# Multivariate Myriad Filters based on Parameter Estimation of the Student-$t$ Distribution

Friederike Laus[*]     Gabriele Steidl[*†]

June 25, 2019

### Abstract

The contribution of this study is twofold: First, we propose an efficient algorithm for the computation of the (weighted) maximum likelihood estimators for the parameters of the multivariate Student-$t$ distribution, which we call generalized multivariate myriad filter. Second, we use the generalized multivariate myriad filter in a nonlocal framework for the denoising of images corrupted by different kinds of noise. The resulting method is very flexible and can handle heavy-tailed noise such as Cauchy noise, as well as the other extreme, namely Gaussian noise. Furthermore, we detail how the limiting case $\nu \to 0$ of the projected normal distribution in two dimensions can be used for the robust denoising of periodic data, in particular for images with circular data corrupted by wrapped Cauchy noise.

## 1. Introduction

Besides mean and median filter, myriad filter form an important class of nonlinear filters, in particular in robust signal and image processing. While in a multivariate setting, the mean filter can be defined componentwise, the generalization of the median to higher dimensions is not canonically, but often the geometric median is used, see, e.g., [44]. In this paper we make a first attempt to establish a multivariate myriad filter based on the multivariate Student-$t$ distribution.

In one dimension, mean, median as well as myriad filters can be derived as maximum likelihood (ML) estimators of the location parameter from a Gaussian, Laplacian respective Cauchy distribution. Concerning a higher dimensional myriad filter, instead of a multivariate Cauchy

[1]Department of Mathematics, Technische Universität Kaiserslautern, Paul-Ehrlich-Str. 31, D-67663 Kaiserslautern, Germany, {friederike.laus,steidl}@mathematik.uni-kl.de.
[2]Fraunhofer ITWM, Fraunhofer-Platz 1, D-67663 Kaiserslautern, Germany

distribution we propose to start with the family of more general Student-$t$ distributions, which possesses an additional degree of freedom parameter $\nu > 0$ that allows to control the robustness of the resulting filter. While the Cauchy distribution is obtained as the special case $\nu = 1$, the Student-$t$ distribution converges for $\nu \to \infty$ to the normal distribution, so that in the limit also mean filters are covered.

The multivariate Student $t$-distribution is frequently used in statistics [22], whereas the multivariate Cauchy distribution is far less common and in contrast to the one-dimensional case usually not considered separately from the Student-$t$ distribution. The parameter(s) of a multivariate Student $t$-distribution are usually estimated via the Maximum Likelihood (ML) method in combination with the EM algorithm [5, 7, 8, 10, 34]. The EM algorithm for the Student-$t$ distribution has been derived, e.g. in [23], For an overview of estimation methods for the multivariate Student $t$-distribution, in particular the EM algorithm and its variants, we refer to [37] and the references therein.

Recently, the Student-$t$ distribution and the closely related Student-$t$ mixture models (SMM) have found interesting applications in various image processing tasks. One of the first papers which suggested a variational approach for denoising of images corrupted by Cauchy noise was [1]. In [24], the authors proposed a unified framework for images corrupted by white noise that can handle (range constrained) Cauchy noise as well. Other recent approaches that consider also the task of deblurring include [11, 54]. Concerning mixture models, in [48] it has been shown that Student-$t$ mixture models are superior to Gaussian mixture models for modeling image patches and the authors proposed an application in image compression. Further applications include robust image segmentation [2, 38, 45] as well as robust registration [16, 55]. In both cases, the SMM is estimated using the EM algorithm derived in [39].

In this paper, we propose an application of the Student-$t$ distribution to robust denoising of images corrupted by different kinds of noise. The initial motivation for this work were the recent papers [26, 35, 43] for Cauchy noise removal. In [35, 43] the authors proposed a variational method consisting of a data term that resembles the noise statistics and a total variation regularization term. Based on a ML approach the authors of [26] introduced a generalized myriad filter which estimates both the location and the scale parameter of the Cauchy distribution. They used this filter in a nonlocal approach, where for each pixel of the image they chose as samples those pixels having a similar neighborhood and replaced the initial pixel by its filtered version. Such a pixelwise treatment assumes the pixels of an image to be independent, which is in practice a rather unrealistic assumption; in fact, in natural images they are usually locally highly correlated. Taking the local dependence structure into account may improve the results of image restoration methods. For instance, in case of denoising images corrupted by additive Gaussian noise this led to the state-of-the-art algorithm of Lebrun et al. [27], who proposed to restore the image patchwise based on a

maximum a posteriori approach.

In the Gaussian setting, their approach is equivalent to minimum mean square error estimation, and more general, the resulting estimator can be seen as a particular instance of a best linear unbiased estimator (BLUE). For denoising images corrupted by additive Cauchy noise, a similar approach addressed in this paper requires to define a multivariate myriad filter. In this paper, we derive a generalized multivariate myriad filter (GMMF) based on ML estimation for the family of Student-t distributions of which the Cauchy distribution forms a special case. In the limiting case $\nu = 0$ and $d = 2$ dimensions this further provides an algorithm for denoising images corrupted by heavy-tailed noise with image values on the circle $\mathbb{S}^1$, such as phase-valued images appearing in inferometric synthetic aperture radar InSAR.

After finishing this paper, we became aware that our algorithm for estimating the parameters of the Student-$t$ distribution can be considered as Jacobi variant of a sophisticated version of the EM algorithm which was heuristically proposed by Kent et al. in [21] and analyzed by van Dyk in [49]. The approach in the present paper is different and does not require the EM framework with special hidden variables.

The paper is organized as follows: In Section 2, we introduce the Student-$t$ distribution and the projected normal distribution. Their likelihood functions are given in Section 3. Then, in Section 4, we recall existence and uniqueness of the (weighted) ML estimators for the location and the scatter parameter of the distribution, where we provide own proofs. We propose an efficient algorithm for computing the ML estimates in Section 5 and prove its convergence. In Section 6, we illustrate how the developed algorithm can be applied in the context of nonlocal (robust) image denoising both for gray-value images and images with values in $\mathbb{S}^1$. Conclusions and directions of future research are addressed in Section 7.

## 2. Student-$t$ and Projected Normal Distribution

In this section, we introduce the multivariate Student-$t$ distribution and the related projected normal distribution and collect some of their properties.

### 2.1. Multivariate Student-$t$ Distribution

The probability density function (pdf) of the $d$-dimensional Student $t$-distribution $T_\nu(\mu, \Sigma)$ with $\nu > 0$ degrees of freedom is given by

$$f_\nu(x|\mu, \Sigma) = \frac{\Gamma\left(\frac{d+\nu}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)(\pi\nu)^{\frac{d}{2}}|\Sigma|^{\frac{1}{2}}} \frac{1}{\left(1 + \frac{1}{\nu}\delta\right)^{\frac{d+\nu}{2}}}, \tag{1}$$

3

where

$$\delta := (x - \mu)^{\mathrm{T}} \Sigma^{-1} (x - \mu)$$

denotes the the *Mahalanobis distance* between $x$ and a distribution with parameters $\mu, \Sigma$ and $\Gamma(s) := \int_0^\infty t^{s-1} \mathrm{e}^{-t} \, \mathrm{d}t$ the *Gamma function*. The smaller the value of $\nu$ for fixed *location* $\mu$ and positive definite *scatter matrix* $\Sigma$, the heavier are the tails of the $T_\nu(\mu, \Sigma)$ distribution. Figure 1 illustrates this behavior for the one-dimensional standard Student-$t$ distribution. For $\nu \to \infty$, we obtain since

$$\frac{\Gamma(\frac{\nu+d}{2})}{\Gamma(\frac{\nu}{2}) \left(\frac{\nu}{2}\right)^{d/2}} \to 1, \qquad (1 + \frac{1}{\nu}\delta)^{(\nu+d)/2} \to \frac{1}{2}\delta \quad \text{as } \nu \to \infty$$

see [30], that $\lim_{\nu\to\infty} f_\nu(x|\mu,\Sigma) = (2\pi)^{-\frac{d}{2}} |\Sigma|^{-\frac{1}{2}} \mathrm{e}^{-\frac{1}{2}\delta}$ so that the Student-$t$ distribution $T_\nu(\mu, \Sigma)$ converges to the normal distribution $\mathcal{N}(\mu, \Sigma)$.
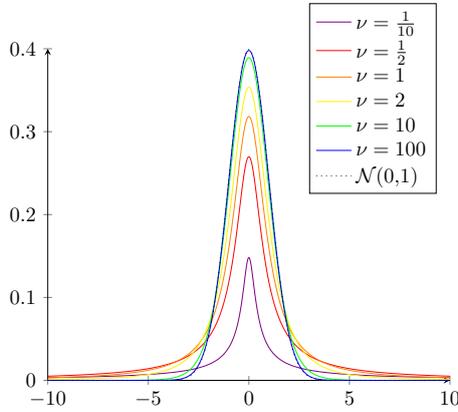


Figure 1: Standard Student-$t$ distribution $T_\nu(0, 1)$ for different values of $\nu$ in comparison with the standard normal distribution $\mathcal{N}(0, 1)$.

The expectation of the Student-$t$ distribution is $\mathbb{E}(X) = \mu$ for $\nu > 1$ and the covariance matrix is given by $\mathrm{Cov}(X) = \frac{\nu}{\nu-2}\Sigma$ for $\nu > 2$, otherwise the quantities are undefined. As the normal distribution, the Student-$t$ distribution belongs to the class of *elliptical distributions*. Some important properties that are needed later on are summarized in the next theorem [22]. In the following, we denote by $\mathrm{Sym}(d)$ the space of symmetric $d \times d$ matrices and by $\mathrm{SPD}(d)$ the cone of symmetric positive definite matrices.

**Theorem 2.1.** *(i) Let $\mu \in \mathbb{R}^d$ and $\Sigma \in \mathrm{SPD}(d)$. Further, let $Z \sim \mathcal{N}(0, \Sigma)$ and $Y \sim \Gamma\left(\frac{\nu}{2}, \frac{\nu}{2}\right)$ be independent, where $\Gamma(\alpha, \beta)$ is the Gamma distribution with parameters $\alpha, \beta > 0$. Then $X = \mu + \frac{Z}{\sqrt{Y}} \sim T_\nu(\mu, \Sigma)$.*

*(ii) Let $X \sim T_\nu(\mu, \Sigma)$, $A \in \mathbb{R}^{d \times d}$ be an invertible matrix and $b \in \mathbb{R}^d$. Then $AX + b \sim T_\nu(A\mu + b, A\Sigma A^{\mathrm{T}})$.*

4

## 2.2. Projected Normal Distribution

For the limiting case $\nu \to 0$, the pdf $f_\nu$ in (1) converges pointwise to zero, i.e.

$$\lim_{\nu \to 0} f_\nu(x|\mu, \Sigma) = \lim_{\nu \to 0} \frac{1}{\Gamma\left(\frac{\nu}{2}\right) \left(1 + \frac{1}{\nu}\delta\right)^{\frac{\nu}{2}}} \frac{\Gamma\left(\frac{d}{2}\right)}{\pi^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}} \delta^{\frac{d}{2}}} = 0$$

However, replacing the first factor by $\frac{1}{2}$ the surface measure $\omega_d = 2\pi^{\frac{d}{2}}/\Gamma\left(\frac{d}{2}\right)$ of the sphere $\mathbb{S}^{d-1} \subset \mathbb{R}^d$ comes into the play and setting $\mu := 0$, we obtain the pdf

$$f_0(x|\Sigma) = \frac{\Gamma\left(\frac{d}{2}\right)}{2\pi^{\frac{d}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}} \delta^{\frac{d}{2}}}$$

of the *projected normal distribution* $\Pi_{\mathcal{N}}(0, \Sigma)$ on $\mathbb{S}^{d-1}$ on $\mathbb{S}^{d-1}$. In the rest of this paper, we will refer to this setting as case $\nu = 0$. More precisely, if $X \sim \mathcal{N}(0, \Sigma)$, then $\frac{X}{\|X\|_2} \sim \Pi_{\mathcal{N}}(0, \Sigma)$. Note that for $X \sim \mathcal{N}(\mu, \Sigma)$ with $\mu \neq 0$ we have again $\frac{X-\mu}{\|X-\mu\|_2} \sim \Pi_{\mathcal{N}}(0, \Sigma)$, but $\frac{X}{\|X\|_2} \sim \Pi_{\mathcal{N}}(\mu, \Sigma)$ with a more sophisticated pdf, see, e.g., [15]. This distribution is also called *angular Gaussian distribution* [17, 32, 52], *off-set normal distribution* [31] or *displaced normal distribution* [18]. It is important to mention that $f_0(x|\Sigma) = f_0(x|\lambda\Sigma)$ for any $\lambda > 0$, so that the positive definite matrix $\Sigma$ is only identifiable up to a positive factor.

