences in Figure 6). Under the planer scenes assumption, trajectories  $\{\mathbf{a}_i \in \mathbb{R}^{2F}\}_{i=1}^n$  from the same motion lies on a 3-dimensional subspace of  $\mathbb{R}^{2F}$ , i.e., Grassmann manifold  $\mathcal{G}_{2F,3}$  (Kanatani, 2002; Subbarao and Meer, 2009). We draw 3 trajectories  $\mathbf{a}_1$ ,  $\mathbf{a}_2$ , and  $\mathbf{a}_3$  from the same motion, and then choose the top 3 left singular vectors from the matrix  $[\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3]$  by applying SVD. Then we create data sets of size  $n = 100$  or  $n = 150$  by drawing 50 samples from each motion.

The ARI values are summarized in Table 3, showing that overall R-LSLDGC achieves higher clustering performance than other methods. Figure 7 plots the average computation time over 20 runs. On the whole,

Table 2: The average and standard error of clustering accuracy measured by ARI (larger is better) over 20 runs for the MNIST data set. Three clusters,  $c_1 = 0$ ,  $c_2 = 1$ , and  $c_3 = \{2, 3, \dots, 9\}$  are picked to construct the clustering task. Bold face denotes the best and comparable methods in terms of the mean ARI according to the t-test at the significance level 5%.

$c_3$	MS	LSDLGC	SC	R-MS	R-LSDLGC
2	.00 (.00)	.08 (.02)	.13 (.02)	.12 (.01)	<b>.45 (.02)</b>
3	.00 (.00)	.06 (.02)	.16 (.02)	.10 (.01)	<b>.41 (.02)</b>
4	.00 (.00)	.06 (.02)	.16 (.02)	.12 (.01)	<b>.39 (.02)</b>
5	.00 (.00)	.06 (.02)	.16 (.01)	.08 (.02)	<b>.48 (.02)</b>
6	.00 (.00)	.06 (.02)	.11 (.02)	.03 (.01)	<b>.47 (.02)</b>
7	.00 (.00)	.06 (.02)	.16 (.02)	.00 (.00)	<b>.42 (.02)</b>
8	.00 (.00)	.08 (.02)	.11 (.01)	.12 (.01)	<b>.35 (.02)</b>
9	.00 (.00)	.06 (.02)	.17 (.02)	.05 (.01)	<b>.42 (.02)</b>

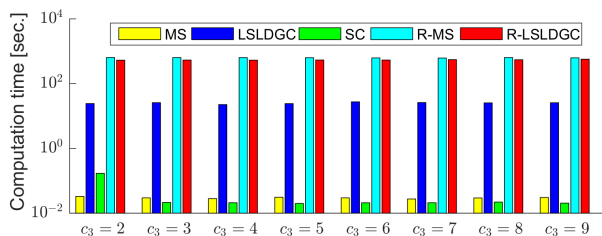


Figure 5: The average computation time of each method over 20 runs for the MNIST data set.

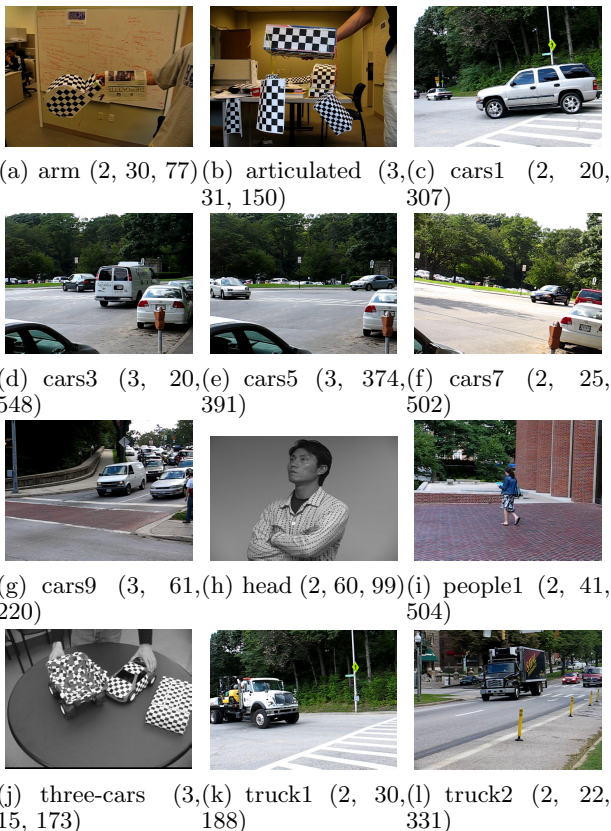


Figure 6: Examples of sequences (at time  $t = 0$ ) from the Hopkins 155 data set. The number in parentheses denotes (#Motions, #Frames, #Trajectories).

