

# The derivatives of the SEDS optimization cost function and constraints with respect to the learning parameters

S. Mohammad Khansari-Zadeh and Aude Billard

## I. INTRODUCTION

This technical report provides supplementary information for the optimization problems defined for Stable Estimator of Dynamical Systems (SEDS) [1]. The source code of SEDS can be downloaded from:

<http://lasa.epfl.ch/sourcecode/>

Reading of this report is *not necessary* for researchers who only want to use SEDS learning algorithm. The report is aimed at helping those persons who want to develop SEDS, or to write their own optimization program. Reading and understanding of [1] is a prerequisite for this document. All the formulations reported here are developed for SEDS models; however, they can also be used for general Gaussian Mixture Model (GMM) formulations. In the case of the latter, they should be slightly modified to consider the general form of GMM. Hopefully, the report should be clear enough to help readers in that.

To facilitate reading of the paper, a list of main variables and mathematical notations is provided in Table I. Furthermore, to have a clean summary of the final results, all the derivatives are summarized in Tables II-VI.

The remainder of this document is structured as follows. Section II gives a recap of the SEDS formulations taken from [1]. Sections III and V provide analytical formulations to compute the derivatives of MSE and Likelihood cost functions with respect to the optimization parameters, respectively. In addition, Sections IV and VI present two alternative optimization problems that automatically satisfy 4 out of 5 constraints of the original optimization problem through a change of variable. Finally, Section VII defines a proper mathematical representation of the optimization constraints, and provides the analytical derivatives of these constraints with respect to the optimization parameters.

## II. SEDS FORMULATION

Let us consider a robot motion that is defined as an autonomous Dynamical System (DS). We formulate this DS as a mixture of Gaussian functions:

$$\hat{\xi} = \hat{f}(\xi) = \sum_{k=1}^K h^k(\xi)(A^k \xi + b^k) \quad (1)$$

where

$$\begin{cases} A^k = \Sigma_{\xi\xi}^k (\Sigma_{\xi\xi}^k)^{-1} \\ b^k = \mu_{\xi}^k - A^k \mu_{\xi}^k \\ h^k(\xi) = \frac{\mathcal{P}^k(\xi) \mathcal{P}(\xi|k)}{\sum_{i=1}^K \mathcal{P}^i(\xi) \mathcal{P}(\xi|i)} \end{cases} \quad (2)$$

S.M. Khansari-Zadeh and A. Billard are with LASA Laboratory, School of Engineering, Ecole Polytechnique Federale de Lausanne, Switzerland {mohammad.khansari,Aude.Billard}@epfl.ch

TABLE I  
NOMENCLATURE

Variable	Type (size)	Description
$d$	Scalar	Dimension of DS
$\xi$	Vector ( $d$ )	Input variable, e.g. position
$\xi^*$	Vector ( $d$ )	Target point
$\dot{\xi}$	Vector ( $d$ )	Output variable, e.g. velocity
$\pi$	Scalar	Prior of the Gaussian function
$\mu$	Vector ( $d$ )	Center of the Gaussian function
$\Sigma$	Matrix ( $2d \times 2d$ )	Covariance matrix of the Gaussian fun.
$f$	Function ( $d \mapsto d$ )	Unknown original DS
$J$	Scalar	Optimization cost function
$\theta$	Structure	Optimization parameters
$L$	Matrix ( $2d \times 2d$ )	Lower triangle matrix
$A$	Matrix ( $d \times d$ )	Matrix of the linear DS
$b$	Vector ( $d$ )	Intersection point of the linear DS
$I$	Matrix	Identity matrix
$\mathbf{0}$	Vector	Zero vector
$K$	Scalar	Number of Gaussian functions
$N$	Scalar	Number of demonstrations

Notation	Description
$(\cdot)$	Estimated value of a variable
$(\cdot)^k$	Of the $k$ -th Gaussian function
$(\cdot)^T$	Transpose of a Vector/matrix
$(\cdot)^{t,n}$	The $t$ -th datapoint of the $n$ -th demonstration
$(\cdot)_i$	The $i$ -th component of a vector
$(\cdot)_{ij}$	The $(i, j)$ -th component of a matrix
$(\text{vec})_{\xi}$	Sub-vector of $\text{vec}$ with indices $1:d$
$(\text{vec})_{\dot{\xi}}$	Sub-vector of $\text{vec}$ with indices $d+1:2d$
$(\text{mat})_{\xi}$	Sub-matrix of $\text{mat}$ with indices $(1:d, 1:d)$
$(\text{mat})_{\dot{\xi}\dot{\xi}}$	Sub-matrix of $\text{mat}$ with indices $(d+1:2d, 1:d)$
$(\cdot)_{1:c, 1:c}$	A slice of a matrix with indices $(1:c, 1:c)$
$\mathbf{0}^{\{i\}}$	A zero vector with the exception that its $i$ -th component is 1
$\mathbf{0}^{\{ij\}}$	A matrix of zeros with the exception that its $(i, j)$ -th component is one.
$\mathbf{0}^{\{\bar{i}\bar{j}\}}$	A matrix of zeros with the exception that its $(i, j)$ and $(j, i)$ -th components are one.
$\text{adj}(\cdot)$	Adjugate of a matrix
$\text{tr}(\cdot)$	Trace of a matrix
$\ln(\cdot)$	The natural logarithm
$\text{Chol}(\cdot)$	Cholesky decomposition of a matrix

$$\mathcal{P}(\xi|k) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_{\xi}^k|}} e^{-\frac{1}{2}(\xi - \mu_{\xi}^k)^T (\Sigma_{\xi}^k)^{-1} (\xi - \mu_{\xi}^k)} \quad (3)$$

The unknown parameters of  $\hat{f}(\xi)$  that should be learned based on demonstrations are the priors  $\pi^k = \mathcal{P}^k$ , means  $\mu^k$  and covariance matrices  $\Sigma^k$  of the  $k = 1..K$  Gaussian functions. Given a set of  $N$  demonstrations  $\{\xi^{t,n}, \dot{\xi}^{t,n}\}_{t=0, n=1}^{T^m, N}$  of the motion, these parameters can be estimated by solving an optimization problem under the constraint of ensuring the model's global asymptotic stability. We consider two different optimization cost functions: 1) log-likelihood, and 2) Mean Square Error (MSE), which we explain next.

