

Closed-Form Information-Theoretic Roughness Measures for Mixture Densities

Uwe D. Hanebeck

Intelligent Sensor-Actuator-Systems Laboratory (ISAS)
 Institute for Anthropomatics and Robotics
 Karlsruhe Institute of Technology (KIT), Germany
 Uwe.Hanebeck@kit.edu

Abstract—We calculate the smoothest mixture density under a variety of prescribed specifications. This includes constraints on certain moments, specifications on density values and/or its derivatives, and prescribed probability masses in certain regions. As a roughness measure, we use Fisher Information (FI) in the space of mixtures \mathcal{M} . For mixtures, FI cannot be calculated in closed form. We define the space \mathcal{R} of root mixtures (RMs) living on the Hilbert sphere. A transformation of FI to \mathcal{R} admits a closed-form solution and yields the desired result in \mathcal{M} . This naturally leads to a tandem processing with two density representations maintained simultaneously in \mathcal{R} and \mathcal{M} . FI is calculated in RM space \mathcal{R} while the constraints are evaluated in mixture space \mathcal{M} .

Keywords—Mixture density, Gaussian mixture, information measure, Fisher information, smooth density, roughness measure, square-root densities, closed-form expression.

I. INTRODUCTION

A. Context

In many applications we are given some specifications on a probability density function (pdf) and want to find a pdf that adds as little information as possible to what is specified. Intuitively, this corresponds to the smoothest pdf under the specifications. The pdf is selected from a certain class of parametric or non-parametric pdfs. The goal is to find the smoothest pdf from this class in an information theoretic sense meeting the specifications.

Simple specifications include prescribing certain moments or confining the pdf to be nonzero on certain intervals and zero in others. More complex specifications assume certain realization/density pairs $(\underline{x}_i, f(\underline{x}_i))$, $i = 1:N$ ¹ or corresponding tolerance bands. Typically, the constraints on the pdf resulting from the specifications do not fully characterize the pdf. Enforcing smoothness can be viewed as a *regularizer* for underdetermined optimization.

W.l.o.g. we can assume the existence of an *unknown* underlying pdf $\tilde{f}(\underline{x})$ with $\underline{x} \in \mathbb{R}^D$, D the number of dimensions, of which we only have a set \mathcal{S} of specifications, see Fig. 1. The set \mathcal{S} is not sufficient to fully describe $\tilde{f}(\underline{x})$. When reconstructing an approximation $f(\underline{x})$ of $\tilde{f}(\underline{x})$, we do not want to artificially add unwarranted information.² We desire the least-informative pdf given the specification.

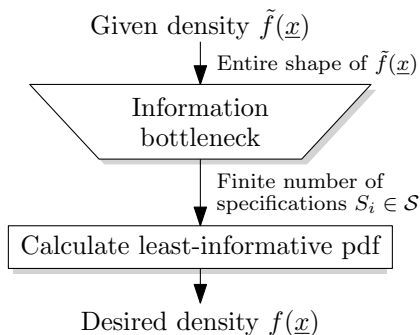


Fig. 1: A finite number of specifications, say, moments or samples, are extracted from the shape of the given density $\tilde{f}(\underline{x})$. Compared to the entire shape of $\tilde{f}(\underline{x})$, the set of specifications \mathcal{S} contains less information. The extraction of specification comprises an information bottleneck.

B. Application Examples

Some selected application examples include (1) increasing the number of parameters of density representation, e.g., increasing the number of components of a mixture pdf, (2) replacing a moment representation of $\tilde{f}(\underline{x})$ with a continuous density representation $f(\underline{x})$, and (3) estimating a continuous pdf representing a given set of samples. Common to these examples is the insufficient specification of the desired pdf, which would yield degenerate results. For example, in the density estimation case, a kernel density estimator with Gaussian kernels would degenerate in the sense of zero kernel widths.

II. PROBLEM FORMULATION

First, we are given a set of specifications $S_i(f) = 0$, $S_i \in \mathcal{S}$ about the desired pdf $f(\underline{x})$. We assume an underlying pdf $\tilde{f}(\underline{x})$ that is unknown, see Fig. 1, that we would like to approximate with the desired pdf $f(\underline{x})$.³ Specifications about $\tilde{f}(\underline{x})$ can be of the form (1) moments or central moments, e.g., mean, covariance, (2) given density values $\tilde{f}(\underline{x}_i)$ at selected \underline{x}_i and/or derivatives, (3) probabilities of \underline{x} being in certain regions, and (4) given samples of \tilde{f} . We also allow ranges or tolerance bands⁴ for certain specifications.

Second, we are given a certain structure of the pdf $f(\underline{x})$. This can be achieved by confining $f(\underline{x})$ to a certain class of densities, say, spherically invariant densities [26]. More concrete, the pdf could be fixed to being of a specific type,

say, a Gaussian density. Here, our focus is on pdfs from the class of mixture densities (MDs) with a given number L of components

$$f(\underline{x}) = \sum_{i=1}^L w_i f_i(\underline{x}) \quad (1)$$

with positive weights $w_i > 0$, $i \in 1:L$, summing up to one

$$\sum_{i=1}^L w_i = 1 \quad , \quad (2)$$

and mixture component pdfs $f_i(\cdot)$, $i \in 1:L$. As a concrete density type, we will consider Gaussian components.

Our desired outcome is a MD $f(\underline{x})$ that does not add any artificial constraints beyond the specifications given. More specific, we require appropriate mixture weights $w_i > 0$, $i \in 1:L$, and parameters of the component pdfs $f_i(\cdot)$, $i \in 1:L$.

For solving this problem, we require a measure of the information content of pdf $f(\underline{x})$ or equivalently of its smoothness/roughness. This measure is then used to find the mixture $f(\underline{x})$ with the least amount of information (or equivalently the smoothest pdf) given the specifications \mathcal{S} .