Table 3: The average and standard error of clustering accuracy measured by ARI (larger is better) over 20 runs on each sequence from the Hopkins 155 data set. The sequence names can be found in Figure 6. Bold face denotes the best and comparable methods in terms of the mean ARI according to the t-test at the significance level 5%.

Seq.	MS	LSDLGC	SC	R-MS	R-LSDLGC
(a)	.05 (.02)	.09 (.01)	.12 (.02)	.26 (.02)	<b>.65 (.04)</b>
(b)	.03 (.02)	.38 (.01)	.15 (.02)	.76 (.01)	<b>.81 (.01)</b>
(c)	.18 (.03)	.08 (.01)	.36 (.04)	.70 (.02)	<b>.75 (.02)</b>
(d)	.21 (.04)	.27 (.02)	.36 (.04)	.72 (.01)	<b>.80 (.01)</b>
(e)	.18 (.03)	.03 (.02)	.25 (.03)	.42 (.01)	<b>.45 (.01)</b>
(f)	.12 (.02)	.03 (.01)	.28 (.03)	<b>.51 (.01)</b>	.45 (.01)
(g)	.08 (.02)	.06 (.01)	.16 (.02)	.63 (.01)	<b>.70 (.01)</b>
(h)	.02 (.01)	.01 (.00)	.03 (.01)	.01 (.00)	<b>.36 (.02)</b>
(i)	.43 (.04)	.34 (.02)	.48 (.03)	.50 (.04)	<b>.59 (.01)</b>
(j)	.11 (.02)	.04 (.01)	.28 (.04)	.65 (.05)	<b>.76 (.01)</b>
(k)	.06 (.02)	.02 (.01)	.28 (.03)	<b>.76 (.04)</b>	<b>.83 (.02)</b>
(l)	.04 (.01)	.06 (.02)	.28 (.03)	.59 (.02)	<b>.64 (.01)</b>

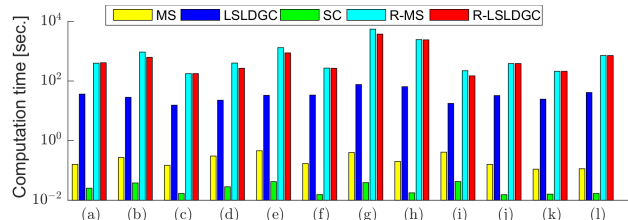


Figure 7: The average computation time of each method over 20 runs on each sequence from the Hopkins 155 data set. The sequence names can be found in Figure 6.

R-LSDLGC is not a computationally efficient method, but its computational time is still comparable to R-MS and it performs much better than R-MS.

## 5 Conclusions

Mean shift is a promising approach to mode-seeking clustering. In this paper, we extended the mean shift clustering algorithm so that Riemannian generalization and direct gradient estimation are both incorporated. Through experiments on Grassmann manifolds, we demonstrated the usefulness of the proposed method. In our future work, we will test the proposed method for other Riemannian manifolds such as Lie groups, the Stiefel manifold, and symmetric positive definite matrices. We will also investigate a computationally efficient approximation scheme for speedup.

## Acknowledgements

We thank Prof. Kazumasa Kuwada and Mr. Ikko Yamane for fruitful discussion. MA was supported by KAKENHI 26280054, HS was supported by KAKENHI 15H06103, TS was supported by KAKENHI 15J09111, and MS was supported by KAKENHI 26280054.