TABLE II  
DERIVATIVES OF THE MSE COST FUNCTION TAKEN FROM SECTION III.

$$\boldsymbol{\theta} = \{\pi^1 \dots \pi^K; \mu_\xi^1 \dots \mu_\xi^K; \Sigma_\xi^1 \dots \Sigma_\xi^K; \Sigma_{\xi\xi}^1 \dots \Sigma_{\xi\xi}^K\}$$

Cost function:  $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} (\hat{\xi}^{t,n} - \xi^{t,n})^T (\hat{\xi}^{t,n} - \xi^{t,n})$

Indices range:  $k \in 1..K, \quad i \in 1..d$

$$\frac{\partial J}{\partial \pi^k} = \frac{1}{\pi^k N} \sum_{n=1}^N \sum_{t=0}^{T^n} h^k(\xi^{t,n}) (\hat{\xi}^{t,n} - \xi^{t,n})^T (A^k \xi^{t,n} - \hat{\xi}^{t,n})$$

$$\frac{\partial J}{\partial \mu_{\xi,i}^k} = \frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} h^k(\xi^{t,n}) \left( (\xi^{t,n} - \mu_{\xi}^k)^T (\Sigma_{\xi}^k)^{-1} \mathbf{0}^{\{i\}} \right) (\hat{\xi}^{t,n} - \xi^{t,n})^T (A^k \xi^{t,n} - \hat{\xi}^{t,n})$$

$$\frac{\partial J}{\partial \Sigma_{\xi,ij}^k} = \frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} h^k(\xi^{t,n}) (\hat{\xi}^{t,n} - \xi^{t,n})^T \left( 0.5 (\xi^{t,n} - \mu_{\xi}^k)^T (\Sigma_{\xi}^k)^{-1} \mathbf{0}^{\{ij\}} (\Sigma_{\xi}^k)^{-1} (\xi^{t,n} - \mu_{\xi}^k) (A^k \xi^{t,n} - \hat{\xi}^{t,n}) \right. \\ \left. - 0.5 \text{tr} \left( (\Sigma_{\xi}^k)^{-1} \mathbf{0}^{\{ij\}} \right) (A^k \xi^{t,n} - \hat{\xi}^{t,n}) - A^k \mathbf{0}^{\{ij\}} (\Sigma_{\xi}^k)^{-1} \xi^{t,n} \right) \quad j \in 1..i$$

$$\frac{\partial J}{\partial \Sigma_{\xi\xi,ij}^k} = \frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} h^k(\xi^{t,n}) (\xi^{t,n} - \mu_{\xi}^k)^T \mathbf{0}^{\{ij\}} (\Sigma_{\xi}^k)^{-1} \xi^{t,n} \quad j \in 1..d$$

TABLE III  
DERIVATIVES OF THE ALTERNATIVE MSE COST FUNCTION TAKEN FROM SECTION IV.

$$\boldsymbol{\theta} = \{\tilde{\pi}^1 \dots \tilde{\pi}^K; \mu_{\xi}^1 \dots \mu_{\xi}^K; L_{\xi}^1 \dots L_{\xi}^K; A^1 \dots A^K\}$$

Cost function:  $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} (\hat{\xi}^{t,n} - \xi^{t,n})^T (\hat{\xi}^{t,n} - \xi^{t,n})$

Indices range:  $k \in 1..K, \quad i \in 1..d$

Change of variables:  $\tilde{\pi}^k = \ln(\pi^k), \quad L_{\xi}^k = \text{Chol}(\Sigma_{\xi}^k)$

$$\frac{\partial J}{\partial \tilde{\pi}^k} = \frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} h^k(\xi^{t,n}) (\hat{\xi}^{t,n} - \xi^{t,n})^T (A^k \xi^{t,n} - \hat{\xi}^{t,n})$$

$$\frac{\partial J}{\partial \mu_{\xi,i}^k} = \frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} h^k(\xi^{t,n}) \left( (\xi^{t,n} - \mu_{\xi}^k)^T (\Sigma_{\xi}^k)^{-1} \mathbf{0}^{\{i\}} \right) (\hat{\xi}^{t,n} - \xi^{t,n})^T (A^k \xi^{t,n} - \hat{\xi}^{t,n})$$

$$\frac{\partial J}{\partial L_{\xi,ij}^k} = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} h^k(\xi^{t,n}) \left( (\xi^{t,n} - \mu_{\xi}^k)^T (\Sigma_{\xi}^k)^{-1} \Phi (\Sigma_{\xi}^k)^{-1} (\xi^{t,n} - \mu_{\xi}^k) - \text{tr} \left( (\Sigma_{\xi}^k)^{-1} \Phi \right) \right) (\hat{\xi}^{t,n} - \xi^{t,n})^T (A^k \xi^{t,n} - \hat{\xi}^{t,n})$$

where  $\Phi = \mathbf{0}^{\{ij\}} (L^k)^T + L^k (\mathbf{0}^{\{ij\}})^T \quad j \in 1..i$

$$\frac{\partial J}{\partial A_{ij}^k} = \frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} h^k(\xi^{t,n}) (\xi^{t,n} - \mu_{\xi}^k)^T \mathbf{0}^{\{ij\}} \xi^{t,n} \quad j \in 1..d$$

Reconstruction of GMM from the optimization parameters:  $\pi^k = e^{\tilde{\pi}^k} / (\sum_{i=1}^K e^{\tilde{\pi}^i}), \quad \Sigma_{\xi}^k = L_{\xi}^k (L_{\xi}^k)^T, \quad \Sigma_{\xi\xi}^k = A^k \Sigma_{\xi}^k$

TABLE IV  
DERIVATIVES OF THE LIKELIHOOD COST FUNCTION TAKEN FROM SECTION V.

$$\boldsymbol{\theta} = \{\pi^1 \dots \pi^K; \mu_{\xi}^1 \dots \mu_{\xi}^K; \Sigma^1 \dots \Sigma^K\}$$

Cost function:  $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \log \mathcal{P}(\xi^{t,n}, \hat{\xi}^{t,n} | \boldsymbol{\theta})$