**Wrapped Cauchy Distribution.** For $d = 2$, there is a relation of the projected normal distribution to the *wrapped Cauchy distribution* [19, 32, 51, 50] which we will use for our applications in Section 6. The density of the real-valued Cauchy distribution $C(a, \gamma)$ is given by

$$g(\vartheta|a, \gamma) = \frac{1}{\pi} \frac{\gamma}{\gamma^2 + (\vartheta - a)^2}, \qquad a \in \mathbb{R}, \ \gamma > 0.$$

The wrapped Cauchy distribution $C(a, \gamma)$ is obtained by wrapping it around the circle, i.e. for $\vartheta \in [-\pi, \pi)$ we have

$$g_w(\vartheta|a, \gamma) = \sum_{k \in \mathbb{Z}} \frac{1}{\pi} \frac{\gamma}{\gamma^2 + (\vartheta + 2k\pi - a)^2}$$

$$= \frac{1}{2\pi} \frac{1 - \rho^2}{1 + \rho^2 - 2\rho\cos(\vartheta - a)}, \qquad a \in [-\pi, \pi),$$

where $\rho = \mathrm{e}^{-\gamma}$ and the second formula follows by Poisson's summation formula using that $\hat{g}(\omega) = \mathrm{e}^{-\gamma|\omega|+\mathrm{i}a\omega}$ is the characteristic function of $g$. We rewrite the density as

$$g_w(\vartheta|a, \gamma) = \frac{1}{2\pi} \frac{1}{\frac{1+\rho^2}{1-\rho^2} - \frac{2\rho}{1-\rho^2}\big(\cos(a)\cos(\vartheta) + \sin(a)\sin(\vartheta)\big)}$$

$$= \frac{1}{2\pi} \frac{\sqrt{1 - \xi_1^2 - \xi_2^2}}{1 - \xi_1 \cos(\vartheta) - \xi_2 \sin(\vartheta)}, \tag{2}$$

where $\xi_1 = \frac{2\rho}{1+\rho^2} \cos(a)$ and $\xi_2 = \frac{2\rho}{1+\rho^2} \sin(a)$.

**Lemma 2.2** (Relation between Projected Normal and Wrapped Cauchy Distribution)**.**

(i) *Let $d = 2$ and $X \sim \Pi_{\mathcal{N}}(0, \Sigma)$ with $\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{pmatrix}$ be a random variable in $\mathbb{S}^1$ with*

*parameterization $X = \begin{pmatrix} \cos(\Phi) \\ \sin(\Phi) \end{pmatrix}$. Then $\Theta := (2\Phi) \bmod 2\pi \sim C_w(a, \rho)$ with parameters*

$$\rho = \mathrm{e}^{-\gamma} = \left( \frac{\mathrm{tr}(\Sigma) - 2\sqrt{|\Sigma|}}{\mathrm{tr}(\Sigma) + 2\sqrt{|\Sigma|}} \right)^{\frac{1}{2}}, \tag{3}$$

$$a = \begin{cases} -\pi & \text{if } \sigma_{11} - \sigma_{22} = 0, \\ \arctan\left( \frac{2\sigma_{12}}{\sigma_{11} - \sigma_{22}} \right) & \text{if } \sigma_{11} - \sigma_{22} > 0, \\ \arctan\left( \frac{2\sigma_{12}}{\sigma_{11} - \sigma_{22}} \right) + \pi & \text{if } \sigma_{11} - \sigma_{22} < 0 \text{ and } \sigma_{12} \geq 0, \\ \arctan\left( \frac{2\sigma_{12}}{\sigma_{11} - \sigma_{22}} \right) - \pi & \text{if } \sigma_{11} - \sigma_{22} < 0 \text{ and } \sigma_{12} < 0. \end{cases} \tag{4}$$

(ii) *Let $\Theta \sim C_w(a, \rho)$ and let $\Xi$ be a discrete random variable with $\mathbb{P}(\Xi = -1) = \mathbb{P}(\Xi = 1) = \frac{1}{4}$, $\mathbb{P}(\Xi = 0) = \frac{1}{2}$ that is independent from $\Theta$. Then $\Phi = \left( \frac{\Theta}{2} + \pi\Xi \right) \bmod 2\pi \sim \Pi_{\mathcal{N}}(0, \Sigma)$, where (up to a positive factor)*

$$\Sigma = \begin{pmatrix} \frac{1}{2} + \frac{\rho}{1+\rho^2} \cos(a) & \frac{\rho}{1+\rho^2} \sin(a) \\ \frac{\rho}{1+\rho^2} \sin(a) & \frac{1}{2} - \frac{\rho}{1+\rho^2} \cos(a) \end{pmatrix}.$$

The proof of this Lemma can be found in the appendix.

Observe that the wrapped Cauchy distribution is unimodular, whereas the projected normal distribution is antipodally symmetric, which causes the mod operation. A similar relation as in statement (i) of the lemma can be found, e.g., in [32] without proof.

## 3. Weighted Likelihood Functions

In this section, we provide weighted log-likelihood functions for the Student-$t$ and projected normal distributions together with the equations characterizing their critical points.

**Student-$t$ distribution.**   Let $\nu > 0$. For $x_i \in \mathbb{R}^d$, $i = 1, \ldots, n$, the likelihood function of the the Student-$t$ distribution is given by

$$\mathcal{L}(\mu, \Sigma | x_1, \ldots, x_n) = \frac{\Gamma\left(\frac{d+\nu}{2}\right)^n}{\Gamma\left(\frac{\nu}{2}\right)^n (\pi\nu)^{\frac{nd}{2}} |\Sigma|^{\frac{n}{2}}} \prod_{i=1}^{n} \frac{1}{\left(1 + \frac{1}{\nu}\delta_i\right)^{\frac{d+\nu}{2}}}$$

and the log-likelihood function by

$$\ell(\mu, \Sigma | x_1, \ldots, x_n) = n \log\left(\Gamma\left(\frac{d+\nu}{2}\right)\right) - n \log\left(\Gamma\left(\frac{\nu}{2}\right)\right) - \frac{nd}{2} \log(\pi\nu)$$
$$- \frac{n}{2} \log|\Sigma| - \frac{d+\nu}{2} \sum_{i=1}^{n} \log\left(1 + \frac{1}{\nu}\delta_i\right),$$

where

$$\delta_i := (x_i - \mu)^{\mathrm{T}} \Sigma^{-1} (x_i - \mu).$$

Ignoring constants, maximizing $\ell$ is equivalent to minimizing the negative, weighted function

$$L(\mu, \Sigma) := (d + \nu) \sum_{i=1}^{n} w_i \log(\nu + \delta_i) + \log|\Sigma| \tag{5}$$

for uniform weights $w_i = \frac{1}{n}$. Note that we allow for different weightings of the summands by introducing weights in the open probability simplex

$$\mathring{\Delta}_n := \left\{ w = (w_1, \ldots, w_n) \in \mathbb{R}^n_{>0} : \sum_{i=1}^{n} w_i = 1 \right\}.$$

Using the relations

$$\frac{\partial \log(|X|)}{\partial X} = X^{-1}, \qquad \frac{\partial a^{\mathrm{T}} X^{-1} b}{\partial X} = -(X^{-\mathrm{T}}) a b^{\mathrm{T}} (X^{-\mathrm{T}}),$$

see [40], the derivatives of $L$ with respect to $\mu$ and $\Sigma$ are given by

$$\frac{\partial L}{\partial \mu}(\mu, \Sigma) = -2(d + \nu) \sum_{i=1}^{n} w_i \frac{\Sigma^{-1}(x_i - \mu)}{\nu + \delta_i},$$
$$\frac{\partial L}{\partial \Sigma}(\mu, \Sigma) = -(d + \nu) \sum_{i=1}^{n} w_i \frac{\Sigma^{-1}(x_i - \mu)(x_i - \mu)^{\mathrm{T}}\Sigma^{-1}}{\nu + \delta_i} + \Sigma^{-1}.$$

Setting them to zero results in the equations

$$0 = \sum_{i=1}^{n} w_i \frac{x_i - \mu}{\nu + \delta_i}, \tag{6}$$

$$I = (d + \nu) \sum_{i=1}^{n} w_i \frac{\Sigma^{-\frac{1}{2}}(x_i - \mu)(x_i - \mu)^{\mathrm{T}}\Sigma^{-\frac{1}{2}}}{\nu + \delta_i} \tag{7}$$

characterizing the *critical points* of $L$. Computing the trace of both sides of (7) and using the linearity and permutation invariance of the trace operator, we obtain

$$d = \mathrm{tr}(I) = (d + \nu) \sum_{i=1}^{n} w_i \frac{\mathrm{tr}\big(\Sigma^{-\frac{1}{2}}(x_i - \mu)(x_i - \mu)^{\mathrm{T}}\Sigma^{-\frac{1}{2}}\big)}{\nu + \delta_i} = (d + \nu) \sum_{i=1}^{n} w_i \frac{\delta_i}{\nu + \delta_i},$$

which yields after division by $\nu > 0$ the relation

$$1 = (d + \nu) \sum_{i=1}^{n} w_i \frac{1}{\nu + \delta_i}. \tag{8}$$

**Projected normal distribution.**  Let $\nu = 0$ and $d \geq 2$. Similarly as above, we obtain for $x_i \in \mathbb{S}^{d-1}$, $i = 1, \ldots, n$, the negative, weighted likelihood function of the $\Pi_{\mathcal{N}}(0, \Sigma)$ distribution as

$$L_0(\Sigma) := d \sum_{i=1}^{n} w_i \log(\delta_i) + \log |\Sigma| . \tag{9}$$

The critical points of $L_0$ are given by the solution of

$$I = d \sum_{i=1}^{n} w_i \frac{\Sigma^{-\frac{1}{2}} x_i x_i^{\mathrm{T}} \Sigma^{-\frac{1}{2}}}{\delta_i}. \tag{10}$$

Note that again $L_0(\lambda\Sigma) = L_0(\Sigma)$, $\lambda > 0$ and if $\Sigma$ fulfills (10) then also $\lambda\Sigma$ does. From the statistical point of view, (9) is only a likelihood function for samples on $\mathbb{S}^{d-1}$. However, later we will also consider the function for arbitrary nonzero points $x_i \in \mathbb{R}^d$.