## References

- B. Clarke, E. Fokoue, and H. H. Zhang. *Principles and Theory for Data Mining and Machine Learning*. Springer, 2009.
- K. P. Murphy. *Machine Learning: A Probabilistic Perspective*. MIT press, 2012.
- K. Fukunaga and L. Hostetler. The estimation of the gradient of a density function, with applications in pattern recognition. *IEEE Transactions on Information Theory*, 21(1):32–40, 1975.
- Y. Cheng. Mean shift, mode seeking, and clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(8):790–799, 1995.
- D. Comaniciu and P. Meer. Mean shift: A robust approach toward feature space analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(5):603–619, 2002.
- M. A. Carreira-Perpinán. A review of mean-shift algorithms for clustering. Technical Report 1503.00687, arXiv, 2015.
- J. Wang, B. Thiesson, Y. Xu, and M. Cohen. Image and video segmentation by anisotropic kernel mean shift. In *European Conference on Computer Vision*, pages 238–249. Springer, 2004.
- W. Tao, H. Jin, and Y. Zhang. Color image segmentation based on mean shift and normalized cuts. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 37(5):1382–1389, 2007.
- D. Comaniciu, V. Ramesh, and P. Meer. Real-time tracking of non-rigid objects using mean shift. In *IEEE Conference on Computer Vision and Pattern Recognition*, volume 2, pages 142–149, 2000.
- R. T. Collins. Mean-shift blob tracking through scale space. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2, pages 234–240, 2003.
- W. M. Boothby. *An Introduction to Differentiable Manifolds and Riemannian Geometry*. Gulf Professional Publishing, 2003.
- O. Tuzel, R. Subbarao, and P. Meer. Simultaneous multiple 3D motion estimation via mode finding on Lie groups. In *Tenth IEEE International Conference on Computer Vision*, volume 1, pages 18–25, 2005.
- R. Subbarao and P. Meer. Nonlinear mean shift for clustering over analytic manifolds. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 1, pages 1168–1175, 2006.
- R. Subbarao and P. Meer. Nonlinear mean shift over Riemannian manifolds. *International Journal of Computer Vision*, 84(1):1–20, 2009.
- H. E. Cetingul and R. Vidal. Intrinsic mean shift for clustering on Stiefel and Grassmann manifolds. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1896–1902, 2009.
- D. D. Cox. A penalty method for nonparametric estimation of the logarithmic derivative of a density function. *Annals of the Institute of Statistical Mathematics*, 37(1):271–288, 1985.
- H. Sasaki, A. Hyvärinen, and M. Sugiyama. Clustering via mode seeking by direct estimation of the gradient of a log-density. In *Machine Learning and Knowledge Discovery in Databases Part III - European Conference*, pages 19–34, 2014.
- B. Pelletier. Kernel density estimation on Riemannian manifolds. *Statistics & probability letters*, 73(3):297–304, 2005.
- G. Arvanitidis, L. K. Hansen, and S. Hauberg. A locally adaptive normal distribution. In *Advances in Neural Information Processing Systems 29*, pages 4251–4259, 2016.
- P. Petersen. *Riemannian Geometry*. Springer New York Inc., New York, NY, USA, 2 edition, 2006.
- J. M. Lee. *Introduction to Smooth Manifolds*. Springer New York Inc., New York, NY, USA, 2012.
- I. Yamane, H. Sasaki, and M. Sugiyama. Regularized multitask learning for multidimensional log-density gradient estimation. *Neural Computation*, 28(7):1388–1410, 2016.
- D. Comaniciu, V. Ramesh, and P. Meer. The variable bandwidth mean shift and data-driven scale selection. In *Proceedings of Eighth IEEE International Conference on Computer Vision*, volume 1, pages 438–445, 2001.
- A. Edelman, T. A. Arias, and S. T. Smith. The geometry of algorithms with orthogonality constraints. *SIAM Journal on Matrix Analysis and Applications*, 20(2):303–353, 1998.
- A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in Neural Information Processing Systems 14*, pages 849–856, 2002.
- L. Hubert and P. Arabie. Comparing partitions. *Journal of Classification*, 2(1):193–218, 1985.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- B. Wang, Y. Hu, J. Gao, Y. Sun, and B. Yin. Low rank representation on Grassmann manifolds. In *Proceedings of 12th Asian Conference on Computer Vision*, pages 81–96, 2014.

- R. Tron and R. Vidal. A benchmark for the comparison of 3-D motion segmentation algorithms. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2007.
- K. Kanatani. Evaluation and selection of models for motion segmentation. In *Proceedings of the European Conference on Computer Vision*, pages 335–349. Springer, 2002.