Indices range:  $k \in 1..K$

$$\frac{\partial J}{\partial \pi^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \left( \frac{\mathcal{P}(\xi^{t,n}, \hat{\xi}^{t,n} | k)}{\mathcal{P}(\xi^{t,n} | k)} - 1 \right)$$

$$\frac{\partial J}{\partial \mu_{\xi,i}^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\mathcal{P}(k) \mathcal{P}(\xi^{t,n}, \hat{\xi}^{t,n} | k)}{\mathcal{P}(\xi^{t,n} | k)} (\mathbf{0}^{\{i\}})^T [\mathbf{I} \quad (A^k)^T] (\Sigma^k)^{-1} (\xi^{t,n}; \hat{\xi}^{t,n}) - \mu^k \quad \forall i \in 1..d$$

$$\frac{\partial J}{\partial \Sigma_{ij}^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\mathcal{P}(k) \mathcal{P}(\xi^{t,n}, \hat{\xi}^{t,n} | k)}{\mathcal{P}(\xi^{t,n} | k)} \left( 0.5 (\xi^{t,n} - \mu^k)^T (\Sigma^k)^{-1} \mathbf{0}^{\{ij\}} (\Sigma^k)^{-1} (\xi^{t,n} - \mu^k) - 0.5 \text{tr} \left( (\Sigma^k)^{-1} \mathbf{0}^{\{ij\}} \right) \right. \\ \left. + (\xi^{t,n} - \mu^k)^T (\Sigma^k)^{-1} \mathcal{S}^k \right) \quad \forall i \in 1..2d, \quad j \in 1..i$$

where  $\mathcal{S}^k = \begin{bmatrix} \mathbf{0} \\ (-A^k [\mathbf{0}^{\{ij\}}]_{\xi} + [\mathbf{0}^{\{ij\}}]_{\xi\xi}) (\Sigma_{\xi}^k)^{-1} \mu_{\xi}^k \end{bmatrix}$

TABLE V  
DERIVATIVES OF THE LIKELIHOOD COST FUNCTION TAKEN FROM SECTION VI.

$$\boldsymbol{\theta} = \{\tilde{\pi}^1.. \tilde{\pi}^K; \mu_{\xi}^1.. \mu_{\xi}^K; L^1.. L^K\}$$

Cost function:  $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \log \mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n} | \boldsymbol{\theta})$

Indices range:  $k \in 1..K$

Change of variables:  $\tilde{\pi}^k = \ln(\pi^k)$ ,  $L^k = \text{Chol}(\Sigma^k)$

$$\frac{\partial J}{\partial \tilde{\pi}^k} = -\frac{e^{\tilde{\pi}^k}}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \left( \frac{\mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n} | k)}{\mathcal{P}^{t,n}} - 1 \right)$$

$$\frac{\partial J}{\partial \mu_{\xi, i}^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\mathcal{P}(k) \mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n} | k)}{\mathcal{P}^{t,n}} (\mathbf{0}^{\{i\}})^T [\mathbf{I} - (A^k)^T] (\Sigma^k)^{-1} ([\xi^{t,n}; \dot{\xi}^{t,n}] - \mu^k) \quad i \in 1..d$$

$$\frac{\partial J}{\partial L_{ij}^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\mathcal{P}(k) \mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n} | k)}{\mathcal{P}^{t,n}} \left( 0.5(\xi^{t,n} - \mu^k)^T (\Sigma^k)^{-1} \Phi (\Sigma^k)^{-1} (\xi^{t,n} - \mu^k) - 0.5 \text{tr}((\Sigma^k)^{-1} \Phi) \right. \\ \left. + (\xi^{t,n} - \mu^k)^T (\Sigma^k)^{-1} \tilde{S}^k \right) \quad i \in 1..2d, \quad j \in 1..i$$

where  $\Phi = \mathbf{0}^{\{ij\}} (L^k)^T + L^k (\mathbf{0}^{\{ij\}})^T$ ,  $\tilde{S}^k = \begin{bmatrix} \mathbf{0} \\ (-A^k \Phi_{\xi} + \Phi_{\xi \xi}) (\Sigma_{\xi}^k)^{-1} \mu_{\xi}^k \end{bmatrix}$

Reconstruction of GMM from the optimization parameters:  $\pi^k = e^{\tilde{\pi}^k} / (\sum_{i=1}^K e^{\tilde{\pi}^i})$ ,  $\Sigma^k = L^k (L^k)^T$

TABLE VI  
CONSTRAINTS FORMULATION AND THEIR DERIVATIVES FOR THE ALTERNATIVE LIKELIHOOD AND MSE COST FUNCTIONS TAKEN FROM SECTION VII.

Indices range:  $k \in 1..K$ ,  $c \in 1..d$

Constraint:  $A^k + (A^k)^T < 0$

The equivalence of the constraint used in the code:  $\mathcal{C}_{(k-1)d+c} : (-1)^{c+1} |B_{1:c, 1:c}| < 0$

$$\frac{\partial \mathcal{C}_{(k-1)d+c}}{\partial \tilde{\pi}^k} = 0 \quad (\text{valid for both the MSE and Likelihood cost functions})$$

$$\frac{\partial \mathcal{C}_{(k-1)d+c}}{\partial \mu_i^k} = 0 \quad i \in 1..d \quad (\text{valid for both the MSE and Likelihood cost functions})$$

The derivatives specific to the MSE cost function:

$$\frac{\partial \mathcal{C}_{(k-1)d+c}}{\partial L_{ij}^k} = 0 \quad i \in 1..d, \quad j \in 1..d$$

$$\frac{\partial \mathcal{C}_{(k-1)d+c}}{\partial A_{ij}^k} = (-1)^{c+1} \text{tr} \left( \text{adj}(B_{1:c, 1:c}) [\mathbf{0}^{\{ij\}}]_{1:c, 1:c} \right) \quad i \in 1..d, \quad j \in 1..d$$

The derivative specific to the Likelihood cost function:

$$\frac{\partial \mathcal{C}_{(k-1)d+c}}{\partial L_{ij}^k} = (-1)^{c+1} \text{tr} \left( \text{adj}(B_{1:c, 1:c}) \mathcal{X}_{1:c, 1:c} \right) \quad i \in 1..2d, \quad j \in 1..i$$

where  $\Phi = \mathbf{0}^{\{ij\}} (L^k)^T + L^k (\mathbf{0}^{\{ij\}})^T$ ,  $\Psi = (-A^k \Phi_{\xi} + \Phi_{\xi \xi} (\Sigma_{\xi}^k)^{-1})$ ,  $\mathcal{X} = \Psi + (\Psi)^T$

### III. MEAN SQUARE ERROR OPTIMIZATION

Mean Square Error (MSE) is a means to quantify the accuracy of estimations based on demonstrations, and it is defined as:

$$\min_{\theta} J(\theta) = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} (\hat{\xi}^{t,n} - \xi^{t,n})^T (\hat{\xi}^{t,n} - \xi^{t,n}) \quad (4)$$

subject to

$$\begin{cases} \text{(a)} & b^k = -A^k \xi^* \\ \text{(b)} & A^k + (A^k)^T < 0 \\ \text{(c)} & \Sigma_{\xi}^k > 0 \\ \text{(d)} & 0 < \pi^k \leq 1 \\ \text{(e)} & \sum_{k=1}^K \pi^k = 1 \end{cases} \quad \forall k \in 1..K \quad (5)$$

where  $\hat{\xi}^{t,n} = \hat{f}(\xi^{t,n})$  are computed from Eq. (1). The optimization parameters for this objective function are:  $\theta = \{\pi^1.. \pi^K; \mu_{\xi}^1.. \mu_{\xi}^K; \Sigma_{\xi}^1.. \Sigma_{\xi}^K; \Sigma_{\xi\xi}^1.. \Sigma_{\xi\xi}^K\}$ . Solving the above optimization requires a user to provide the derivative of the cost function w.r.t. the optimization parameters. These derivatives are provided next.