In the next section, we review common smoothness and information measures and their suitability for the task at hand.

III. STATE-OF-THE-ART

In this section, we briefly review the state-of-the-art in calculating least-informative pdfs. Given a set of specifications on the pdf, a least-informative pdf does not add further information. A pdf with little information content intuitively exhibits a high level of smoothness or equivalently a low level of roughness. We will focus on methods based on curvature, entropy, and FI as a basis for the derivation of the proposed new method in Sec. IV.

A. Mean Curvature

At first look, mean curvature would be a good candidate for quantifying smoothness. The problem is that integrating over the curvature is numerically difficult. Hence, curvature is often approximated by the second derivative of the considered pdf. Expressions of this type are used as roughness penalties for smoothing splines [27, p. 177].

For simplicity, we only consider the scalar case, which is already quite complex. The local curvature of f is defined as⁵

$$k(f(x)) = \frac{f''(x)}{\left(1 + (f'(x))^2\right)^{3/2}} \quad . \quad (3)$$

The mean squared curvature is given as the integral over the local curvature as

$$k(f) = \int_{\mathbb{R}} k^2(f(x)) dx \quad , \quad (4)$$

which cannot be solved in closed form for densities of interest here. For this reason, the local curvature is often approximated with the second derivative of f as $k(f(x)) = f''(x)$. In higher dimensions, this can be generalized to $k(f(\underline{x})) = \nabla^2 f(\underline{x})$ with

Hessian operator $\nabla^2 = \nabla \cdot \nabla^\top$. The resulting mean squared curvature

$$k(f) = \int_{\mathbb{R}^D} |\nabla^2 f(\underline{x})|^2 d\underline{x} \quad (5)$$

with Frobenius norm $|\cdot|$ can often be calculated in closed form, e.g., for Gaussian or Gaussian Mixture (GM) densities $f(\underline{x})$. However, although this simplified mean curvature might be useful for, e.g., spline smoothing, it is not a good indicator of smoothness for pdfs. The reason is that in a typical reference situation, i.e., minimization of the roughness of a density with zero mean and unit variance, minimization of the curvature does not yield a Gaussian density as would be expected. On the other hand, maximization of the entropy or minimization of the FI yields the expected Gaussian density in this case as we will see in Subsec. III-B and Subsec. III-C.

Summary: The mean curvature of a general density according to (3) and (4) are difficult to compute even in the univariate case. This is exacerbated for mixture densities. For this reason, the curvature is often approximated by the second derivative of $f(x)$ in 1D. In the multivariate case, the definition of curvature is even more complicated, which is another reason for approximation by the Hessian of $f(\underline{x})$. Unfortunately, the simplified mean curvature, although easily computable for mixtures, does not yield the expected results in simple reference cases.

B. Entropy

Most methods for finding the least-informative pdf under given specifications use the principle of maximum entropy. It was conceived by Jaynes in 1957 [14, 15]. It received its share of criticism [7], but was eventually adopted by the research community. One example is the use of maximum entropy for the trigonometric moment problem and orthogonal polynomials [16].

For a continuous density $f(\underline{x})$, the so-called differential entropy is defined as

$$E\{-\log(f(\underline{x}))\} = - \int_{\underline{x} \in \mathbb{R}^D} f(\underline{x}) \log(f(\underline{x})) d\underline{x} \quad . \quad (6)$$

The differential entropy has many nice properties. However, it is an ad-hoc generalization of the famous Shannon entropy for discrete random variables to the continuous setting. It can assume negative values and is not invariant under a change of variables.

A table with entropy expressions for standard pdfs is provided in [17]. Expressions for the entropy in multivariate settings are given in [6], but this does not include mixture distributions. Because of the logarithm in the expression of the differential entropy used in maximum entropy methods, computation is complicated beyond simple cases such as Gaussian densities. This is especially a problem for mixture densities as these lead to logarithms of sums. As a result, entropy is often calculated via numerical integration, e.g., Monte Carlo, which is especially complex in multivariate settings.

This led to an extensive development of approximations and bounds, especially for mixture densities. Differential entropy for GMs is approximated in [12] by a Taylor-series expansion of the logarithm of the GM around each component mean. For large component variances, splitting of components is required to maintain a desired accuracy. In [11], a deterministic sample representation is approximated by a piecewise constant density to facilitate the computation of the relative entropy. A piecewise constant approximation of a GM is proposed in [30]. Sharp bounds on the so-called entropy concavity deficit are derived in [19], i.e., the difference between the mixture entropy and the sum of component entropies. In [21], lower and upper bounds on the differential entropy for the special case of GMs are provided when all components have identical variances and only differ in their means and weights. For a symmetric GM with two components with equal weights and equal variances, lower and upper bounds on the differential entropy are given in [20].

Summary: As it requires integration over a logarithm of $f(\underline{x})$, differential entropy can be calculated in closed form only in rare cases. In particular, MDs do not admit a closed-form solution. In that case, one has to rely on the approximations or bounds mentioned above.

C. Fisher Information

Using FI for finding the smoothest continuous density has first been proposed in [9]. It was used as a roughness penalty in maximum likelihood density estimation [10] based on orthogonal Hermite polynomials. Similarly, FI is employed in [29] for wavelet-based density estimation. As a roughness measure, it prevents the density estimate to come too close to the Dirac functions representing the observations. In [13], FI is used for the robust estimation of a location parameter. This method has been extended to minimizing FI over mixtures in [1].