## 4. Weighted Maximum Likelihood Estimators

In this section, we are interested in the minimizers of $L$ and $L_0$. We prove

- for $\nu > 0$ and fixed $\mu$ that $L$ has a unique critical point $\Sigma \in \mathrm{SPD}(d)$ and this point is a minimizer;

- for $\nu \geq 1$ that $L$ has a unique critical point $(\mu, \Sigma) \in \mathbb{R}^d \times \mathrm{SPD}(d)$ and this point is a minimizer;

- for $\nu = 0$ ($\mu = 0$) that $L_0$ has a unique critical point $\Sigma \in \mathrm{SPD}(d)$ with $\mathrm{tr}\,\Sigma = 1$ and this point is a minimizer. Moreover, all critical points are given by $\lambda\Sigma$, $\lambda > 0$.

8

Under several assumptions, these results are known even in the more general context of $M$-estimator for ellipical distributions. In [33], Maronna established sufficient conditions for the existence and uniqueness of a joint minimizer for uniform weights and $M$-estimators whose cost function fulfills certain properties. More results can be found in [20]. For the projected normal distribution, the claims were proved under certain assumptions by Tyler in [46], see also [12, 14, 47]. For an overview we refer to the survey paper of Dümbgen et al. [13].

In this paper, we incorporate different weights into the log-likelihood function which allows for multiple samples and is moreover useful in the nonlocal denoising approach in Section 6. We give direct existence and uniqueness proofs in order to make the paper self-contained.

## 4.1. Estimation of Scatter

First, we consider the estimation of the scatter matrix $\Sigma$ only. We start with the case $\nu > 0$, where the location parameter $\mu$ is known. If $\mu \neq 0$, we might transform the samples to $y_i = x_i - \mu$, $i = 1, \ldots, n$, so that we can assume w.l.o.g. $\mu = 0$ and use the notation $L(\Sigma) = L(0, \Sigma)$.

We make the following assumption on the samples $x_1, \ldots, x_n \in \mathbb{R}^d$ and weights $w \in \mathring{\Delta}_n$:

**Assumption 4.1.** *(i) Any subset of $\leq d$ samples $x_i$, $i \in \{1, \ldots, n\}$ is linearly independent.*

*(ii) $(d-1)w_{max} < \frac{\nu+d-1}{\nu+d}$, where $w_{max} := \max\{w_i : i = 1, \ldots, n\}$.*

The linear independence assumption (i) holds $\lambda^d$-a.s. when sampling from a continuous distribution. The interpretation behind the constraints is that the mass of the (empirical) distribution determined by $x_1, \ldots, x_n$ is not allowed to be concentrated in lower dimensional subspaces, which would cause the resulting distribution to be degenerated.

**Lemma 4.2.** *Let $x_i \in \mathbb{R}^d$, , $i = 1, \ldots, n$ and $w \in \mathring{\Delta}_n$ fulfill Assumption 4.1. Further, let $V \subset \mathbb{R}^d$ be a linear subspace with $0 \leq \dim(V) \leq d-1$ and $\mathcal{I}_V := \{i \in \{1, \ldots, n\}\colon x_i \in V\}$. Then it holds*

$$\sum_{i \in \mathcal{I}_V} w_i < \frac{\nu + \dim(V)}{\nu + d} \tag{11}$$

*and $n \geq d$.*

*Proof.* If $V = \{0\}$, we have $\mathcal{I}_V = \emptyset$, so that the statement holds true. Next, let $d \geq 2$ and $1 \leq k = \dim(V) \leq d-1$. By Assumption 4.1 (i), it holds $|\mathcal{I}_V| \leq k$. Using that by assumption $w_{\max} < \frac{\nu+d-1}{(d-1)(\nu+d)}$, we obtain

$$\sum_{i \in \mathcal{I}_V} w_i < k \frac{\nu+d-1}{(d-1)(\nu+d)} = \frac{\frac{k}{d-1}\nu + k}{\nu+d} \leq \frac{\nu+k}{\nu+d} = \frac{\nu + \dim(V)}{\nu+d}.$$

Assume that $n < d$ and let $V = \operatorname{span}\{x_1, \ldots, x_n\}$. Then $\dim(V) \leq d-1$, and we obtain the contradiction

$$1 = \sum_{i \in \mathcal{I}_V} w_i \leq \frac{\nu + \dim(V)}{\nu + d} \leq \frac{\nu + d - 1}{d + \nu} < 1. \qquad \square$$

Next, we show the existence of a minimizer of $L$. Since $\operatorname{SPD}(d)$ is an open cone, any minimizer of $L$ is also critical point of $L$.

**Theorem 4.3** (Existence of Scatter, $\nu > 0$). *Let $x_i \in \mathbb{R}^d$, $i = 1, \ldots, n$ and $w \in \mathring{\Delta}_n$ fulfill Assumption 4.1. Then it holds*

$$\operatorname*{argmin}_{\Sigma \in \operatorname{SPD}(d)} L(\Sigma) \neq \emptyset.$$

*Proof.* Let $\{\Sigma_r\}_{r \in \mathbb{N}} \subseteq \operatorname{SPD}(d)$ be a sequence in $\operatorname{SPD}(d)$, where $\lambda_{1r} \geq \ldots \geq \lambda_{dr} > 0$ denote the eigenvalues of $\Sigma_r$ and $e_{1r}, \ldots, e_{dr}$ the corresponding orthonormal eigenvectors. We prove that $L(\Sigma_r) \overset{r \to \infty}{\to} +\infty$ if one of the following situations is met:

(i) $\lambda_{1r} \overset{r \to \infty}{\to} +\infty$ and $\lambda_{dr} \geq c > 0$ for all $r \in \mathbb{N}$,

(ii) $\lambda_{dr} \overset{r \to \infty}{\to} 0$.

Then $L$ attains its minimum in $\operatorname{SPD}(d)$, and since $L$ is continuously differentiable, it is necessarily a critical point. By definition of $L$ we have

$$L(\Sigma_r) = (d + \nu) \sum_{i=1}^{n} w_i \log\left(\nu + x_i^{\mathrm{T}} \Sigma_r^{-1} x_i\right) + \sum_{j=1}^{d} \log(\lambda_{rj}).$$

In case (i), the first sum is bounded below, while the second one tends to infinity, so that $L(\Sigma_r) \overset{r \to \infty}{\to} +\infty$.

In case (ii), let $0 \leq p \leq d-1$ such that $\lambda_{1r} \geq \ldots \geq \lambda_{rp} \geq c > 0$ for all $r \in \mathbb{N}$ and $\lambda_{rp+1} \geq \ldots \geq \lambda_{dr} \overset{r \to \infty}{\to} 0$. If $p = 0$, then all $\lambda_{rj}$ tend to zero. Since the sphere $\mathbb{S}^{d-1}$ is compact, there exist subsequences (w.l.o.g. again denoted by $e_{rj}$) such that $\lim_{r \to \infty} e_{rj} = e_j \in \mathbb{S}^{d-1}$ for $j = 1, \ldots, d$. Set $S_0 = \{0\}$, and for $k = 1, \ldots, d$ further $S_k \coloneqq \operatorname{span}\{e_1, \ldots, e_k\}$ and

$$W_k \coloneqq S_k \setminus S_{k-1} = \left\{y \in \mathbb{R}^d \colon \langle y, e_k \rangle \neq 0, \ \langle y, e_l \rangle = 0 \text{ for } l = k+1, \ldots, d\right\}.$$

By Assumption 4.1(i) we know that $\dim(S_k) = k$. Let

$$\widetilde{I}_k \coloneqq \left\{i \in \{1, \ldots, n\} \colon x_i \in S_k\right\} \qquad \text{and} \qquad I_k \coloneqq \left\{i \in \{1, \ldots, n\} \colon x_i \in W_k\right\}.$$

Using $S_k = W_k \dot{\cup} S_{k-1}$, we obtain $\tilde{I}_k = I_k \dot{\cup} \tilde{I}_{k-1}$ for $k = 1, \ldots, d$. According to Assumption 4.1(i) it holds $|I_k| \leq |\tilde{I}_k| \leq \dim(S_k) = k$ for $k = 1, \ldots, d-1$. Introducing the functions

$$L_j(\Sigma_r) = (d + \nu) \sum_{i \in I_j} w_i \log\left(\nu + x_i^{\mathrm{T}} \Sigma_r^{-1} x_i\right) + \log(\lambda_{rj}),$$

we can split

$$L(\Sigma_r) = \sum_{j=1}^{p} L_j(\Sigma_r) + \sum_{j=p+1}^{d} L_j(\Sigma_r).$$

By assumption the first sum remains bounded from below as $r \to \infty$ so that we have to consider the second one. We have

$$\sum_{j=p+1}^{d} L_j(\Sigma_r) = \sum_{j=p+1}^{d} (d+\nu) \sum_{i\in I_j} w_i \log \left( \nu + x_i^{\mathrm{T}} \Sigma_r^{-1} x_i \right) + \log(\lambda_{rj})$$

$$= \sum_{j=p+1}^{d} \left( (d+\nu) \sum_{i\in I_j} w_i \log \left( \lambda_{rj}(\nu + x_i^{\mathrm{T}} \Sigma_r^{-1} x_i) \right) + \left( 1 - (d+\nu) \sum_{i\in I_j} w_i \right) \log(\lambda_{rj}) \right). \quad (12)$$

Since

$$y^{\mathrm{T}} \Sigma_r^{-1} y = \sum_{j=1}^{d} \frac{1}{\lambda_{rj}} \langle y, e_{rj} \rangle^2 \geq \frac{1}{\lambda_{rk}} \langle y, e_{rk} \rangle^2,$$

and $\langle y, e_{rk} \rangle \overset{r\to\infty}{\to} \langle y, e_k \rangle \neq 0$ for $y \in W_k$, we obtain for all $r$ sufficiently large

$$\liminf_{r\to\infty} \lambda_{rk}\, y^{\mathrm{T}} \Sigma_r^{-1} y \geq \langle y, e_k \rangle^2 > 0.$$

This implies that the first summand in (12) is bounded from below. Concerning the second summand in (12) we prove for sufficiently large $r$ by induction for $k \geq p+1$ that

$$\sum_{j=k}^{d} \left( 1 - (d+\nu) \sum_{i\in I_j} w_i \right) \log(\lambda_{rj}) \geq \left( (d+\nu) \sum_{i\in\tilde{I}_{k-1}} w_i - (\nu + k - 1) \right) \log(\lambda_{rk}). \quad (13)$$

Then, for $k = p+1$, we conclude by Lemma 4.2 that the factor of $\log(\lambda_{rp+1})$ on the right-hand side is negative, so that the whole sum tends to $+\infty$ as $r \to \infty$.