#### A. Derivatives w.r.t. Priors $\pi^k$

$$\frac{\partial J}{\partial \pi^k} = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \hat{\xi}^{t,n}} \frac{\partial \hat{\xi}^{t,n}}{\partial \pi^k} \quad \forall k \in 1..K \quad (6)$$

The partial derivatives  $\frac{\partial J}{\partial \hat{\xi}^{t,n}}$  and  $\frac{\partial \hat{\xi}^{t,n}}{\partial \pi^k}$  can be computed from Eqs. (7) and (8), respectively:

$$\frac{\partial J}{\partial \hat{\xi}^{t,n}} = \frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} (\hat{\xi}^{t,n} - \xi^{t,n})^T \quad (7)$$

$$\frac{\partial \hat{\xi}^{t,n}}{\partial \pi^k} = \frac{h^k(\xi^{t,n})}{\pi^k} (A^k \xi^{t,n} - \hat{\xi}^{t,n}) \quad (8)$$

Substituting Eqs. (7) and (8) into Eq. (6) yields:

$$\frac{\partial J}{\partial \pi^k} = \frac{1}{\pi^k N} \sum_{n=1}^N \sum_{t=0}^{T^n} h^k(\xi^{t,n}) (\hat{\xi}^{t,n} - \xi^{t,n})^T (A^k \xi^{t,n} - \hat{\xi}^{t,n}) \quad (9)$$

#### B. Derivatives w.r.t. Means $\mu_{\xi}^k$

Since  $\mu_{\xi}^k$  is a  $d$ -dimensional vector, we need to compute the derivative w.r.t. each component of  $\mu_{\xi}^k$  separately:

$$\frac{\partial J}{\partial \mu_{\xi,i}^k} = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \hat{\xi}^{t,n}} \frac{\partial \hat{\xi}^{t,n}}{\partial \mu_{\xi,i}^k} \quad \forall i \in 1..d, k = 1..K \quad (10)$$

The partial derivative  $\frac{\partial J}{\partial \hat{\xi}^{t,n}}$  is given by Eq. (7), and  $\frac{\partial \hat{\xi}^{t,n}}{\partial \mu_{\xi,i}^k}$  is:

$$\frac{\partial \hat{\xi}^{t,n}}{\partial \mu_{\xi,i}^k} = h^k(\xi^{t,n}) \left( (\xi^{t,n} - \mu_{\xi}^k)^T (\Sigma_{\xi}^k)^{-1} \mathbf{0}^{\{i\}} \right) (A^k \xi^{t,n} - \hat{\xi}^{t,n}) \quad (11)$$

where  $\mathbf{0}^{\{i\}}$  has the dimension of  $d$ .

#### C. Derivatives w.r.t. Means $\mu_{\xi}^k$

By substituting directly the constraint Eq. (5)-(a) into Eq. (1), the partial derivative  $\frac{\partial \hat{\xi}^{t,n}}{\partial \mu_{\xi,i}^k}$  is always zero because  $\hat{f}(\xi)$  no longer depends on  $\mu_{\xi}^k$ . Therefore,  $\mu_{\xi,i}^k$  can be dropped from the list of the optimization parameters. In fact, at each iteration  $\mu_{\xi}^k$  is exploited to satisfy this constraint, and its value can be directly computed from Eq. (5)-(a).

#### D. Derivatives w.r.t. $\Sigma_{\xi}^k$

Since  $\Sigma_{\xi}^k$  is a  $d \times d$  matrix, we will compute the derivative w.r.t. its each component separately. Since  $\Sigma_{\xi}^k$  is a symmetric matrix, we compute the derivatives only for the components on the lower triangle matrix.

$$\frac{\partial J}{\partial \Sigma_{\xi,ij}^k} = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \hat{\xi}^{t,n}} \frac{\partial \hat{\xi}^{t,n}}{\partial \Sigma_{\xi,ij}^k} \quad \begin{cases} \forall i \in 1..d \\ \forall j \in 1..i \\ \forall k \in 1..K \end{cases} \quad (12)$$

The partial derivative  $\partial \hat{\xi}^{t,n} / \partial \Sigma_{\xi,ij}^k$  is:

$$\begin{aligned} \frac{\partial \hat{\xi}^{t,n}}{\partial \Sigma_{\xi,ij}^k} &= -h^k(\xi^{t,n}) A^k \mathbf{0}^{\{\bar{i}\bar{j}\}} (\Sigma_{\xi}^k)^{-1} \xi^{t,n} + \\ &\frac{h^k(\xi^{t,n})}{2} \left( (\xi^{t,n} - \mu_{\xi}^k)^T (\Sigma_{\xi}^k)^{-1} \mathbf{0}^{\{\bar{i}\bar{j}\}} (\Sigma_{\xi}^k)^{-1} (\xi^{t,n} - \mu_{\xi}^k) \right. \\ &\quad \left. - \text{tr}((\Sigma_{\xi}^k)^{-1} \mathbf{0}^{\{\bar{i}\bar{j}\}}) \right) (A^k \xi^{t,n} - \hat{\xi}^{t,n}) \end{aligned} \quad (13)$$

where  $\mathbf{0}^{\{\bar{i}\bar{j}\}}$  has the dimension of  $d \times d$ .

#### E. Derivatives w.r.t. $\Sigma_{\xi\xi}^k$

The partial derivatives of the cost function w.r.t. the components of  $\Sigma_{\xi\xi}^k$  are

$$\frac{\partial J}{\partial \Sigma_{\xi\xi,ij}^k} = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \hat{\xi}^{t,n}} \frac{\partial \hat{\xi}^{t,n}}{\partial \Sigma_{\xi\xi,ij}^k} \quad \begin{cases} \forall i \in 1..d \\ \forall j \in 1..d \\ \forall k \in 1..K \end{cases} \quad (14)$$

The partial derivative  $\partial \hat{\xi}^{t,n} / \partial \Sigma_{\xi\xi,ij}^k$  is:

$$\frac{\partial \hat{\xi}^{t,n}}{\partial \Sigma_{\xi\xi,ij}^k} = h^k(\xi^{t,n}) \mathbf{0}^{\{ij\}} (\Sigma_{\xi}^k)^{-1} \xi^{t,n} \quad (15)$$

where  $\mathbf{0}^{\{ij\}}$  has the dimension of  $d \times d$ .