For deriving FI, we start with the so-called score [28, p. 18] given by

$$s(\underline{x}; \underline{\theta}) = \frac{\partial}{\partial \underline{\theta}} \log(f(\underline{x}; \underline{\theta})) \quad (7)$$

for some density f depending on a vector parameter $\underline{\theta}$ or equivalently

$$s(\underline{x}; \underline{\theta}) = \frac{\partial f(\underline{x}; \underline{\theta})}{\partial \underline{\theta}} \frac{1}{f(\underline{x}; \underline{\theta})} = \frac{\nabla f(\underline{x}; \underline{\theta})}{f(\underline{x}; \underline{\theta})} . \quad (8)$$

Following [9, p. 29], we define roughness as the difference between the original density $f(\underline{x})$ and a copy $f(\underline{x} + \underline{\theta})$ shifted by a small displacement $\underline{\theta}$. Thus, we have $f(\underline{x}; \underline{\theta}) = f(\underline{x} + \underline{\theta})$ and the score is

$$s(\underline{x}; \underline{\theta}) = \frac{\partial f(\underline{x} + \underline{\theta})}{\partial \underline{\theta}} \frac{1}{f(\underline{x} + \underline{\theta})} = \frac{\partial f(\underline{x} + \underline{\theta})}{\partial \underline{x}} \frac{1}{f(\underline{x} + \underline{\theta})} . \quad (9)$$

With $\underline{\theta} \rightarrow \underline{0}$, the score (8) can be written as [22, p. 2]

$$s(\underline{x}) = \frac{\partial f(\underline{x})}{\partial \underline{x}} \frac{1}{f(\underline{x})} = \frac{\nabla f(\underline{x})}{f(\underline{x})} \quad (10)$$

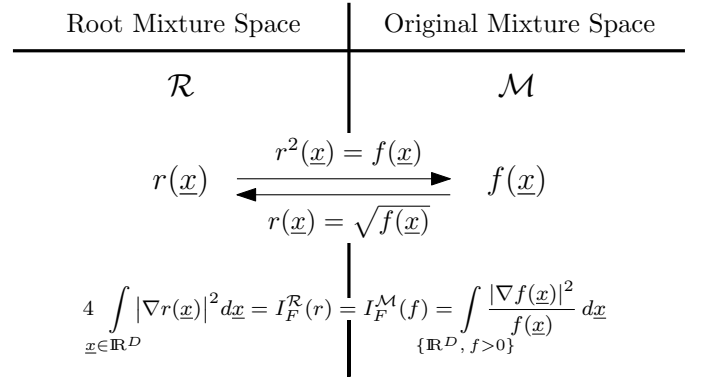


Fig. 2: The two spaces used for calculating least-informative pdfs fulfilling given specifications.

or the score (7) as

$$s(\underline{x}) = \frac{\partial}{\partial \underline{x}} \log(f(\underline{x})) . \quad (11)$$

Integrating over the entire domain in the original space of mixtures \mathcal{M} gives the FI [28, p. 15, (1)]

$$I_F^{\mathcal{M}}(f) = \int_{\{\mathbb{R}^D, f > 0\}} \left| \frac{\nabla f(\underline{x})}{f(\underline{x})} \right|^2 f(\underline{x}) d\underline{x} \quad (12)$$

or alternatively

$$I_F^{\mathcal{M}}(f) = \int_{\{\mathbb{R}^D, f > 0\}} \frac{|\nabla f(\underline{x})|^2}{f(\underline{x})} d\underline{x} . \quad (13)$$

We exclude regions where the density approaches zero.

The FI is often simpler to compute than the differential entropy. However, the division by $f(\underline{x})$ in (13) still makes computation difficult especially for mixture densities. In [1], although minimizing FI over mixtures has been considered, computational issues were not addressed.

Summary: FI in the form (12) or (13) cannot be calculated in closed form for MDs because of the division by $f(\underline{x})$.

IV. CLOSED-FORM FISHER INFORMATION FOR MIXTURES

A. Key Idea

The key idea for obtaining a closed-form expression for the FI of a mixture density is to work in the space of square-root densities, see Fig. 2. In particular, we define RMs that yield the original mixture upon squaring. The desired FI of the original mixture can be calculated in closed form in the RM space. We propose to simultaneously maintain both density representations and perform tandem processing in both spaces. The FI is calculated in the RM space \mathcal{R} . The specifications on the pdf $f(\underline{x})$ are calculated in the original mixture space \mathcal{M} .

B. Root Densities

Square roots of densities $r(\underline{x}) = \sqrt{f(\underline{x})}$ have already been considered in [10]. The reason for doing so was to ensure that the squared density $f(\underline{x}) = r^2(\underline{x})$ is non-negative. For the same reason, densities based on squared Fourier series are

introduced in [2], generalized for multivariate densities in [3], and improved in [25].

Here, we use square-root pdfs, or in short root densities (RDs), for an entirely different reason. Our considered density representations, e.g., Gaussian mixtures, are non-negative anyway. Instead, we want to obtain a closed-form expression for the FI in (13).

The unit-mass constraint for the pdf $f(\underline{x})$

$$\int_{\mathbb{R}^D} f(\underline{x}) d\underline{x} = 1 \quad (14)$$

is equivalent to the constraint

$$\int_{\mathbb{R}^D} r^2(\underline{x}) d\underline{x} = 1 \quad (15)$$

for the RD $r(\underline{x})$. This means that $r(\underline{x})$ is restricted to the infinite-dimensional unit sphere S^∞ . This manifold is also called the Hilbert sphere [5, p. 2]. When we restrict the RDs $r(\underline{x})$ to be non-negative, i.e., $r(\underline{x}) \geq 0 \forall \underline{x}$, $r(\underline{x})$ is restricted to the ‘‘positive orthant’’ of the Hilbert sphere.