Based on the relation $\tilde{I}_k = I_k \cup \tilde{I}_{k-1}$ we write

$$\sum_{i\in\tilde{I}_k} w_i - \sum_{i\in I_k} w_i = \sum_{i\in\tilde{I}_{k-1}} w_i. \quad (14)$$

For the induction basis $k = d$, since $\lambda_{dr} < 1$ for sufficiently large $r$, we have to show

$$1 - (d+\nu) \sum_{i\in I_d} w_i \leq (d+\nu) \sum_{i\in\tilde{I}_{d-1}} w_i - (\nu + d - 1).$$

This follows directly from (14), since

$$1 - (d+\nu)\sum_{i\in I_d} w_i = 1 - (d+\nu)\Big(\sum_{i\in \tilde{I}_d} w_i - \sum_{i\in \tilde{I}_{d-1}} w_i\Big) = 1 - (d+\nu)\Big(1 - \sum_{i\in \tilde{I}_{d-1}} w_i\Big)$$

$$= (d+\nu)\sum_{i\in \tilde{I}_{d-1}} w_i - (\nu + d - 1).$$

Now, assume that (13) holds for some $k+1$ with $d \geq k+1 > p+1$, i.e.,

$$\sum_{j=k+1}^{d}\Big(1 - (d+\nu)\sum_{i\in I_j} w_i\Big)\log(\lambda_{rj}) \geq (d+\nu)\Big(\sum_{i\in \tilde{I}_k} w_i - \frac{\nu+k}{d+\nu}\Big)\log(\lambda_{rk+1}).$$

Then we get

$$\sum_{j=k}^{d}\Big(1 - (d+\nu)\sum_{i\in I_j} w_i\Big)\log(\lambda_{rj})$$

$$= \sum_{j=k+1}^{d}\Big(1 - (d+\nu)\sum_{i\in I_j} w_i\Big)\log(\lambda_{rj}) + \Big(1 - (d+\nu)\sum_{i\in I_k} w_i\Big)\log(\lambda_{rk})$$

$$\geq (d+\nu)\Big(\sum_{i\in \tilde{I}_k} w_i - \frac{\nu+k}{d+\nu}\Big)\log(\lambda_{rk+1}) + \Big(1 - (d+\nu)\sum_{i\in I_k} w_i\Big)\log(\lambda_{rk})$$

and by Lemma 4.2 and since $\lambda_{rk+1} \leq \lambda_{rk} < 1$ finally

$$\geq (d+\nu)\Big(\sum_{i\in \tilde{I}_k} w_i - \frac{\nu+k}{d+\nu}\Big)\log(\lambda_{rk}) + \Big(1 - (d+\nu)\sum_{i\in I_k} w_i\Big)\log(\lambda_{rk})$$

$$= \Big((d+\nu)\sum_{i\in \tilde{I}_{k-1}} w_i - (\nu + k - 1)\Big)\log(\lambda_{rk}).$$

This finishes the proof. $\qquad\square$

The end of the proof of Theorem 4.3 reveals that the condition on the weights stated in (11) is sufficient for the existence of a minimizer. The next lemma shows that the non strong inequality is necessary for the existence of a critical point.

**Lemma 4.4.** *Let* $x_i \in \mathbb{R}^d$, $i = 1, \ldots, n$ *fulfill Assumption 4.1 (i) and assume there exists a critical point of* $L$. *Then, for all linear subspaces* $V \subset \mathbb{R}^d$ *with* $0 \leq \dim(V) \leq d-1$ *it holds*

$$\sum_{i\in \mathcal{I}_V} w_i \leq \frac{\nu + \dim(V)}{d+\nu},$$

where $\mathcal{I}_V = \{i \in \{1, \ldots, n\} : x_i \in V\}$.

*Proof.* Condition (7) with $\mu = 0$ can be alternatively written as

$$I = (d + \nu) \sum_{i=1}^{n} w_i \frac{R^{-1} x_i \, x_i^{\mathrm{T}} R^{-\mathrm{T}}}{\nu + \delta_i} \tag{15}$$

with the Cholesky decomposition $\Sigma = RR^{\mathrm{T}}$. W.l.o.g. we might assume that $\Sigma = I$ is the critical point, otherwise we can transform the samples to $y_i = R^{-1} x_i$, where $RR^{\mathrm{T}} = \Sigma$. The idea of the proof is to project the samples onto the orthogonal complement of $V$. More precisely, let $k = \dim(V) < d$ and choose an orthonormal basis $v_1, \ldots, v_k$ of $V$. Set $W = (v_1, \ldots, v_k)$ so that $P = WW^{\mathrm{T}}$ is the orthogonal projection onto $V$. Now, for $\Sigma = I$ and $\mu = 0$, equation (7) reads as

$$I = (d + \nu) \sum_{i=1}^{n} w_i \frac{x_i x_i^{\mathrm{T}}}{\nu + x_i^{\mathrm{T}} x_i}.$$

Multiplying both sides with $I - P$ and taking the trace, yields

$$d - k = (d + \nu) \sum_{i=1}^{n} w_i \frac{x_i^{\mathrm{T}}(I - P)x_i}{\nu + x_i^{\mathrm{T}} x_i}.$$

We split the sum into a sum over $i \in \mathcal{I}_V$ and $i \in \mathcal{I}_V^c = \{1, \ldots, n\} \setminus \mathcal{I}_V$ and get

$$d - k = (d + \nu) \sum_{i \in \mathcal{I}_V} w_i \frac{x_i^{\mathrm{T}}(I - P)x_i}{\nu + x_i^{\mathrm{T}} x_i} + (d + \nu) \sum_{i \in \mathcal{I}_V^c} w_i \frac{x_i^{\mathrm{T}}(I - P)x_i}{\nu + x_i^{\mathrm{T}} x_i}$$

$$= (d + \nu) \sum_{i \in \mathcal{I}_V^c} w_i \frac{x_i^{\mathrm{T}}(I - P)x_i}{\nu + x_i^{\mathrm{T}} x_i}.$$

With $x_i^{\mathrm{T}}(I - P)x_i \leq x_i^{\mathrm{T}} x_i$ we obtain further

$$d - k = (d + \nu) \sum_{i \in \mathcal{I}_V^c} w_i \frac{x_i^{\mathrm{T}}(I - P)x_i}{\nu + x_i^{\mathrm{T}} x_i} \leq (d + \nu) \sum_{i \in \mathcal{I}_V^c} w_i \frac{x_i^{\mathrm{T}} x_i}{\nu + x_i^{\mathrm{T}} x_i} \leq (d + \nu) \sum_{i \in \mathcal{I}_V^c} w_i$$

$$= d + \nu - (d + \nu) \sum_{i \in \mathcal{I}_V} w_i.$$

Rearranging yields

$$(d + \nu) \sum_{i \in \mathcal{I}_V} w_i \leq \nu + k$$

and finally
$$\sum_{i \in \mathcal{I}_V} w_i \leq \frac{\nu + k}{d + \nu} = \frac{\nu + \dim(V)}{d + \nu}.$$

$\square$

Now, we turn to the question whether $L$ has a unique critical point. If this is the case, then clearly the unique critical point of $L$ is the minimizer of $L$.

**Theorem 4.5** (Uniqueness of Scatter, $\nu > 0$). *Let $x_i \in \mathbb{R}^d$, $i = 1, \ldots, n$ and $w \in \mathring{\Delta}_n$ fulfill Assumption 4.1. Then $L$ has a unique critical point.*

*Proof.* Let $\Sigma_1$ and $\Sigma_2$ fulfill (15). W.l.o.g. we can assume that $\Sigma_1 = I$, since otherwise we can use the Cholesky factorization $\Sigma_1 = R_1 R_1^{\mathrm{T}}$, transform the data to $y_i = R_1^{-1} x_i$ and replace $\Sigma_2$ by $R_1^{-1} \Sigma_2 R_1^{-\mathrm{T}}$.

Let $\lambda_1 \geq \ldots \geq \lambda_d$ denote the eigenvalues of $\Sigma_2$. We show that $\lambda_1 \leq 1$ and $\lambda_d \geq 1$ which implies $I = \Sigma_2$. Assume that $\lambda_1 > 1$. By (15) it holds

$$\Sigma_2 = (d + \nu) \sum_{i=1}^{n} w_i \frac{x_i x_i^{\mathrm{T}}}{\nu + x_i^{\mathrm{T}} \Sigma_2^{-1} x_i}. \tag{16}$$

Using the Courant-Fisher min-max principle we have

$$x_i^{\mathrm{T}} \Sigma_2^{-1} x_i \geq \lambda_1^{-1} x_i^{\mathrm{T}} x_i,$$

and with $\lambda_1 > 1$ we can estimate

$$\frac{d + \nu}{\nu + x_i^{\mathrm{T}} \Sigma_2^{-1} x_i} \leq \frac{d + \nu}{\nu + \lambda_1^{-1} x_i^{\mathrm{T}} x_i} = \lambda_1 \frac{d + \nu}{\lambda_1 \nu + x_i^{\mathrm{T}} x_i} < \lambda_1 \frac{d + \nu}{\nu + x_i^{\mathrm{T}} x_i}, \qquad i = 1, \ldots, n.$$

Inserting this in (16) and regarding Assumption 4.1(i) yields the contradiction

$$\Sigma_2 = (d + \nu) \sum_{i=1}^{n} w_i \frac{x_i x_i^{\mathrm{T}}}{\nu + x_i^{\mathrm{T}} \Sigma_2^{-1} x_i} \prec \lambda_1 (d + \nu) \sum_{i=1}^{n} w_i \frac{x_i x_i^{\mathrm{T}}}{\nu + x_i^{\mathrm{T}} x_i} = \lambda_1 I,$$

where $A \prec B$ means that $B - A$ is positive definite. Similarly we can show that $\lambda_d \geq 1$ and we are done. $\square$

Now we turn to the case $\nu = 0$, i.e. to the projected normal distribution and the function $L_0$ in (9).

**Theorem 4.6** (Existence of Scatter, $\nu = 0$). *Let $d \geq 2$ and $x_i \in \mathbb{R}^d$, $i = 1, \ldots, n$ fulfill*

*Assumption 4.1 (i). Let $w \in \mathring{\Delta}_n$ satisfy $w_{\max} < 1/d$. Then, it holds*

$$\operatorname*{argmin}_{\Sigma \in \mathrm{SPD}(d)} L_0(\Sigma) \neq \emptyset.$$

Note that for $x_i \in \mathbb{S}^{d-1}$, $i = 1, \dots, n$, the Assumption 4.1 i) is fulfilled if the points are pairwise distinct and not antipodal. Further, we see that $w \in \mathring{\Delta}_n$ and $w_{\max} < 1/d$ imply $n > d$.

*Proof.* By definition of $L_0$ we can restrict the domain of $L_0$ to the bounded set

$$D := \{\Sigma \in \mathrm{SPD}(d)\colon \operatorname{tr}(\Sigma) = d\} \subset \mathrm{Sym}(d).$$

Let $\{\Sigma_r\}_{r \in \mathbb{N}} \subseteq D$ be a sequence with $\lim_{r \to \infty} \Sigma_r = \Sigma \in \partial D$. We will show that $L_0(\Sigma_r) \overset{r \to \infty}{\to} \infty$. Let $\lambda_{1r} \geq \dots \geq \lambda_{dr} > 0$ denote the eigenvalues of $\Sigma_r$. Further, let $\lambda_1 \geq \dots \geq \lambda_p > 0$, $\lambda_{p+1} = \dots = \lambda_d = 0$ be the eigenvalues of $\Sigma$, where $p = \operatorname{rank}(\Sigma) < d$. Note that $p \geq 1$ since $\operatorname{tr}(\Sigma) = d$. Now, the same argumentation as in the proof of case (ii) of Theorem 4.3 yields

$$L_0(\Sigma_r) = \sum_{j=1}^{d} L_j(\Sigma_r) = \sum_{j=1}^{p} L_j(\Sigma_r) + \sum_{j=p+1}^{d} L_j(\Sigma_r)$$

$$\geq \sum_{j=1}^{p} L_j(\Sigma_r) + d \sum_{j=p+1}^{d} \sum_{i \in I_j} w_i \log(\lambda_{rj} x_i^{\mathrm{T}} \Sigma_r^{-1} x_i) + \left(d \sum_{i \in \tilde{I}_p} w_i - p\right) \log(\lambda_{rp+1}^{-1}),$$

which tends to $+\infty$ for $r \to \infty$. Consequently, since $\{\Sigma_r\}_{r \in \mathbb{N}}$ is bounded, it contains a convergent subsequence, whose limit is by the above argumentation an inner point. $\square$

The next lemma shows, that the solution of (10) is no longer unique. However, there exists a unique solution $\Sigma$ with trace $d$, and all other solutions are of the form $\lambda \Sigma$ with $\lambda > 0$.