#### IV. ALTERNATIVE MSE OPTIMIZATION

Though the MSE optimization provided in Section III is sufficient to estimate a stable DS, its performance can be significantly increased through a change of optimization parameters. Let us define:

$$\begin{cases} \hat{\pi}^k = \ln(\pi^k) \\ L_\xi^k = \text{Chol}(\Sigma_\xi^k) \end{cases} \quad (16)$$

where  $L_\xi^k$  is a  $d \times d$  lower triangle matrix. Since  $\Sigma_\xi^k$  are positive definite matrix, their Cholesky decomposition  $L_\xi^k$  always exist. Furthermore, as it was pointed out before, by substituting Eq. (5)-(a) into Eq. (1), we can define the evolution of motion with:

$$\hat{\xi} = \hat{f}(\xi) = \sum_{k=1}^K h^k(\xi) A^k (\xi - \xi^*) \quad (17)$$

Considering Eqs. (16) and (17) and defining the optimization parameters to be  $\theta = \{\hat{\pi}^1 \dots \hat{\pi}^K; \mu_\xi^1 \dots \mu_\xi^K; L_\xi^1 \dots L_\xi^K; A^1 \dots A^K\}$ , the alternative MSE optimization can be expressed as:

$$\min_{\theta} J(\theta) = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} (\hat{\xi}^{t,n} - \xi^{t,n})^T (\hat{\xi}^{t,n} - \xi^{t,n}) \quad (18)$$

subject to

$$A^k + (A^k)^T < 0 \quad \forall k \in 1..K \quad (19)$$

where  $\hat{\xi}^{t,n} = \hat{f}(\xi^{t,n})$  are computed from Eq. (17). Once the optimization finished, the parameters of GMM can be reconstructed as follows:

$$\begin{cases} \pi^k = e^{\hat{\pi}^k} / (\sum_{i=1}^K e^{\hat{\pi}^i}) \\ \Sigma_\xi^k = L_\xi^k (L_\xi^k)^T \\ \Sigma_{\xi\xi}^k = A^k \Sigma_\xi^k \end{cases} \quad (20)$$

In fact the proposed change of parameters allows us to automatically satisfy the last three optimization constraints of Eq. (5). The first constraint of Eq. (5) is also removed since it is directly considered in Eq. (17). The derivatives of the new optimization problem are provided in the following subsections.

##### A. Derivatives w.r.t. Priors $\pi^k$

$$\frac{\partial J}{\partial \hat{\pi}^k} = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \hat{\xi}^{t,n}} \frac{\partial \hat{\xi}^{t,n}}{\partial \pi^k} \frac{\partial \pi^k}{\partial \hat{\pi}^k} \quad \forall k \in 1..K \quad (21)$$

The partial derivatives  $\partial J / \partial \hat{\xi}^{t,n}$  and  $\partial \hat{\xi}^{t,n} / \partial \pi^k$  are given by Eqs. (7) and (8), and the derivative  $\partial \pi^k / \partial \hat{\pi}^k$  is simply:

$$\frac{\partial \pi^k}{\partial \hat{\pi}^k} = e^{\hat{\pi}^k} \quad (22)$$

##### B. Derivatives w.r.t. Means $\mu_\xi^k$

These derivative can be similarly computed from Eq. (10).

##### C. Derivatives w.r.t. $L^k$

$L_\xi^k$  is a  $d \times d$  lower triangle matrix. The partial derivatives of the cost function w.r.t. its parameters are:

$$\frac{\partial J}{\partial L_{\xi,ij}^k} = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \hat{\xi}^{t,n}} \frac{\partial \hat{\xi}^{t,n}}{\partial L_{\xi,ij}^k} \quad \begin{cases} \forall i \in 1..d \\ \forall j \in 1..i \\ \forall k \in 1..K \end{cases} \quad (23)$$

The partial derivative  $\partial \hat{\xi}^{t,n} / \partial L_{\xi,ij}^k$  is:

$$\frac{\partial \hat{\xi}^{t,n}}{\partial L_{\xi,ij}^k} = \frac{h^k(\xi^{t,n})}{2} \left( (\xi^{t,n} - \mu_\xi^k)^T (\Sigma_\xi^k)^{-1} \Phi (\Sigma_\xi^k)^{-1} (\xi^{t,n} - \mu_\xi^k) - \text{tr} \left( (\Sigma_\xi^k)^{-1} \Phi \right) \right) (A^k \xi^{t,n} - \hat{\xi}^{t,n}) \quad (24)$$

where  $\Phi = \mathbf{0}^{\{ij\}} (L_\xi^k)^T + L_\xi^k (\mathbf{0}^{\{ij\}})^T$ , and has the dimension of  $d \times d$ .

##### D. Derivatives w.r.t. $A^k$

The partial derivatives of the cost function w.r.t. the components of  $A^k$  are

$$\frac{\partial J}{\partial A_{ij}^k} = \frac{1}{2N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \hat{\xi}^{t,n}} \frac{\partial \hat{\xi}^{t,n}}{\partial A_{ij}^k} \quad \begin{cases} \forall i \in 1..d \\ \forall j \in 1..d \\ \forall k \in 1..K \end{cases} \quad (25)$$

The partial derivative  $\partial \hat{\xi}^{t,n} / \partial A_{ij}^k$  is:

$$\frac{\partial \hat{\xi}^{t,n}}{\partial A_{ij}^k} = h^k(\xi^{t,n}) \mathbf{0}^{\{ij\}} \xi^{t,n} \quad (26)$$

where  $\mathbf{0}^{\{ij\}}$  has the dimension of  $d \times d$ .