C. Fisher Information for Root Densities

For RDs, the expression for the FI can be simplified. In the univariate case, $D = 1$, with

$$r'(\underline{x}) = \frac{d}{d\underline{x}} r(\underline{x}) = \frac{d}{d\underline{x}} \sqrt{f(\underline{x})} = -\frac{f'(\underline{x})}{2\sqrt{f(\underline{x})}}, \quad (16)$$

we have

$$(r'(\underline{x}))^2 = \frac{(f'(\underline{x}))^2}{4f(\underline{x})}. \quad (17)$$

Finally, we can rewrite the FI in the original density space in (13) in the RD space \mathcal{M} as

$$I_F^{\mathcal{R}}(r) = 4 \int_{\underline{x} \in \mathbb{R}^D} (r'(\underline{x}))^2 d\underline{x} \text{ with } r = \sqrt{f}. \quad (18)$$

In the multivariate case, we have

$$I_F^{\mathcal{R}}(r) = 4 \int_{\underline{x} \in \mathbb{R}^D} |\nabla r(\underline{x})|^2 d\underline{x} \text{ with } r = \sqrt{f}. \quad (19)$$

These expressions for $I_F^{\mathcal{R}}(r)$ in the RD space \mathcal{R} can be calculated in closed form, e.g., for RMs. In fact, they correspond to the simplified expressions for mean curvature in (5) that are, however, given in the original density space.

D. Root Mixtures

We define RDs in the case of MDs of the form (1) as

$$r(\underline{x}) = \sum_{i=1}^R v_i r_i(\underline{x}). \quad (20)$$

These square-root MDs will be abbreviated as RMs.

1) *Conversion of Root Mixture to Mixture:* The conversion of a RM in (20) to a mixture in (1) is unique, always exists, and done in closed form according to

$$\begin{aligned} f(\underline{x}) &= r^2(\underline{x}) \\ &= \left(\sum_{i=1}^R v_i r_i(\underline{x}) \right) \left(\sum_{i=1}^R v_i r_i(\underline{x}) \right) \\ &= \sum_{i=1}^R \sum_{j=1}^R v_i v_j r_i(\underline{x}) r_j(\underline{x}). \end{aligned} \quad (21)$$

This expression contains redundant terms as $r_i(\underline{x}) r_j(\underline{x}) = r_j(\underline{x}) r_i(\underline{x})$ for $i \neq j$ and can be written as

$$f(\underline{x}) = \sum_{i=1}^R v_i^2 r_i^2(\underline{x}) + 2 \sum_{i=1}^R \sum_{j=i+1}^R v_i v_j r_i(\underline{x}) r_j(\underline{x}) \quad (22)$$

or⁶

$$f(\underline{x}) = \sum_{i=1}^R \sum_{j=i}^R v_i v_j c_{i,j} r_i(\underline{x}) r_j(\underline{x}) \quad (23)$$

with

$$c_{i,j} = \begin{cases} 1, & i = j \\ 2, & i \neq j \end{cases}. \quad (24)$$

The MD $f(\underline{x})$ in (23) contains a total of

$$L = R \cdot (R + 1) / 2 \quad (25)$$

components⁷ and is non-negative by definition (although its weights may be negative). However, without additional constraints on $r(\underline{x})$, the unit mass constraint may be violated.

For deriving specific constraints, we need to consider specific mixture densities. Here, we will consider mixtures of the form (1) with the special case of Gaussian components

$$f_i(\underline{x}) = \frac{1}{\sqrt{|2\pi \Sigma_i|}} \exp \left\{ -\frac{1}{2} (\underline{x} - \underline{x}_i)^\top \cdot \Sigma_i^{-1} \cdot (\underline{x} - \underline{x}_i) \right\} \quad (26)$$

with $\underline{x} \in \mathbb{R}^D$ abbreviated as⁸

$$f_i(\underline{x}) = N(\underline{x}; \underline{x}_i, \Sigma_i). \quad (27)$$

The corresponding Gaussian RM for $\underline{x} \in \mathbb{R}^D$ is given by (20) with Gaussian components

$$r_i(\underline{x}) = \frac{1}{\sqrt{|2\pi \mathbf{P}_i|}} \exp \left\{ -\frac{1}{2} (\underline{x} - \underline{\rho}_i)^\top \cdot \mathbf{P}_i^{-1} \cdot (\underline{x} - \underline{\rho}_i) \right\} \quad (28)$$

with weights v_i , mean vectors $\underline{\rho}_i$, and covariance matrices \mathbf{P}_i .

According to (23), the conversion of a Gaussian RM with R components and parameters v_i , $\underline{\rho}_i$, \mathbf{P}_i to a GM results in $L = R \cdot (R + 1) / 2$ components with w_i , \underline{x}_i , Σ_i given in Appendix VII. Please note that the conversion is especially simple for Gaussian components as the product of Gaussians is again Gaussian.

Basic Constraints on Parameters of $r(\underline{x})$: The GM $f(\underline{x})$ resulting from the conversion of the Gaussian RM $r(\underline{x})$ does not necessarily have positive weights that sum to one. For this reason, basic constraints on the parameters of $r(\underline{x})$ in RM space \mathcal{R} are derived from the constraints

$$w_i > 0, \quad i = 1:L, \quad \sum_{i=1}^L w_i = 1 \quad (29)$$

on the GM $f(\underline{x})$ in the original mixture space \mathcal{M} .

Number of Parameters: The number of parameters of the Gaussian RM $r(\underline{x})$ and the GM $f(\underline{x})$ are equal as the squaring operation in (23) does not add parameters. In the multivariate case, the parameters of $r(\underline{x})$ with R components are given by⁹

$$\underbrace{v_1, v_2, \dots, v_R}_R, \underbrace{\rho_1, \rho_2, \dots, \rho_R}_{R \cdot D}, \underbrace{\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_R}_{R \cdot \frac{D \cdot (D+1)}{2}}. \quad (30)$$

With the above mentioned basic constraint we lose one degree of freedom in the weights, so that we end up with

$$(R-1) + R \cdot D + R \cdot \frac{D \cdot (D+1)}{2} = R \cdot \frac{D^2 + 3D + 2}{2} - 1 \quad (31)$$

parameters.