**Theorem 4.7** (Uniqueness of Scatter up to a factor, $\nu = 0$)**.** *Let $d \geq 2$ and $x_i \in \mathbb{R}^d$, $i = 1, \dots, n$ fulfill Assumption 4.1 (i). Let $w \in \mathring{\Delta}_n$ satisfy $w_{\max} < 1/d$. Then, it holds*

$$\sum_{i=1}^{n} w_i \frac{\Sigma^{-\frac{1}{2}} x_i x_i^{\mathrm{T}} \Sigma^{-\frac{1}{2}}}{x_i^{\mathrm{T}} \Sigma^{-1} x_i} = \sum_{i=1}^{n} w_i \frac{S^{-\frac{1}{2}} x_i x_i^{\mathrm{T}} S^{-\frac{1}{2}}}{x_i^{\mathrm{T}} \Sigma^{-1} x_i} \tag{17}$$

*if and only if $S = \lambda \Sigma$ for some $\lambda > 0$. The critical points of $L_0$ are unique up to a positive factor.*

*Proof.* Clearly, if $S = \lambda \Sigma$ then (17) holds true.
To show the reverse direction, we may as in the proof of Theorem 4.5 w.l.o.g. assume that $\Sigma = I$. Let $\lambda_1$ denote the largest eigenvalue of $S^{-1}$ and assume that it has multiplicity $k < d$.

Let $e_1, \ldots, e_k$ be the corresponding orthonormal eigenvectors. Further, let $P = \sum_{i=1}^{k} e_i e_i^{\mathrm{T}}$ be the associated orthogonal projector onto $\mathrm{span}\{e_1, \ldots, e_k\}$. Equality (17) reads for $\Sigma = I$ as

$$\sum_{i=1}^{n} w_i \frac{x_i x_i^{\mathrm{T}}}{x_i^{\mathrm{T}} x_i} = \sum_{i=1}^{n} w_i \frac{S^{-\frac{1}{2}} x_i x_i^{\mathrm{T}} S^{-\frac{1}{2}}}{x_i^{\mathrm{T}} S^{-1} x_i}.$$

Multiplying both sides with $P$ and taking the trace gives

$$\sum_{i=1}^{n} w_i \frac{x_i^{\mathrm{T}} P x_i}{x_i^{\mathrm{T}} x_i} = \lambda_1 \sum_{i=1}^{n} w_i \frac{x_i^{\mathrm{T}} P x_i}{x_i^{\mathrm{T}} S^{-1} x_i}.$$

Since $\frac{1}{\lambda_1} x_i^{\mathrm{T}} S^{-1} x_i \leq x_i^{\mathrm{T}} x_i$ with equality if and only if $P x_i = x_i$, this implies $P x_i = 0$ or $P x_i = x_i$ for all $i = 1, \ldots, n$. This contradicts the assumptions on the points $x_i$ and $w_{\max}$ unless $P = I$. Thus, $S = \lambda_1 I$. $\qquad\square$

## 4.2. Estimation of Location and Scatter

In order to show the existence and uniqueness of a joint minimizer of the likelihood function $L(\mu, \Sigma)$, we use the established technique to consider the location and scatter estimation problem as a higher dimensional centered scatter only problem. We need the following auxiliary lemma, see, e.g. [13].

**Lemma 4.8.** *Let $\lambda > 0$, $\mu \in \mathbb{R}^d$ and $\Sigma \in \mathrm{SPD}(d)$. Then it holds*

$$A_\lambda = \lambda \begin{pmatrix} \Sigma + \mu\mu^{\mathrm{T}} & \mu \\ \mu^{\mathrm{T}} & 1 \end{pmatrix} \in \mathrm{SPD}(d+1). \tag{18}$$

*Conversely, every $A \in \mathrm{SPD}(d+1)$ can be written in the form (18) with uniquely determined $\lambda > 0$, $\mu \in \mathbb{R}^d$ and $\Sigma \in \mathrm{SPD}(d)$.*

The inverse of the matrix $A_\lambda$ in (18) is given by

$$A_\lambda^{-1} = \frac{1}{\lambda} \begin{pmatrix} \Sigma^{-1} & -\Sigma^{-1}\mu \\ -\mu^{\mathrm{T}}\Sigma^{-1} & 1 + \mu^{\mathrm{T}}\Sigma^{-1}\mu \end{pmatrix}.$$

It fulfills

$$(x^{\mathrm{T}}\ 1)\, A_\lambda^{-1} \begin{pmatrix} x \\ 1 \end{pmatrix} = \frac{1}{\lambda}\big(1 + (x - \mu)^{\mathrm{T}}\Sigma^{-1}(x - \mu)\big) \tag{19}$$

and $|A| = \lambda^{d+1}|\Sigma|$. Using the last two relations and setting

$$A := \begin{pmatrix} \Sigma + \mu\mu^{\mathrm{T}} & \mu \\ \mu^{\mathrm{T}} & 1 \end{pmatrix}, \quad z_i := \begin{pmatrix} x_i \\ 1 \end{pmatrix}, \quad i = 1, \ldots, n, \tag{20}$$

we can rewrite the function $L(\mu, \Sigma)$ with $\nu > 1$ in (5) as

$$L(\mu, \Sigma) = (d + \nu) \sum_{i=1}^{n} w_i \log\big(\nu - 1 + z_i^{\mathrm{T}} A^{-1} z_i\big) + \log |A|$$

$$= (\tilde{d} + \tilde{\nu}) \sum_{i=1}^{n} w_i \log\big(\tilde{\nu} + z_i^{\mathrm{T}} A^{-1} z_i\big) + \log |A| =: \tilde{L}(A), \qquad (21)$$

where $\tilde{d} := d + 1$, $\tilde{\nu} := \nu - 1$ and the tilde on top of $L$ shows that the function is considered in dimension $d + 1$ now.

The points $x_i \in \mathbb{R}^d$, $i = 1, \ldots, k$ are called *affinely independent* if $\sum_{i=1}^{k} \lambda_i x_i = 0$ and $\sum_{i=1}^{k} \lambda_i = 0$ implies $\lambda_1 = \ldots = \lambda_k = 0$. It can be immediately seen that the points $x_i \in \mathbb{R}^d$, $i = 1, \ldots, k$ are affinely independent if and only if the points $z_i = (x_i^{\mathrm{T}} 1)^{\mathrm{T}} \in \mathbb{R}^{d+1}$, $i = 1, \ldots, k$ are linearly independent.

Let $\nu > 1$. By Theorem 4.5 we know that $\tilde{L}(A)$ has a unique minimizer in $\mathrm{SPD}(d+1)$ if the following conditions according to Assumption 4.1 are fulfilled:

**Assumption 4.9.**     *(i) Any subset of $\leq d + 1$ samples $x_i$, $i \in \{1, \ldots, n\}$ is affinely independent.*

    *(ii) $d w_{max} < \frac{\nu + d - 1}{\nu + d}$.*

Note that $w \in \mathring{\Delta}_n$ and (ii) imply $n \geq d + 1$ for $\nu > 1$ and $n > d + 1$ for $\nu = 1$.

By the next lemma, see, e.g. [13], we can restrict our attention to matrices of the form (20) when minimizing $\tilde{L}$ over $\mathrm{SPD}(d+1)$.

**Lemma 4.10.** *Let $\nu > 1$ and $x_i \in \mathbb{R}^d$, $i = 1, \ldots, n$ and $w \in \mathring{\Delta}_n$ fulfill Assumption 4.9. For $z_i := (x_i^{\mathrm{T}} 1)^{\mathrm{T}}$, $i = 1, \ldots, n$, let the function $\tilde{L}$ be defined by (21). Then the critical point of $\tilde{L}$ has the $(d+1, d+1)$-th entry 1.*

Now the existence and uniqueness of a minimizer of $L$ follows immediately from Lemmata 4.8 and 4.10 and (21).

**Theorem 4.11** (Existence and Uniqueness of Location and Scatter, $\nu > 1$). *Let $x_i \in \mathbb{R}^d$, $i = 1, \ldots, n$ and $w \in \mathring{\Delta}_n$ fulfill Assumption 4.9. Then $L$ has a unique critical point and this point is a minimizer of $L$.*

*Proof.* Let $A$ be the critical point/minimizer of $\tilde{L}$. By Lemma 4.10 it has the form (20) and by (21), we see that the corresponding $(\mu, \Sigma)$ is the minimizer of $L$ and thus a critical point of $L$.

Conversely, let $(\mu, \Sigma)$ be a critical point of $L$, i.e. fulfill (6) and (7). Define $A$ by (20). We have to show that $A$ fulfills the critical point condition

$$A = (\tilde{d} + \tilde{\nu}) \sum_{i=1}^{n} w_i \frac{z_i z_i^{\mathrm{T}}}{\nu - 1 + z_i^{\mathrm{T}} A^{-1} z_i}, \tag{22}$$

where $z_i := (x_i^{\mathrm{T}} \, 1)^{\mathrm{T}}$, $i = 1, \ldots, n$. Since $\tilde{L}$ has only one critical point, this would imply the assertion. By (19) and definition of $z_i$, this can be written as

$$A = (d + \nu) \sum_{i=1}^{n} w_i \frac{\begin{pmatrix} x_i x_i^{\mathrm{T}} & x_i \\ x_i^{\mathrm{T}} & 1 \end{pmatrix}}{\nu + \delta_i}.$$

Considering the different blocks of this matrix, we obtain with (6), (7) and (8) that

$$(d + \nu) \sum_{i=1}^{n} w_i \frac{x_i x_i^{\mathrm{T}}}{\nu + \delta_i} = \Sigma + (d + \nu) \sum_{i=1}^{n} w_i \frac{x_i \mu^{\mathrm{T}} + \mu x_i^{\mathrm{T}} - \mu \mu^t T}{\nu + \delta_i} = \Sigma + \mu \mu^{\mathrm{T}},$$

$$(d + \nu) \sum_{i=1}^{n} w_i \frac{x_i}{\nu + \delta_i} = \mu, \quad \text{and} \quad (d + \nu) \sum_{i=1}^{n} w_i \frac{1}{\nu + \delta_i} = 1.$$

Thus, (22) holds indeed true. □

The multivariate Cauchy distribution, i.e. the case $\nu = 1$ requires a special consideration.

**Theorem 4.12** (Existence and Uniqueness of Location and Scatter, $\nu = 1$). *Let $x_i \in \mathbb{R}^d$, $i = 1, \ldots, n$ and $w \in \mathring{\Delta}_n$ fulfill Assumption 4.9. Then $L$ has a unique critical point and this point is a minimizer of $L$.*

*Proof.* As in (21) we see also for $\nu = 1$ with $A$ and $(\mu, \Sigma)$ related by (20) that $L(\mu, \Sigma) = \tilde{L}_0(A)$, where $\tilde{L}_0$ resembles $L_0$ with respect to dimension $d + 1$. Let $A$ be the unique critical point of $\tilde{L}_0$ with $a_{d+1,d+1} = 1$. It is also a minimizer. Then, by the above relation, $(\mu, \Sigma)$ obtained from $A$ by (20) is a minimizer of $L$ and also a critical point. Conversely, let $(\mu, \Sigma)$ be any critical point of $L$. Then we can show as in the proof of Theorem 4.11 that the related $A$ in (20) is the unique critical point of $\tilde{L}_0$ with $a_{d+1,d+1} = 1$. This finishes the proof. □

## 5. Efficient Minimization Algorithm

In this section, we propose an efficient algorithm to compute the ML estimates of the multivariate Student-$t$ distribution and prove its convergence. After finishing this paper we became aware that a Gauss-Seidel variant of our algorithm was already suggested by Kent et al. [21] to speed up the classical EM algorithm without any convergence considerations. In

[49], see also [36], van Dyk gave an interpretation of this algorithm from an EM theoretical point of view using a sophisticated choice of the hidden variables, such that convergence is ensured by general results for the EM algorithm. However, our derivation of the Jacobi variant of the algorithm and its convergence proof do not rely on an EM setting.