#### V. LIKELIHOOD OPTIMIZATION

The likelihood optimization is defined as:

$$\min_{\theta} J(\theta) = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \log \mathcal{P}(\xi^{t,n}, \hat{\xi}^{t,n} | \theta) \quad (27)$$

subject to the same constrains as given by Eq. (5). In Eq. (27),  $\mathcal{P}(\xi^{t,n}, \hat{\xi}^{t,n} | \theta)$  is computed from:

$$\mathcal{P}(\xi^{t,n}, \hat{\xi}^{t,n}; \theta) = \sum_{k=1}^K \mathcal{P}(k) \mathcal{P}(\xi^{t,n}, \hat{\xi}^{t,n} | k) \quad \begin{cases} \forall n \in 1..N \\ t \in 0..T^n \end{cases} \quad (28)$$

where  $\mathcal{P}(k) = \pi^k / (\sum_{i=1}^K \pi^i)$  is the prior and  $\mathcal{P}(\xi^{t,n}, \hat{\xi}^{t,n} | k)$  is the conditional probability density function given by:

$$\mathcal{P}(\xi^{t,n}, \hat{\xi}^{t,n} | k) = \mathcal{N}(\xi^{t,n}, \hat{\xi}^{t,n}; \mu^k, \Sigma^k) = \frac{1}{\sqrt{(2\pi)^{2d} |\Sigma^k|}} e^{-\frac{1}{2} ([\xi^{t,n}, \hat{\xi}^{t,n}] - \mu^k)^T (\Sigma^k)^{-1} ([\xi^{t,n}, \hat{\xi}^{t,n}] - \mu^k)} \quad (29)$$

The optimization parameters for this objective function are:  $\theta = \{\pi^1 \dots \pi^K; \mu_\xi^1 \dots \mu_\xi^K; \Sigma^1 \dots \Sigma^K\}$ . Next we compute these derivatives with respect to  $\theta$ .

### A. Derivatives w.r.t. Priors $\pi^k$

$$\frac{\partial J}{\partial \pi^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \mathcal{P}^{t,n}} \frac{\partial \mathcal{P}^{t,n}}{\partial \pi^k} \quad \forall k \in 1..K \quad (30)$$

where for simplicity we shorten the notation  $\mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n}; \theta)$  to  $\mathcal{P}^{t,n}$ . The partial derivatives  $\frac{\partial J}{\partial \mathcal{P}^{t,n}}$  and  $\frac{\partial \mathcal{P}^{t,n}}{\partial \pi^k}$  can be computed from Eqs. (31) and (32), respectively:

$$\frac{\partial J}{\partial \mathcal{P}^{t,n}} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{1}{\mathcal{P}^{t,n}} \quad (31)$$

$$\frac{\partial \mathcal{P}^{t,n}}{\partial \pi^k} = \mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n}|k) - \mathcal{P}^{t,n} \quad (32)$$

Substituting Eqs. (31) and (32) into Eq. (30) yields:

$$\frac{\partial J}{\partial \pi^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \left( \frac{\mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n}|k)}{\mathcal{P}^{t,n}} - 1 \right) \quad (33)$$

### B. Derivatives w.r.t. Means $\mu_{\xi}^k$

Special attention should be considered in computing derivatives with respect to  $\mu_{\xi}^k$ . As it is already discussed in Section III, there is a direct relation between  $\mu_{\xi}^k$  and  $\mu_{\xi}^k$  through the constraint Eq. (5)-(a). By substituting the corresponding value of  $\mu_{\xi}^k$  into the cost function given by Eq. (27), the optimization no longer depends on  $\mu_{\xi}^k$ . Hence, we can drop  $\mu_{\xi}^k$  from the optimization parameters and the constraint Eq. (5)-(a) is always satisfied. However, this substitution should be considered when computing the derivatives with respect to  $\mu_{\xi}^k$ :

$$\frac{\partial J}{\partial \mu_{\xi,i}^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \mathcal{P}^{t,n}} \left( \frac{\partial \mathcal{P}^{t,n}}{\partial \mu_{\xi,i}^k} + \sum_{j=1}^d \frac{\partial \mathcal{P}^{t,n}}{\partial \mu_{\xi,j}^k} \frac{\partial \mu_{\xi,j}^k}{\partial \mu_{\xi,i}^k} \right) \quad (34)$$

The partial derivative  $\frac{\partial J}{\partial \mathcal{P}^{t,n}}$  is given by Eq. (31), and  $\partial \mathcal{P}^{t,n} / \partial \mu_{\xi,i}^k$  is:

$$\frac{\partial \mathcal{P}^{t,n}}{\partial \mu_{\xi,i}^k} = (\mathbf{0}^{\{i\}})^T (\Sigma^k)^{-1} ([\xi^{t,n}; \dot{\xi}^{t,n}] - \mu^k) \mathcal{P}(k) \quad \forall i \in 1..d \quad (35)$$

where  $\mathbf{0}^{\{i\}}$  is a vector of dimension  $2d$ .

The partial derivative  $\partial \mathcal{P}^{t,n} / \partial \mu_{\xi,j}^k$  can be computed similarly to Eq. (34); however by replacing  $\mathbf{0}^{\{i\}}$  with  $\mathbf{0}^{\{i+d\}}$ .

The derivative  $\frac{\partial \mu_{\xi,j}^k}{\partial \mu_{\xi,i}^k}$  can be computed by differentiating Eq. (5)-(a) with respect to  $\mu_{\xi,i}^k$ :

$$\frac{\partial \mu_{\xi,j}^k}{\partial \mu_{\xi,i}^k} = A_{ji}^k \quad \forall i \in 1..d, j \in 1..d \quad (36)$$

Thanks to matrix multiplication, we can significantly simplify the multiplications by substituting Eqs. (35), (36), and (31) into Eq. (34), and compute  $\frac{\partial J}{\partial \mu_{\xi}^k}$ :

$$\frac{\partial J}{\partial \mu_{\xi}^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\mathcal{P}(k) \mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n}|k)}{\mathcal{P}^{t,n}} [\mathbf{I} \quad (A^k)^T] (\Sigma^k)^{-1} ([\xi^{t,n}; \dot{\xi}^{t,n}] - \mu^k) \quad \forall i \in 1..d \quad (37)$$

where  $\mathbf{I}$  has the dimension of  $d \times d$ . Note that  $\frac{\partial J}{\partial \mu_{\xi}^k}$  is now a vector of dimension  $d$ , and each  $\frac{\partial J}{\partial \mu_{\xi,i}^k}$  is in fact one element of this vector.

### C. Derivatives w.r.t. Means $\mu_{\xi}^k$

By substituting directly the constraint Eq. (5)-(a) into Eq. (1), the partial derivative  $\frac{\partial \mathcal{P}^{t,n}}{\partial \mu_{\xi,i}^k}$  is always zero because  $\hat{f}(\xi)$  no longer depends on  $\mu_{\xi}^k$ . Therefore,  $\mu_{\xi}^k$  can be dropped from the list of the optimization parameters. For more information see Section V-B.