Limiting the Number of Components: The conversion from an RM to a mixture becomes impractical for large R as L increases quadratically with R . However, there are several simple and effective ways to alleviate this problem. First, we note that there is redundancy in the resulting mixture in (23) as the number of parameters of $f(\underline{x})$ stays the same as in $r(\underline{x})$, but the number of components is larger. This redundancy could be exploited in a subsequent component reduction of $f(\underline{x})$. However, this is computationally costly. A better alternative is to exploit the decreasing overlap between distant mixture components. For mixtures with finitely supported components, say, triangular mixtures, the overlap becomes zero after a certain distance. For components with infinite support, e.g., in the case of GMs, the overlap quickly decreases in practical settings. Hence, not every combination of components has to be multiplied. This typically results in $L \ll R^2$.

The extreme case would be to assume no overlaps between mixture components at all. In this case, which is realistic for GMs, the numbers of components would stay exactly the same, i.e., $L = R$. A more practical approach is to only consider close neighbors or just the nearest neighbors. However, closeness does not only depend on the locations of the individual components but also on their extent. This calls for an adaptive approach that only includes those pairs of mixture components in the expansion (23) that yield new components with significant weights.

2) *Conversion of Mixture to Root Mixture:* An RM $r(\underline{x})$ corresponding to a given MD $f(\underline{x})$ does not necessarily exist. This is due to the smaller number of parameters in $r(\underline{x})$ compared to $f(\underline{x})$ so that not every possible $f(\underline{x})$ can be represented by $r(\underline{x})$. In addition, even when it exists, an appropriate RM is non-unique due to, e.g., squared weights. Again, we can exploit

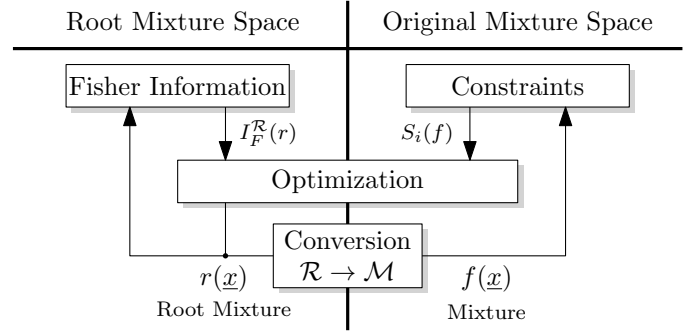


Fig. 3: Overview of optimization via tandem processing in RM space \mathcal{R} and original mixture space \mathcal{M} .

small overlaps between distant components when performing the conversion. This reduces the differences in numbers of parameters between the RM and the original mixture.

E. Tandem Processing for Optimization

The RM $r(\underline{x})$ in root mixture space \mathcal{R} is convenient as it allows a closed-form calculation of the FI. In addition, the nonnegativity of the corresponding mixture $f(\underline{x}) = r^2(\underline{x})$ in mixture space \mathcal{M} is automatically guaranteed. However, the unit integral constraint for $f(\underline{x})$ has to be explicitly ensured by an appropriate constraint in mixture space \mathcal{M} . Also, the specifications on $f(\underline{x})$ are formulated in mixture space \mathcal{M} .

During the optimization procedure, i.e., minimization of FI, see Fig. 3, we maintain both density representations in RM space \mathcal{R} and in mixture space \mathcal{M} . FI calculation is done in RM space \mathcal{R} in tandem with the constraint evaluation in mixture space \mathcal{M} .

The optimization problem can be written as

$$\begin{aligned} \min_{f \in \mathcal{M}} \quad & I_F^{\mathcal{M}}(f) \\ \text{s.t.} \quad & S_i(f) = 0, \quad S_i \in \mathcal{S}, \end{aligned} \quad (32)$$

with the FI $I_F^{\mathcal{M}}(f)$ from (12). This is equivalent to

$$\begin{aligned} \min_{r \in \mathcal{R}} \quad & I_F^{\mathcal{R}}(r) \\ \text{s.t.} \quad & f = r^2 \\ & S_i(f) = 0, \quad S_i \in \mathcal{S}, \end{aligned} \quad (33)$$

with the FI $I_F^{\mathcal{R}}(r)$ from (18) and $I_F^{\mathcal{M}}(f) = I_F^{\mathcal{R}}(r)$.

V. EXAMPLES

We will now give two example applications of the proposed new optimization method. In the examples, we assume a GM $f(\underline{x})$ of the form (1) with L components of the form (26). We start with a simple example of the smoothest GM with zero mean and unit variance, see Example V.1. In Example V.2, we prescribe density values at certain locations.

Example V.1 (Zero mean and unit variance). *Besides the constraints on the weights (larger than zero and unit sum), the specifications are only zero mean and unit variance. We know that of all pdfs with given covariance matrix, the Gaussian density minimizes the FI¹⁰ [24, Lemma 1, p. 184] and [23]. Hence, we would expect the GM $f(\underline{x})$ to approach a Gaussian*