We assume that $\nu \geq 0$ in case of estimating only $\Sigma$ and $\nu \geq 1$ when estimating both $\mu$ and $\Sigma$. Conditions (6) and (7) can be reformulated as fixed-point equations

$$\mu = \frac{\sum\limits_{i=1}^{n} w_i \frac{1}{\nu+\delta_i} x_i}{\sum\limits_{i=1}^{n} w_i \frac{1}{\nu+\delta_i}}, \tag{23}$$

$$\Sigma = (d + \nu) \sum\limits_{i=1}^{n} w_i \frac{(x_i - \mu)(x_i - \mu)^{\mathrm{T}}}{\nu + \delta_i}. \tag{24}$$

Based on this, we propose the following Algorithm 1 which is a Picard iteration for $\mu$ and a *scaled version* of the Picard iteration for $\Sigma$. For $d = 1$, Algorithm 1 coincides with the generalized myriad filtering considered in [26].

---

**Algorithm 1** Generalized Multivariate Myriad Filter (GMMF)

---

**Input:** $x_1, \ldots, x_n \in \mathbb{R}^d$, $w \in \mathring{\Delta}_n$

**Initialization:** $\mu_0 = \frac{1}{n} \sum\limits_{i=1}^{n} x_i$, $\Sigma_0 = \frac{1}{n} \sum\limits_{i=1}^{n} (x_i - \mu_0)(x_i - \mu_0)^{\mathrm{T}}$

**for** $r = 0, \ldots$ **do**

$$\delta_{i,r} = (x_i - \mu_r)^{\mathrm{T}} \Sigma_r^{-1} (x_i - \mu_r)$$

$$\mu_{r+1} = \frac{\sum\limits_{i=1}^{n} w_i \frac{1}{\nu+\delta_{i,r}} x_i}{\sum\limits_{i=1}^{n} w_i \frac{1}{\nu+\delta_{i,r}}}$$

$$\Sigma_{r+1} = \frac{\sum\limits_{i=1}^{n} w_i \frac{(x_i-\mu_r)(x_i-\mu_r)^{\mathrm{T}}}{\nu+\delta_{i,r}}}{\sum\limits_{i=1}^{n} w_i \frac{1}{\nu+\delta_{i,r}}}$$

---

The classical EM algorithm for maximizing the log-likelihood function [29] with equal weights $w_i = \frac{1}{n}$, $i = 1, \ldots, n$, leads to the same iteration for $\mu$, but the iteration of $\Sigma$ reads as

$$\Sigma_{r+1} = \frac{\nu + d}{n} \sum\limits_{i=1}^{n} \frac{(x_i - \mu_{r+1})(x_i - \mu_{r+1})^{\mathrm{T}}}{\nu + \delta_{i,r}}.$$

In the variant of Kent et al., this expression is divided by $\sum_{i=1}^{n} w_i \frac{1}{\nu+\delta_{i,r}}$. Despite the nonuniform weights, our Algorithm 1 differs from this update by taking $\mu_r$ instead of $\mu_{r+1}$ on the right-hand side and can therefore be seen as its Jacobi variant.

In the following, we focus on estimating both $\mu$ and $\Sigma$, where we assume that $\nu \geq 1$. However, the algorithm without the iteration in $\mu$ and its convergence proof work also for estimating only $\Sigma$ with $\nu \geq 0$. In case $\nu = 0$, we have to assume that $x_i \in \mathbb{S}^{d-1}$, $i = 1, \ldots, n$. Then the iteration for $\Sigma$ is the well-known Tyler $M$-estimator [46] and all iterates have trace 1.

## 5.1. Convergence Analysis

Since Algorithm 1 cannot be interpreted as an EM algorithm, we prove its convergence in the following. We note that the following Lemmata 5.1 and 5.2 on the monotone descent of the log likelihood function could be also deduced from convergence results of the EM algorithm, but one obtains different estimates. However, parts of those proofs are needed to prove Theorem 5.3.

Using again $\Sigma = RR^{\mathrm{T}}$, we introduce the notation $y_i := R^{-1}(x_i - \mu)$, $i = 1, \ldots, n$ and

$$S_0(\mu, \Sigma) := (d+\nu) \sum_{i=1}^{n} w_i \frac{1}{\nu + \delta_i},$$

$$S_1(\mu, \Sigma) := (d+\nu) \sum_{i=1}^{n} w_i \frac{y_i}{\nu + \delta_i},$$

$$S_2(\mu, \Sigma) := (d+\nu) \sum_{i=1}^{n} w_i \frac{y_i y_i^{\mathrm{T}}}{\nu + \delta_i},$$

where the subscripts 0, 1 and 2 stand for 'scalar', 'vector' and 'matrix', respectively. Further we abbreviate $S_{jr} := S_j(\mu_r, \Sigma_r)$, $j = 0, 1, 2$ and use this notation also if $\mu$ or $\Sigma$ is fixed. By the Cholesky decomposition $\Sigma_r = R_r R_r^{\mathrm{T}}$, the iterations in Algorithm 1 read as

$$\mu_{r+1} = \mu_r + R_r \frac{S_{1r}}{S_{0r}} = \mu_r - \frac{1}{2 S_{0r}} \Sigma_r \frac{\partial L}{\partial \mu}(\mu_r, \Sigma_r),$$

$$\Sigma_{r+1} = R_{r+1} R_{r+1}^{\mathrm{T}} = R_r \frac{S_{2r}}{S_{0r}} R_r^{\mathrm{T}} = R_r \frac{S_{2r}}{\frac{d+\nu}{\nu} - \frac{1}{\nu} \mathrm{tr}(S_{2r})} R_r^{\mathrm{T}} = R_r \frac{\nu S_{2r}}{d + \nu - \mathrm{tr}(S_{2r})} R_r^{\mathrm{T}}$$

$$= \frac{1}{S_{0r}} \left( \Sigma_r - \Sigma_r \frac{\partial L}{\partial \Sigma}(\mu_r, \Sigma_r) \Sigma_r \right). \tag{25}$$

We will need the difference of two iterates

$$L_\nu(\mu_{r+1}, \Sigma_{r+1}) - L_\nu(\mu_r, \Sigma_r) = (d+\nu) \sum_{i=1}^{n} w_i \log \left( \frac{\nu + \delta_{i,r+1}}{\nu + \delta_{i,r}} \frac{|\Sigma_{r+1}|^{\frac{1}{d+\nu}}}{|\Sigma_r|^{\frac{1}{d+\nu}}} \right). \tag{26}$$

We start with the convergence analysis of our algorithm for fixed $\mu$, where we might assume w.l.o.g. that $\mu = 0$

**Lemma 5.1.** *Let $\nu \geq 0$, $\mu = 0$ and $x_i \neq 0$, $i = 1, \ldots, n$. Then $\{\Sigma_r\}_{r \in \mathbb{N}}$ defined by the iterations in Algorithm 1 satisfy*

$$
\begin{aligned}
L(\Sigma_{r+1}) - L(\Sigma_r) &\leq 0 \quad \text{for} \quad \nu > 0, \\
L_0(\Sigma_{r+1}) - L_0(\Sigma_r) &\leq 0 \quad \text{for} \quad \nu = 0,
\end{aligned}
$$

*with equality if and only if $\Sigma_{r+1} = \Sigma_r$.*

*Proof.* We consider $\nu > 0$. The proof for $L_0$ follows exactly the same lines with $\nu = 0$. By concavity of the logarithm and (26) we have

$$
L(\Sigma_{r+1}) - L(\Sigma_r) \leq (d + \nu) \log \underbrace{\left( \sum_{i=1}^n w_i \frac{\nu + \delta_{i,r+1}}{\nu + \delta_{i,r}} \frac{|\Sigma_{r+1}|^{\frac{1}{d+\nu}}}{|\Sigma_r|^{\frac{1}{d+\nu}}} \right)}_{=\Upsilon},
$$

so that it suffices to show that $\Upsilon \leq 1$. Using properties of the determinant it holds

$$
\frac{|\Sigma_{r+1}|^{\frac{1}{d+\nu}}}{|\Sigma_r|^{\frac{1}{d+\nu}}} = \frac{\left| R_r \frac{S_{2r}}{S_{0r}} R_r^{\mathrm{T}} \right|^{\frac{1}{d+\nu}}}{|R_r R_r^{\mathrm{T}}|^{\frac{1}{d+\nu}}} = S_{0r}^{-\frac{d}{d+\nu}} |S_{2r}|^{\frac{1}{d+\nu}}.
$$

Next, we consider the term

$$
\sum_{i=1}^n w_i \frac{\nu + \delta_{i,r+1}}{\nu + \delta_{i,r}} = \sum_{i=1}^n w_i \frac{\delta_{i,r+1}}{\nu + \delta_{i,r}} + \frac{\nu}{d + \nu} S_{0r}.
$$

Using

$$
\delta_{i,r+1} = \mathrm{tr}\left( x_i^{\mathrm{T}} \Sigma_{r+1}^{-1} x_i \right) = \mathrm{tr}\left( x_i^{\mathrm{T}} R_r^{-\mathrm{T}} S_{0r} S_{2r}^{-1} R_r^{-1} x_i \right) = S_{0r} \, \mathrm{tr}\left( S_{2r}^{-1} R_r^{-1} x_i x_i^{\mathrm{T}} R_r^{-\mathrm{T}} \right)
$$

and the linearity of the trace, the sum simplifies to

$$
\begin{aligned}
\sum_{i=1}^n w_i \frac{\delta_{i,r+1}}{\nu + \delta_{i,r}} &= \sum_{i=1}^n w_i \frac{\mathrm{tr}\left( \delta_{i,r+1} \right)}{\nu + \delta_{i,r}} = S_{0r} \sum_{i=1}^n w_i \frac{\mathrm{tr}\left( S_{2r}^{-1} R_r^{-1} x_i x_i^{\mathrm{T}} R_r^{-\mathrm{T}} \right)}{\nu + \delta_{i,r}} \\
&= S_{0r} \, \mathrm{tr}\left( S_{2r}^{-1} \underbrace{\sum_{i=1}^n w_i \frac{\left( R_r^{-1} x_i x_i^{\mathrm{T}} R_r^{-\mathrm{T}} \right)}{\nu + \delta_{i,r}}}_{=\frac{1}{d+\nu} S_{2r}} \right) \\
&= \frac{1}{d + \nu} S_{0r} \, \mathrm{tr}(I) = \frac{d}{d + \nu} S_{0r}.
\end{aligned}
$$

Thus, we obtain

$$\Upsilon = \sum_{i=1}^{n} w_i \frac{\nu + \delta_{i,r+1}}{\nu + \delta_{i,r}} \frac{|\Sigma_{r+1}|^{\frac{1}{d+\nu}}}{|\Sigma_r|^{\frac{1}{d+\nu}}} = S_{0r} S_{0r}^{-\frac{d}{d+\nu}} |S_{2r}|^{\frac{1}{d+\nu}} = (S_{0r}^{\nu} |S_{2r}|)^{\frac{1}{d+\nu}} .$$