### D. Derivatives w.r.t. $\Sigma^k$

Similar to Section V-B, we need to consider the effect of substitution of  $\mu_{\xi}^k$  when computing the derivatives of  $\Sigma^k$ . All  $\Sigma^k$  are  $2d \times 2d$  symmetric matrices, hence we compute the derivatives only for the components on the lower triangle matrix.

$$\frac{\partial J}{\partial \Sigma_{ij}^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \mathcal{P}^{t,n}} \left( \frac{\partial \mathcal{P}^{t,n}}{\partial \Sigma_{ij}^k} + \frac{\partial \mathcal{P}^{t,n}}{\partial \Sigma_{ij}^k} \right) \left. \vphantom{\frac{\partial J}{\partial \Sigma_{ij}^k}} \right|_{\mu_{\xi}^k} \quad \left\{ \begin{array}{l} \forall i \in 1..2d \\ \forall j \in 1..i \\ \forall k \in 1..K \end{array} \right. \quad (38)$$

where  $\left. \frac{\partial \mathcal{P}^{t,n}}{\partial \Sigma_{ij}^k} \right|_{\mu_{\xi}^k}$  corresponds to the portion of derivatives due to the effect of  $\mu_{\xi}^k$ , and can be computed from:

$$\left. \frac{\partial \mathcal{P}^{t,n}}{\partial \Sigma_{ij}^k} \right|_{\mu_{\xi}^k} = \sum_{l=1}^d \sum_{m=1}^d \frac{\partial \mathcal{P}^{t,n}}{\partial \mu_{\xi,l}^k} \frac{\partial \mu_{\xi,l}^k}{\partial A_{lm}^k} \frac{\partial A_{lm}^k}{\partial \Sigma_{ij}^k} \quad (39)$$

The partial derivative  $\partial \mathcal{P}^{t,n} / \partial \Sigma_{ij}^k$  is:

$$\frac{\partial \mathcal{P}^{t,n}}{\partial \Sigma_{ij}^k} = 0.5 \left( (\xi^{t,n} - \mu^k)^T (\Sigma^k)^{-1} \mathbf{0}^{\{\overline{ij}\}} (\Sigma^k)^{-1} (\xi^{t,n} - \mu^k) - \text{tr}((\Sigma^k)^{-1} \mathbf{0}^{\{\overline{ij}\}}) \right) \mathcal{P}(k) \mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n}|k) \quad (40)$$

where  $\mathbf{0}^{\{\overline{ij}\}}$  has the dimension of  $2d \times 2d$ .

The partial derivative  $\left. \frac{\partial \mathcal{P}^{t,n}}{\partial \Sigma_{ij}^k} \right|_{\mu_{\xi}^k}$  could be significantly simplified if it is computed in the matrix form (because we can drop the both summations on  $l$  and  $m$ ):

$$\left. \frac{\partial \mathcal{P}^{t,n}}{\partial \Sigma_{ij}^k} \right|_{\mu_{\xi}^k} = \mathcal{P}(k) \mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n}|k) (\xi^{t,n} - \mu^k)^T (\Sigma^k)^{-1} \mathcal{S}^k \quad (41)$$

where  $\mathcal{S}^k$  is a vector of dimension  $2d$  and is equal to:

$$\mathcal{S}^k = \left[ \begin{array}{c} \mathbf{0} \\ \left( -A^k [\mathbf{0}^{\{\overline{ij}\}}]_{\xi} + [\mathbf{0}^{\{\overline{ij}\}}]_{\dot{\xi}} \right) (\Sigma_{\xi}^k)^{-1} \mu_{\xi}^k \end{array} \right] \quad (42)$$



In Eq. (42),  $\mathbf{0}$  is a zero column vector of dimension  $d$ , and  $\mathbf{0}_{\xi}^{\{ij\}}$  and  $\mathbf{0}_{\xi\xi}^{\{ij\}}$  are partitions of  $\mathbf{0}^{\{ij\}}$ . Finally, by substituting Eqs. (40), (41), and (31) into Eq. (38) we have:

$$\frac{\partial J}{\partial \Sigma_{ij}^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\mathcal{P}(k) \mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n}|k)}{\mathcal{P}^{t,n}} \left( 0.5(\xi^{t,n} - \mu^k)^T (\Sigma^k)^{-1} \mathbf{0}^{\{ij\}} (\Sigma^k)^{-1} (\xi^{t,n} - \mu^k) - 0.5 \text{tr}((\Sigma^k)^{-1} \mathbf{0}^{\{ij\}}) + (\xi^{t,n} - \mu^k)^T (\Sigma^k)^{-1} \mathcal{S}^k \right) \quad (43)$$

## VI. ALTERNATIVE LIKELIHOOD OPTIMIZATION

Similarly to Section IV, we can define an alternative likelihood optimization so that 4 out of 5 optimization constraints can be automatically satisfied through a change of variable:

$$\begin{cases} \tilde{\pi}^k = \ln(\pi^k) \\ L^k = \text{Chol}(\Sigma^k) \end{cases} \quad (44)$$

where  $L^k$  are  $2d \times 2d$  lower triangle matrices. Since  $\Sigma^k$  are positive definite matrices, their Cholesky decomposition always exist. The alternative likelihood optimization can be expressed as:

$$\min_{\theta} J(\theta) = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \log \mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n}|\theta) \quad (45)$$

subject to

$$A^k + (A^k)^T < 0 \quad \forall k \in 1..K \quad (46)$$

where  $\theta = \{\tilde{\pi}^1 .. \tilde{\pi}^K; \mu_{\xi}^1 .. \mu_{\xi}^K; L^1 .. L^K\}$ . Once the optimization finished, the parameters of GMM can be reconstructed as follows:

$$\begin{cases} \pi^k = e^{\tilde{\pi}^k} / (\sum_{i=1}^K e^{\tilde{\pi}^i}) \\ \Sigma^k = L^k (L^k)^T \end{cases} \quad (47)$$

In fact the proposed change of parameters allows us to automatically satisfy the last three optimization constraints of Eq. (5). The first constraint of Eq. (5) is also removed since it is directly considered in Eq. (17). The derivatives of the new optimization problem are provided in the following subsections.

### A. Derivatives w.r.t. Priors $\pi^k$

$$\frac{\partial J}{\partial \tilde{\pi}^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \mathcal{P}^{t,n}} \frac{\partial \mathcal{P}^{t,n}}{\partial \pi^k} \frac{\partial \pi^k}{\partial \tilde{\pi}^k} \quad \forall k \in 1..K \quad (48)$$

The partial derivatives  $\partial J / \partial \mathcal{P}^{t,n}$ ,  $\partial \mathcal{P}^{t,n} / \partial \pi^k$  and  $\partial \pi^k / \partial \tilde{\pi}^k$  are given by Eqs. (31), (32), and (22), respectively.

### B. Derivatives w.r.t. Means $\mu_{\xi}^k$

These derivative can be similarly computed from Eq. (34).