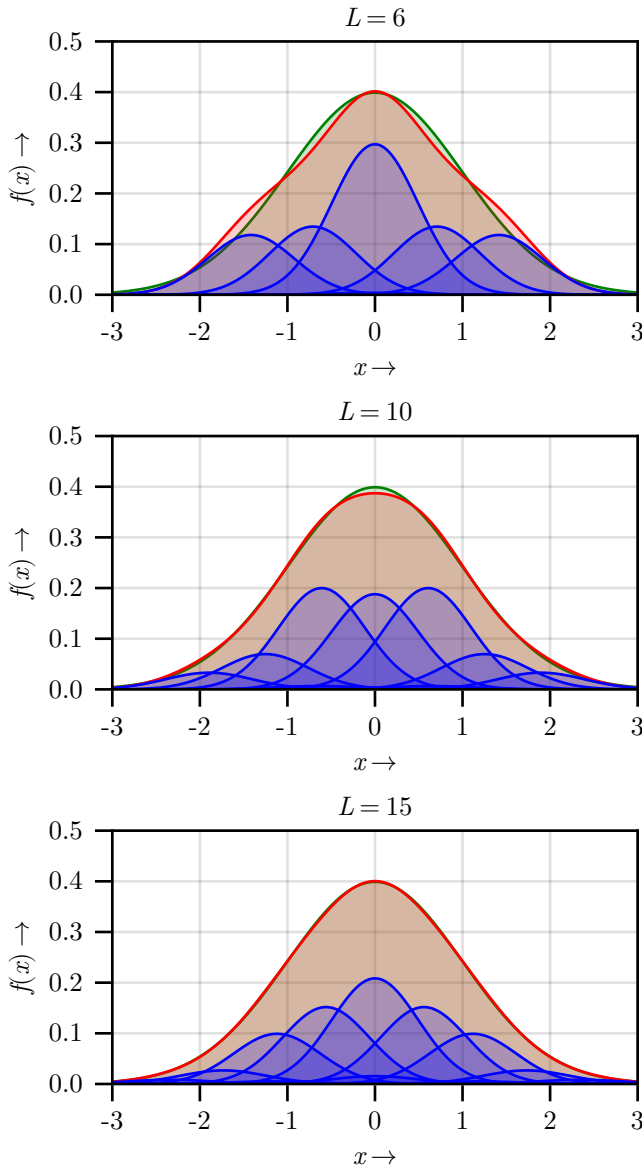


Fig. 4: Red: Smoothest GM $f(\underline{x})$ under specifications zero mean and unit variance in Example V.1 the first example for different numbers of components L . Blue: Corresponding mixture components $f_i(\underline{x})$, $i \in 1:L$. Green: Gaussian density with zero mean and unit variance as reference.

shape when the number of components grows. The results are shown in the original mixture space \mathcal{M} in Fig. 4 for $R \in \{3, 4, 5\}$ RM components. With (23), we obtain $L = R \cdot (R + 1)/2 \in \{6, 10, 15\}$ mixture components. For $L = 6$, the number of parameters in the GM $f(\underline{x})$ is not sufficient to approach the expected Gaussian density. For $L = 10$, $f(\underline{x})$ is closer and for $L = 15$, $f(\underline{x})$ is visually indistinguishable from the Gaussian. For $R = 5$ and $L = 15$, it is apparent, that not all components significantly contribute to the density shape as some components are very small. **◀Example▶**

Example V.2 (Prescribed values and derivatives). In this example, we use the weight constraints and zero mean as basic constraints. In addition, we specify the following

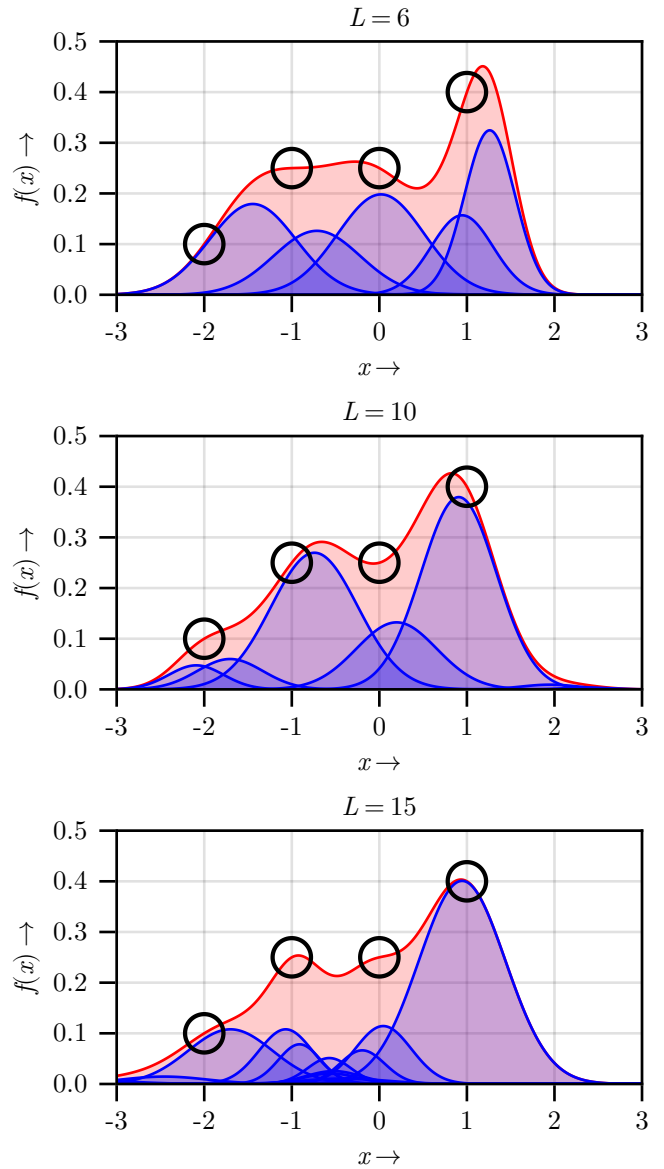


Fig. 5: Red: Smoothest GM $f(\underline{x})$ under zero mean and value specifications in Example V.2 for different numbers of components L . Blue: Corresponding mixture components $f_i(\underline{x})$, $i \in 1:L$.