We have $\nu S_{0r} + \mathrm{tr}(S_{2r}) = d + \nu$. (If $\nu = 0$ and $n > d$, we are ready here since $\mathrm{tr}(S_{2r}) = d$ implies by the arithmetic-geometric mean property that $|S_{2r}| \leq (d/n)^n$.) For $\nu > 0$, we can express $S_{0r}$ in terms of $\mathrm{tr}(S_{2r})$ as

$$S_{0r} = \frac{d+\nu}{\nu} - \frac{1}{\nu} \mathrm{tr}(S_{2r}).$$

Next we consider maximizing the function

$$g \colon \mathrm{SPD}(d) \to \mathbb{R}, \qquad g(X) = \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \mathrm{tr}(X) \right)^{\nu} |X|$$

under the constraint $0 \leq \mathrm{tr}(X) \leq d + \nu$. Note that for $X \in \mathrm{SPD}(d)$ we always have $\mathrm{tr}(X), |X| > 0$. Further, if $\mathrm{tr}(X) = 0$ or $\mathrm{tr}(X) = d + \nu$, we set $g(X) = 0$. Since $g(I) = 1$, the maximum is not attained at the boundary, but inside the open set. The derivative of $g$ with respect to $X$ is given by

$$\nabla g(X) = -\nu \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \mathrm{tr}(X) \right)^{\nu-1} \frac{1}{\nu} I |X| + \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \mathrm{tr}(X) \right)^{\nu} |X| X^{-1}$$

$$= \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \mathrm{tr}(X) \right)^{\nu-1} |X| \left[ \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \mathrm{tr}(X) \right) X^{-1} - I \right].$$

The necessary condition for a critical point $\hat{X}$ of $g$ reads as

$$\left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \mathrm{tr}(\hat{X}) \right) I = \hat{X}.$$

To verify that $\hat{X} = I$ is a maximizer, we have a look at the Hessian $\nabla^2 g$ of $g$ and show that it is negative definite for $\hat{X} = I$. We compute

$$D(\nabla g)(X)[H] = |X| \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \mathrm{tr}(|X|) \right)^{\nu-2} \left[ \frac{\nu-1}{\nu} \mathrm{tr}(H) I \right.$$

$$- \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \mathrm{tr}(|X|) \right) \left( \mathrm{tr}(X^{-1}H) I + \mathrm{tr}(H) X^{-1} \right)$$

$$\left. + \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \mathrm{tr}(|X|) \right)^2 \mathrm{tr}(X^{-1}H) X^{-1} - X^{-1} H X^{-1} \right]$$

and further

$$\langle \nabla^2 g(X)[H], H \rangle = |X| \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \operatorname{tr}(|X|) \right)^{\nu-2} \left[ \frac{\nu-1}{\nu} \operatorname{tr}(H)^2 \right.$$

$$- 2 \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \operatorname{tr}(|X|) \right) \operatorname{tr}(H) \operatorname{tr}(X^{-1}H)$$

$$\left. + \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \operatorname{tr}(|X|) \right)^2 \left( \operatorname{tr}(X^{-1}H)^2 - \operatorname{tr}((X^{-1}H)^2) \right) \right].$$

At the critical point $\hat{X} = I$ we obtain

$$\langle \nabla^2 g(I)[H], H \rangle = \frac{\nu-1}{\nu} \operatorname{tr}(H)^2 - 2\operatorname{tr}(H)^2 + \operatorname{tr}(H)^2 - \operatorname{tr}(H^2)$$

$$= -\frac{1}{\nu} \operatorname{tr}(H)^2 - \operatorname{tr}(H^2) < 0,$$

so that $\hat{X} = I$ is indeed a maximizer. The corresponding maximum of $g$ is given by $g(I) = \left( \frac{d+\nu}{\nu} - \frac{1}{\nu} \operatorname{tr}(I) \right)^\nu |I| = 1$ and therewith finally

$$\Upsilon = (S_{0r}^\nu |S_{2r}|)^{\frac{1}{d+1}} = g(S_{2r})^{\frac{1}{d+1}} \le 1.$$

By (25), we have equality in the theorem if and only if $\Sigma_{r+1} = \Sigma_r$. $\quad\square$

Next, we analyze the difference of two iterates for fixed $\Sigma$.

**Lemma 5.2.** *For fixed $\Sigma \in \mathrm{SPD}(d)$ and $\nu > 0$, let $\{\mu_r\}_{r \in \mathbb{N}}$ be defined by Algorithm 1. Then it holds*

$$L(\mu_{r+1}, \Sigma) - L(\mu_r, \Sigma) \le 0$$

*with equality if and only if $\mu_{r+1} = \mu_r$.*

*Proof.* By concavity of the logarithm and (26) we have

$$L(\mu_{r+1}, \Sigma) - L(\mu_r, \Sigma) = (d+\nu) \sum_{i=1}^n w_i \log \left( \frac{\nu + \delta_{i,r+1}}{\nu + \delta_{i,r}} \right)$$

$$\le (d+\nu) \log \underbrace{\left( \sum_{i=1}^n w_i \frac{\nu + \delta_{i,r+1}}{\nu + \delta_{i,r}} \right)}_{=\Upsilon},$$

so that it suffices to show that $\Upsilon \le 1$. We compute

$$\Upsilon = \sum_{i=1}^n w_i \frac{\nu + (x_i - \mu_r + \mu_r - \mu_{r+1})^{\mathrm{T}} \Sigma^{-1} (x_i - \mu_r + \mu_r - \mu_{r+1})}{\nu + \delta_{i,r}}$$

23

$$= \sum_{i=1}^{n} w_i \frac{\nu + \delta_{i,r}}{\nu + \delta_{i,r}}$$

$$+ \frac{2\langle R^{-1}(x_i - \mu_r), \overbrace{(R^{-1}(\mu_r - \mu_{r+1}))}^{-\frac{S_{1r}}{S_{0r}}}\rangle + (\mu_r - \mu_{r+1})^{\mathrm{T}}\Sigma^{-1}(\mu_r - \mu_{r+1})}{\nu + \delta_{i,r}}$$

$$= 1 - 2\sum_{i=1}^{n} w_i \frac{\left\langle R^{-1}(x_i - \mu_r), \frac{S_{1r}}{S_{0r}}\right\rangle}{\nu + \delta_{i,r}} + \sum_{i=1}^{n} w_i \frac{\left\|\frac{S_{1r}}{S_{0r}}\right\|_2^2}{\nu + \delta_{i,r}}$$

$$= 1 - 2\underbrace{\left\langle \sum_{i=1}^{n} w_i \frac{R^{-1}(x_i - \mu_r)}{\nu + \delta_{i,r}}, \frac{S_{1r}}{S_{0r}}\right\rangle}_{=\frac{1}{d+\nu}S_{1r}} + \left\|\frac{S_{1r}}{S_{0r}}\right\|_2^2 \underbrace{\sum_{i=1}^{n} w_i \frac{1}{\nu + \delta_{i,r}}}_{=\frac{1}{d+\nu}S_{0r}}$$

$$= 1 - \frac{1}{d+\nu}\frac{\|S_{1r}\|_2^2}{S_{0r}} \le 1,$$

with equality if and only if $S_{1r} = 0$, that is, $\mu_{r+1} = \mu_r$ and $\mu_r$ is a critical point of $L(\cdot, \Sigma)$. $\quad\square$

Combining the results of Lemma 5.2 and 5.1 we obtain the following lemma.

**Lemma 5.3.** *For $\nu > 0$, let $\{\mu_r, \Sigma_r\}_{r\in\mathbb{N}}$ be defined by Algorithm 1. Then it holds*

$$L(\mu_{r+1}, \Sigma_{r+1}) - L(\mu_r, \Sigma_r) \le 0$$

*with equality if and only if $(\mu_{r+1}, \Sigma_{r+1}) = (\mu_r, \Sigma_r)$.*

*Proof.* By concavity of the logarithm and (26) we have

$$L(\mu_{r+1}, \Sigma_{r+1}) - L(\mu_r, \Sigma_r) \le (d + \nu)\log\Bigg(\underbrace{\sum_{i=1}^{n} w_i \frac{\nu + \delta_{i,r+1}}{\nu + \delta_{i,r}}\frac{|\Sigma_{r+1}|^{\frac{1}{d+\nu}}}{|\Sigma_r|^{\frac{1}{d+\nu}}}}_{=\Upsilon}\Bigg),$$

and it suffices to show that $\Upsilon \le 1$. As in the proof of Lemma 5.1 we have

$$\frac{|\Sigma_{r+1}|^{\frac{1}{d+\nu}}}{|\Sigma_r|^{\frac{1}{d+\nu}}} = S_{0r}^{-\frac{d}{d+\nu}}|S_{2r}|^{\frac{1}{d+\nu}}.$$

Next, we consider the term

$$\sum_{i=1}^{n} w_i \frac{\nu + \delta_{i,r+1}}{\nu + \delta_{i,r}} = \sum_{i=1}^{n} w_i \frac{\delta_{i,r+1}}{\nu + \delta_{i,r}} + \frac{\nu}{d+\nu}S_{0r}.$$

Combining the computations in the proofs of Lemma 5.2 and Lemma 5.1, we get

$$\sum_{i=1}^{n} w_i \frac{\delta_{i,r+1}}{\nu + \delta_{i,r}}$$

$$= \sum_{i=1}^{n} w_i \frac{\delta_{i,r} + 2\langle R_{r+1}^{-1}(x_i - \mu_r), R_{r+1}^{-1}(\mu_r - \mu_{r+1})\rangle + \delta_{i,r+1}}{\nu + \delta_{i,r}}$$

$$= \sum_{i=1}^{n} w_i \frac{S_{0r} \operatorname{tr}\left(S_{2r}^{-1} R_r^{-1}(x_i - \mu_r)(x_i - \mu_r)^{\mathrm{T}} R_r^{-\mathrm{T}}\right)}{\nu + \delta_{i,r}} - 2\sum_{i=1}^{n} w_i \frac{(x_i - \mu_r)^{\mathrm{T}} R_r^{-\mathrm{T}} S_{0r} S_{2r}^{-1} R_r^{-1} R_r \frac{S_{1r}}{S_{0r}}}{\nu + \delta_{i,r}}$$

$$+ \sum_{i=1}^{n} w_i \frac{\frac{S_{1r}^{\mathrm{T}}}{S_{0r}} R_r^{\mathrm{T}} R_r^{-\mathrm{T}} S_{0r} S_{2r}^{-1} R_r^{-1} R_r \frac{S_{1r}}{S_{0r}}}{\nu + \delta_{i,r}}$$

$$= \frac{d}{d+\nu} S_{0r} - \frac{2}{d+\nu} S_{1r}^{\mathrm{T}} S_{2r}^{-1} S_{1r} + \frac{S_{1r}^{\mathrm{T}} S_{2r}^{-1} S_{1r}}{S_{0r}} \sum_{i=1}^{n} w_i \frac{1}{\nu + \delta_{i,r}}$$

$$= \frac{d}{d+\nu} S_{0r} - \frac{1}{d+\nu} S_{1r}^{\mathrm{T}} S_{2r}^{-1} S_{1r}.$$

Since $S_{2r} \in \mathrm{SPD}(d)$ and consequently also $S_{2r}^{-1} \in \mathrm{SPD}(d)$, we obtain

$$\Upsilon = \left( \left( \frac{d}{d+\nu} + \frac{\nu}{d+\nu} \right) S_{0r} - \frac{1}{d+\nu} \underbrace{S_{1r}^{\mathrm{T}} S_{2r}^{-1} S_{1r}}_{\geq 0} \right) S_{0r}^{-\frac{d}{d+\nu}} |S_{2r}|^{\frac{1}{d+\nu}} \leq (S_{0r}^{\nu} |S_{2r}|)^{\frac{1}{d+\nu}} \leq 1. \quad \square$$