### C. Derivatives w.r.t. $L^k$

$L^k$  is a  $2d \times 2d$  lower triangle matrix. The partial derivatives of the cost function with respect to the optimization parameters are:

$$\frac{\partial J}{\partial L_{ij}^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\partial J}{\partial \mathcal{P}^{t,n}} \frac{\partial \mathcal{P}^{t,n}}{\partial L_{ij}^k} \quad \begin{cases} \forall i \in 1..2d \\ \forall j \in 1..i \\ \forall k \in 1..K \end{cases} \quad (49)$$

The partial derivative  $\partial \mathcal{P}^{t,n} / \partial L_{ij}^k$  is:

$$\frac{\partial J}{\partial L_{ij}^k} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T^n} \frac{\mathcal{P}(k) \mathcal{P}(\xi^{t,n}, \dot{\xi}^{t,n}|k)}{\mathcal{P}^{t,n}} \left( 0.5(\xi^{t,n} - \mu^k)^T (\Sigma^k)^{-1} \Phi (\Sigma^k)^{-1} (\xi^{t,n} - \mu^k) - 0.5 \text{tr}((\Sigma^k)^{-1} \Phi) + (\xi^{t,n} - \mu^k)^T (\Sigma^k)^{-1} \tilde{\mathcal{S}}^k \right) \quad (50)$$

where  $\Phi = \mathbf{0}^{\{ij\}} (L^k)^T + L^k (\mathbf{0}^{\{ij\}})^T$ , and has the dimension of  $2d \times 2d$ . The  $2d$  dimension vector  $\tilde{\mathcal{S}}^k$  is:

$$\tilde{\mathcal{S}}^k = \begin{bmatrix} \mathbf{0} \\ (-A^k \Phi_{\xi} + \Phi_{\xi\xi}) (\Sigma_{\xi}^k)^{-1} \mu_{\xi}^k \end{bmatrix} \quad (51)$$

where  $\mathbf{0}$  is a zero column vector of dimension  $d$ .

## VII. OPTIMIZATION CONSTRAINTS AND THEIR DERIVATIVE

In this section we provide formulations for the optimization problems defined in Sections IV and VI, where the only constraint is the negative definiteness of matrices  $A^k$ . To ensure this constraint, we first need to define a method to mathematically determine whether a matrix is negative definite. There are several ways to ensure whether a symmetric matrix  $B$  is negative definite, among which the two most famous ones are 1) all eigenvalues of  $B$  are strictly negative, 2) using Sylvester's criterion. In our work, we use Sylvester's criterion because it provides us with an analytical formulation to verify negative definiteness (compared to computing eigenvalues which is an iterative procedure).

Sylvester's criterion states that a Hermitian matrix  $B$  is negative-definite if and only if the determinant of all  $i$ -th order leading principal minors<sup>1</sup> are negative if  $i$  is odd and positive if  $i$  is even [2]. Each  $d \times d$  symmetric matrix has  $d$  principal minors. By defining  $B^k = A^k + (A^k)^T$ , the optimization constraint given by Eq. (46) is equal to:

$$\mathcal{C}_{(k-1)d+c} : (-1)^{c+1} |B_{1:c,1:c}| < 0 \quad \begin{cases} \forall c \in 1..d \\ \forall k \in 1..K \end{cases} \quad (52)$$

where we use  $\mathcal{C}_{(k-1)d+c}$  to refer to the  $((k-1)d+c)$ -th constraint. Thus for a GMM model composed of  $K$  Gaussian functions, there are  $K \times d$  constraints that should be satisfied

<sup>1</sup>The  $i$ -th principal minor of a  $d \times d$  symmetric matrix  $B$  is a quadratic upper-left part of  $B$ , which consists of matrix elements in rows and columns from 1 to  $i$ .

during the optimization. The derivative of these constraints with respect to  $\pi^k$  and  $\mu^k$  are always zero, irrespective of which cost function is used:

$$\frac{\partial \mathcal{C}_{(k-1)d+c}}{\partial \tilde{\pi}^k} = 0 \quad \begin{cases} \forall c \in 1..d \\ \forall k \in 1..K \end{cases} \quad (53)$$

$$\frac{\partial \mathcal{C}_{(k-1)d+c}}{\partial \mu_i^k} = 0 \quad \begin{cases} \forall c \in 1..d \\ \forall i \in 1..2d \\ \forall k \in 1..K \end{cases} \quad (54)$$

For the MSE optimization defined by Eq. (18) we have:

$$\frac{\partial \mathcal{C}_{(k-1)d+c}}{\partial L_{ij}^k} = 0 \quad \begin{cases} \forall c \in 1..d \\ \forall i \in 1..d \\ \forall j \in 1..i \\ \forall k \in 1..K \end{cases} \quad (55)$$

$$\frac{\partial \mathcal{C}_{(k-1)d+c}}{\partial A_{ij}^k} = (-1)^{c+1} \text{tr} \left( \text{adj}(B_{1:c,1:c}) [\mathbf{0}^{\{\bar{ij}\}}]_{1:c,1:c} \right) \begin{cases} \forall c \in 1..d \\ \forall i \in 1..d \\ \forall j \in 1..d \\ \forall k \in 1..K \end{cases} \quad (56)$$

where  $\mathbf{0}^{\{\bar{ij}\}}$  has the dimension of  $d \times d$ . For the likelihood optimization defined by Eq. (45) we have:

$$\frac{\partial \mathcal{C}_{(k-1)d+c}}{\partial L_{ij}^k} = (-1)^{c+1} \text{tr} \left( \text{adj}(B_{1:c,1:c}) \mathcal{X}_{1:c,1:c} \right) \begin{cases} \forall c \in 1..d \\ \forall i \in 1..2d \\ \forall j \in 1..i \\ \forall k \in 1..K \end{cases} \quad (57)$$

where  $\mathcal{X}$  is a  $d \times d$  symmetric matrix defined by:

$$\Phi = \mathbf{0}^{\{ij\}} (L^k)^T + L^k (\mathbf{0}^{\{ij\}})^T \quad (58)$$

$$\Psi = (-A^k \Phi_\xi + \Phi_{\xi\xi}) (\Sigma_\xi)^{-1} \quad (59)$$

$$\mathcal{X} = \Psi + (\Psi)^T \quad (60)$$

where  $\Phi$  is a  $2d \times 2d$  matrix.

## REFERENCES

- [1] S.-M. Khansari-Zadeh and A. Billard, "Imitation learning of globally stable non-linear point-to-point robot motions using nonlinear programming," in *Proceeding of the 2010 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2010, pp. 2676–2683.
- [2] G. T. Gilbert, "Positive definite matrices and Sylvester's criterion," *The American Mathematical Monthly (Mathematical Association of America)*, vol. 98(1), pp. 44–46, 1991.