$(x, f(x))$ value pairs: $(-2.0, 0.1)$, $(-1.0, 0.25)$, $(0.0, 0.25)$, and $(1.0, 0.4)$. The results are shown in Fig. 5 for $R \in \{3, 4, 5\}$ RM components and $L \in \{6, 10, 15\}$ mixture components. For all L , the constraints are matched. For increasing L , the density $f(x)$ becomes smoother. **◀Example▶**

VI. CONCLUSIONS

We derive a closed-form expression for the FI of a GM density $f(\underline{x})$ given in a mixture space \mathcal{M} . This is achieved by introducing an RM space \mathcal{R} and calculating the FI there. The resulting expression consists of a combination of higher-order moments of Gaussian densities, containing nonlinear terms of the component means and covariances. We consider very general specifications of the shape, probability mass distribution, and (central) moments of $f(\underline{x})$. These specifications can also typically be described by nonlinear equations containing

component means and covariances. Standard optimization procedures for minimization under equality constraints can be employed to find the least-informative GM density $f(\underline{x})$ under the given specifications.

VII. APPENDIX

Product of Two Multivariate Gaussian Densities: Given two multivariate Gaussian densities

$$f_i(\underline{x}) = N(\underline{x}; \underline{x}_i, \mathbf{\Sigma}_i) \quad (34)$$

or

$$f_i(\underline{x}) = \frac{1}{\sqrt{|2\pi\mathbf{\Sigma}_i|}} \exp \left\{ -\frac{1}{2}(\underline{x} - \underline{x}_i)^\top \mathbf{\Sigma}_i^{-1}(\underline{x} - \underline{x}_i) \right\} \quad (35)$$

for $i = 1, 2$, the product $f_3(\underline{x}) = f_1(\underline{x}) \cdot f_2(\underline{x})$ is given by

$$f_3(\underline{x}) = c_3 \cdot N(\underline{x}; \underline{x}_3, \mathbf{\Sigma}_3) \quad (36)$$

with factor

$$c_3 = N(0; \underline{x}_1 - \underline{x}_2, \mathbf{\Sigma}_1 + \mathbf{\Sigma}_2) , \quad (37)$$

new mean vector

$$\underline{x}_3 = (\mathbf{\Sigma}_1^{-1} + \mathbf{\Sigma}_2^{-1})^{-1} (\mathbf{\Sigma}_1^{-1} \underline{x}_1 + \mathbf{\Sigma}_2^{-1} \underline{x}_2) , \quad (38)$$

and new covariance matrix

$$\mathbf{\Sigma}_3 = (\mathbf{\Sigma}_1^{-1} + \mathbf{\Sigma}_2^{-1})^{-1} . \quad (39)$$

Please note that $f_3(\underline{x})$ is an unnormalized density, where the factor c_3 reflects the distance between the original means. Obviously, the new density $f_3(\underline{x})$ as the result of the multiplication of $f_1(\underline{x})$ and $f_2(\underline{x})$ contains less probability mass, the farther $f_1(x)$ and $f_2(x)$ are apart. The density can easily be normalized by omitting the factor c_3 . However, when multiplying two Gaussian mixtures, the factors have to be considered as they take care of the relative weights.

ENDNOTES

- 1: We use the notation $i:j$ to denote integer sequence $\{i, i+1, \dots, j-1, j\}$ from [8].
- 2: By selecting the density $f(\underline{x})$ from a certain class of densities, we inadvertently add information.
- 3: The underlying pdf $\tilde{f}(\underline{x})$ is mainly a vehicle for explanation. However, in some cases, the goal might indeed be to reconstruct $\tilde{f}(\underline{x})$ from given moments or samples. It is important to note that $f(\underline{x})$ is completely unknown and not used anywhere in the procedure of calculating the desired pdf $f(\underline{x})$.
- 4: Ranges or tolerance bands could be given for moments, values, or even samples. An example would be a mean constraint of the form $E\{\mathbf{x}\} \in [-0.2, 0.2]$. Another example is a value constraint of the form $f(x_i) \in [y, \bar{y}]$.
- 5: Defining curvature in the multivariate case is much more complicated, see for example [18].
- 6: Please note that the summations range in $i \in [1, R]$ and $j \in [i, R]$ as we exploit symmetry.

- 7: When we prespecify the minimum number of components L we desire of a GM in the original mixture space \mathcal{M} , the number of RM components in the RM space \mathcal{R} for a full expansions according to (23) are given by

$$R = \left\lceil \frac{\sqrt{8 \cdot L + 1} - 1}{2} \right\rceil , \quad (40)$$

where $\lceil \cdot \rceil$ denotes the next largest integer.

- 8: For a multivariate normal random vector $\underline{x} \in \mathbb{R}^D$ we denote a Gaussian density by

$$f(\underline{x}) = N(\underline{x}; \underline{m}, \mathbf{C}) \quad (41)$$

with realization $\underline{x} \in \mathbb{R}^D$, mean vector $\underline{m} \in \mathbb{R}^D$, and covariance matrix \mathbf{C} .

- 9: The covariance matrices \mathbf{P}_i are positive-definite and symmetric. Thus, each \mathbf{P}_i is specified by $D \cdot (D+1)/2$ parameters.
- 10: For a given covariance matrix, the Gaussian density also maximizes the relative entropy [4, p. 411, Example 12.2.1, (12.14)].

REFERENCES

- [1] Peter J. Bickel and John R. Collins. "Minimizing Fisher Information over Mixtures of Distributions". *Sankhyā: The Indian Journal of Statistics, Series A (1961-2002)* 45.1 (1983), pp. 1–19.
- [2] Dietrich Brunn, Felix Sawo, and Uwe D. Hanebeck. "Efficient Nonlinear Bayesian Estimation based on Fourier Densities". *Proceedings of the 2006 IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems (MFI 2006)*. Heidelberg, Germany, Sept. 2006, pp. 312–322.
- [3] Dietrich Brunn, Felix Sawo, and Uwe D. Hanebeck. "Nonlinear Multidimensional Bayesian Estimation with Fourier Densities". *Proceedings of the 2006 IEEE Conference on Decision and Control (CDC 2006)*. San Diego, California, USA, Dec. 2006, pp. 1303–1308.
- [4] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley Series in Telecommunications and Signal Processing. USA: Wiley-Interscience, June 2006.
- [5] Xiongtao Dai. *Statistical Inference on the Hilbert Sphere with Application to Random Densities*. Jan. 2, 2021, preprint.
- [6] G.A Darbellay and I Vajda. "Entropy Expressions for Multivariate Continuous Distributions". *IEEE Transactions on Information Theory* 46.2 (Mar. 2000), pp. 709–712.
- [7] Kenneth Friedman and Abner Shimony. "Jaynes's Maximum Entropy Prescription and Probability Theory". *Journal of Statistical Physics* 3.4 (Dec. 1, 1971), pp. 381–384.

- [8] Gene H. Golub and Charles F. Van Loan. *Matrix Computations (3rd Ed.)* Baltimore, MD, USA: Johns Hopkins University Press, 1996.
- [9] Irving J. Good. “Non-Parametric Roughness Penalty for Probability Densities”. *Nature Physical Science* 229.1 (1 Jan. 1971), pp. 29–30.
- [10] Irving J. Good and Ray A. Gaskins. “Nonparametric Roughness Penalties for Probability Densities”. *Biometrika* 58.2 (Aug. 1, 1971), pp. 255–277.
- [11] Uwe D. Hanebeck. “Truncated Moment Problem for Dirac Mixture Densities with Entropy Regularization”. *arXiv preprint: Systems and Control (cs.SY)* (Aug. 2014).
- [12] Marco F. Huber et al. “On Entropy Approximation for Gaussian Mixture Random Vectors”. *Proceedings of the 2008 IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems (MFI 2008)*. Seoul, Republic of Korea, Aug. 2008, pp. 181–188.
- [13] Peter J. Huber. “Robust Estimation of a Location Parameter”. *The Annals of Mathematical Statistics* 35.1 (1964), pp. 73–101.
- [14] Edwin T. Jaynes. “Information Theory and Statistical Mechanics”. *Physical Review* 106.4 (May 15, 1957), pp. 620–630.
- [15] Edwin T. Jaynes. “Information Theory and Statistical Mechanics. II”. *Physical Review* 108.2 (Oct. 15, 1957), pp. 171–190.
- [16] Henry J. Landau. “Maximum Entropy and the Moment Problem”. *Bulletin (New Series) of the American Mathematical Society* 16.1 (Jan. 1987), pp. 47–77.
- [17] Aida V. Lazo and Pushpa N. Rathie. “On the Entropy of Continuous Probability Distributions”. *IEEE Transactions on Information Theory* 24.1 (Jan. 1978), pp. 120–122.
- [18] John M. Lee. *Riemannian Manifolds: An Introduction to Curvature*. Graduate Texts in Mathematics 176. New York: Springer, 1997. 224 pp.
- [19] James Melbourne et al. “The Differential Entropy of Mixtures: New Bounds and Applications”. *IEEE Transactions on Information Theory* 68.4 (Apr. 2022), pp. 2123–2146.
- [20] Joseph V. Michalowicz, Jonathan M. Nichols, and Frank Bucholtz. “Calculation of Differential Entropy for a Mixed Gaussian Distribution”. *Entropy* 10.3 (3 Sept. 2008), pp. 200–206.
- [21] Kamyar Moshksar and Amir K. Khandani. “Arbitrarily Tight Bounds on Differential Entropy of Gaussian Mixtures”. *IEEE Transactions on Information Theory* 62.6 (June 2016), pp. 3340–3354.
- [22] Ágnes Nagy. “Fisher Information and Density Functional Theory”. *International Journal of Quantum Chemistry* 122.8 (2022), e26679.
- [23] Abbas Pak. “An Alternative Proof For the Minimum Fisher Information of Gaussian Distribution”. *Journal of Applied Mathematics, Statistics and Informatics* 14.2 (Dec. 1, 2018), pp. 5–10.
- [24] Sangwoo Park, Erchin Serpedin, and Khalid Qaraqe. “Gaussian Assumption: The Least Favorable but the Most Useful”. *IEEE Signal Processing Magazine* 30.3 (May 2013), pp. 183–186.
- [25] Florian Pfaff, Gerhard Kurz, and Uwe D. Hanebeck. “Multivariate Angular Filtering Using Fourier Series”. *Journal of Advances in Information Fusion* 11.2 (Dec. 2016), pp. 206–226.
- [26] M. Rangaswamy, D. Weiner, and A. Ozturk. “Non-Gaussian Random Vector Identification Using Spherically Invariant Random Processes”. *IEEE Transactions on Aerospace and Electronic Systems* 29.1 (1993), pp. 111–124.
- [27] Christian H. Reinsch. “Smoothing by Spline Functions”. *Numerische Mathematik* 10.3 (Oct. 1, 1967), pp. 177–183.
- [28] Giuseppe Toscani. “Score Functions, Generalized Relative Fisher Information and Applications”. *Ricerche di Matematica* 66.1 (June 1, 2017), pp. 15–26.
- [29] Marina Vannucci and Brani Vidakovic. “Preventing the Dirac Disaster: Wavelet Based Density Estimation”. *Journal of the Italian Statistical Society* 6.2 (Aug. 1, 1997), p. 145.
- [30] Chenchen Zhang and Yuan Luo. “Approximating the Differential Entropy of Gaussian Mixtures”. *GLOBECOM 2017 - 2017 IEEE Global Communications Conference*. GLOBECOM 2017 - 2017 IEEE Global Communications Conference. Dec. 2017, pp. 1–6.