**Theorem 5.4** (Convergence of Algorithm 1).    *(i) Let $\nu \geq 1$ and $x_i \in \mathbb{R}^d$, $i = 1, \ldots, n$ and $w \in \mathring{\Delta}_n$ fulfill Assumption 4.9. Then the sequence $\{(\mu_r, \Sigma_r)\}_{r \in \mathbb{N}}$ generated by Algorithm 1 converges to the minimizer of $L$.*

*(ii) Let $\nu > 0$, $\mu$ fixed and $x_i \in \mathbb{R}^d$, $i = 1, \ldots, n$ and $w \in \mathring{\Delta}_n$ fulfill Assumption 4.1. Then the sequence $\{\Sigma_r\}_{r \in \mathbb{N}}$ generated by Algorithm 1 converges to the minimizer of $L$.*

*(iii) Let $\nu = 0$, $\mu = 0$ and $x_i \in \mathbb{S}^{d-1}$, $i = 1, \ldots, n$ be pairwise different non antipodal points. Then the sequence $\{\Sigma_r\}_{r \in \mathbb{N}}$ generated by Algorithm 1 converges to the minimizer of $L_0$ with trace 1.*

*Proof.* We restrict our attention to (i), the other parts follows the same lines. Let $\{(\mu_r, \Sigma_r)\}_{r \in \mathbb{N}}$ be the sequence of iterates generated by Algorithm 1. Consider the mapping

$$T(\mu, \Sigma) = \left( \mu + \Sigma^{\frac{1}{2}} \frac{S_1(\mu, \Sigma)}{S_0(\mu, \Sigma)}, \Sigma^{\frac{1}{2}} \frac{S_1(\mu, \Sigma)}{S_0(\mu, \Sigma)} \Sigma^{\frac{1}{2}} \right).$$

Then, according to (23), (24) and Theorem 4.11, $(\mu, \Sigma) = T(\mu, \Sigma)$ is a fixed point of $T$ if and only if it is the unique critical point/minimizer of $L$. Consider the case $(\mu_{r+1}, \Sigma_{r+1}) \neq (\mu_r, \Sigma_r)$ for all $r \in \mathbb{N}$. We show that the sequence $\{(\mu_r, \Sigma_r)\}_{r \in \mathbb{N}}$ is bounded: the update of $\mu_r$ is just a convex combination of the samples $x_1, \ldots, x_n$. By construction, $\Sigma_r$, is also a weighted average and the sequence remains bounded since the value $\mu_r \in \operatorname{conv}\{x_1, \ldots, x_n\}$ stays bounded, $(x_i - \mu_r)(x_i - \mu_r)^{\mathrm{T}}$, $i = 1, \ldots, n$ is bounded as well, and consequently also $\Sigma_{r+1}$. By Lemma 5.3, we see that the sequence $L_r := L(\mu_r, \Sigma_r)$ is a strictly decreasing, bounded below sequence such that it converges to some $\hat{L}$. Further, $\{(\mu_r, \Sigma_r)\}_{r \in \mathbb{N}}$ contains a convergent

subsequence $\{(\mu_{r_s}, \Sigma_{r_s})\}_{s \in \mathbb{N}}$, which converges to some $(\hat{\mu}, \hat{\Sigma})$. By the continuity of $L$ and $T$ we obtain

$$
\begin{aligned}
L(\hat{\mu}, \hat{\Sigma}) &= \lim_{s \to \infty} L(\mu_{r_s}, \Sigma_{r_s}) = \lim_{s \to \infty} L_{r_s} = \lim_{s \to \infty} L_{r_s+1} \\
&= \lim_{s \to \infty} L(\mu_{r_s+1}, \Sigma_{r_s+1}) \\
&= \lim_{s \to \infty} L\big(T(\mu_{r_s}, \Sigma_{r_s})\big) = L\big(T(\hat{\mu}, \hat{\Sigma})\big).
\end{aligned}
$$

This implies $(\hat{\mu}, \hat{\Sigma}) = T(\hat{\mu}, \hat{\Sigma})$, so that $(\hat{\mu}, \hat{\Sigma})$ is a fixed point of $T$ and consequently the critical point. Since this point is unique, not only a subsequence, but the whole sequence $\{(\mu_r, \Sigma_r)\}_{r \in \mathbb{N}}$ converges to $(\hat{\mu}, \hat{\Sigma})$, which finishes the proof. $\qquad \square$

## 5.2. Simulation Study

Next we evaluate the numerical performance, in particular the speed of convergence, of the proposed GMMF Algorithm 1 compared to the EM algorithm. Actually, the EM algorithm [29] in our notation reads as Algorithm 1 except for the iteration with respect to $\Sigma$ which is given by

$$
\Sigma_{r+1} = \frac{1}{\nu + d} \sum_{i=1}^{n} w_i \frac{(x_i - \mu_r)(x_i - \mu_r)^{\mathrm{T}}}{\nu + \delta_{i,r}}.
$$

The convergence of the EM algorithm under quite general assumptions has been established in [53]. However, it is well known that the EM algorithm might suffer from slow convergence, which is also what we observed in the following Monte Carlo simulation: we draw $n = 100$ i.i.d. samples of a $T_\nu(\mu, \Sigma)$ distribution for different degrees of freedom $\nu \in \{1, 5, 10, 100\}$ and run Algorithm 1 respective the EM Algorithm to compute the joint ML-estimate $(\hat{\mu}, \hat{\Sigma})$. Both algorithms are initialized with sample mean and sample covariance and we used the relative difference between two iterates $(\mu_r, \Sigma_r)$ and $(\mu_{r+1}, \Sigma_{r+1})$ as stopping criterion, that is

$$
\frac{\sqrt{\|\mu_{r+1} - \mu_r\|_2^2 + \|\Sigma_{r+1} - \Sigma_r\|_F^2}}{\sqrt{\|\mu_r\|_2^2 + \|\Sigma_r\|_F^2}} < 10^{-6}.
$$

This experiment is repeated $N = 10.000$ times and afterward, we calculated the average number of iterations $\overline{\text{iter}}$ and $\overline{\text{iter}}_{\text{EM}}$ needed to reach the tolerance criterion together with their standard deviations. The results are given in Table 1, where we chose $d = 2$, $\mu = 0$ and different values for $\Sigma$. First, we notice that the average number of iterations is in general higher for the EM Algorithm, and further, it does merely not depend on $(\mu, \Sigma)$, but only on the degree of freedom $\nu$. Here, the smaller the value of $\nu$, the larger on the one hand the number of iterations for both algorithms, and on the other hand the larger the gain in speed of Algorithm 1 compared to the EM Algorithm. The fact that only very few iterations are

needed for large $\nu$ can be explained by the fact that for $\nu \to \infty$ the Student-$t$ distribution converges to the normal distribution, so that on the one hand, the estimation becomes in some sense easier, and on the other hand, we can expect the initialization with sample mean and sample covariance matrix, which are the maximum likelihood estimates for the normal distribution, to be already close to the actual parameters to be estimated. Additionally, we examined the speed of convergence for other location parameters $\mu$ as well as fixing $\mu$ and estimating only $\Sigma$. Here, the results are qualitatively and quantitatively similar.

Table 1: Comparison of GMMF Algorithm 1 and the EM Algorithm.

| $\Sigma$ | $\nu$ | $\overline{\text{iter}} \pm \sigma(\overline{\text{iter}})$ | $\overline{\text{iter}}_{\text{EM}} \pm \sigma(\overline{\text{iter}}_{\text{EM}})$ |
|---|---|---|---|
| $\begin{pmatrix} 0.1 & 0 \\ 0 & 0.1 \end{pmatrix}$ | 1 | $20.7582 \pm 1.5430$ | $60.6318 \pm 3.9313$ |
| | 2 | $16.0843 \pm 1.1242$ | $33.5389 \pm 2.0370$ |
| | 5 | $11.166 \pm 0.8121$ | $16.8973 \pm 0.9948$ |
| | 10 | $8.5245 \pm 0.6450$ | $11.1186 \pm 0.6534$ |
| | 100 | $4.1066 \pm 0.3086$ | $4.9072 \pm 0.2915$ |
| $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ | 1 | $20.3536 \pm 1.5899$ | $60.8843 \pm 3.9302$ |
| | 2 | $15.7742 \pm 1.1840$ | $33.6515 \pm 2.0373$ |
| | 5 | $10.9528 \pm 0.8513$ | $16.9305 \pm 0.9957$ |
| | 10 | $8.3487 \pm 0.6646$ | $11.1186 \pm 0.6534$ |
| | 100 | $4.0654 \pm 0.2472$ | $4.9040 \pm 0.2953$ |
| $\begin{pmatrix} 5 & 0 \\ 0 & 5 \end{pmatrix}$ | 1 | $20.2702 \pm 1.6145$ | $60.9139 \pm 3.9326$ |
| | 2 | $15.7136 \pm 1.0099$ | $33.6644 \pm 2.0381$ |
| | 5 | $10.9100 \pm 0.8699$ | $16.9343 \pm 0.9957$ |
| | 10 | $8.3181 \pm 0.6738$ | $11.1191 \pm 0.6540$ |
| | 100 | $4.0627 \pm 0.2424$ | $4.9035 \pm 0.2960$ |
| $\begin{pmatrix} 10 & 0 \\ 0 & 10 \end{pmatrix}$ | 1 | $20.2592 \pm 1.6179$ | $28.0073 \pm 2.0546$ |
| | 2 | $15.7055 \pm 1.2136$ | $33.6662 \pm 2.0386$ |
| | 5 | $10.9050 \pm 0.8725$ | $16.9346 \pm 0.9959$ |
| | 10 | $8.3137 \pm 0.6757$ | $11.1195 \pm 0.6537$ |
| | 100 | $4.0623 \pm 0.2417$ | $4.9036 \pm 0.2958$ |
| $\begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix}$ | 1 | $20.2091 \pm 1.6384$ | $27.9290 \pm 2.1314$ |
| | 2 | $15.6265 \pm 1.2344$ | $33.6569 \pm 2.0452$ |
| | 5 | $10.8407 \pm 0.8841$ | $16.9230 \pm 0.9954$ |
| | 10 | $8.2607 \pm 0.6844$ | $11.1092 \pm 0.6534$ |
| | 100 | $4.0573 \pm 0.2324$ | $4.8908 \pm 0.3126$ |

## 5.3. ML Estimation of Wrapped Cauchy Distribution

By the relation between the wrapped Cauchy and projected normal distribution in Lemma 2.2, we can use our GMMF Algorithm 1 with $\mu_r \equiv \mu = 0$ for the ML estimation of the wrapped Cauchy distribution. By (3) and (4) and (28), the iteration for $\Sigma$ in Algorithm 1 is of the form

$$\Sigma_{r+1} = \begin{pmatrix} \sigma_{11,r+1} & \sigma_{12,r+1} \\ \sigma_{12,r+1} & \sigma_{22,r+1} \end{pmatrix} = \frac{\sum\limits_{i=1}^{n} w_i \frac{x_i x_i^{\mathrm{T}}}{x_i^{\mathrm{T}} \Sigma_r^{-1} x_i}}{\sum\limits_{i=1}^{n} w_i \frac{1}{x_i^{\mathrm{T}} \Sigma_r^{-1} x_i}}$$

$$= \frac{\sum\limits_{i=1}^{n} w_i \frac{1}{\left(1-\xi_{1,r}\cos(2\phi_i)-\xi_{2,r}\sin(2\phi_i)\right)} \begin{pmatrix} \cos^2(\phi_i) & \cos(\phi_i)\sin(\phi_i) \\ \cos(\phi_i)\sin(\phi_i) & \sin^2(\phi_i) \end{pmatrix}}{\sum\limits_{i=1}^{n} w_i \frac{1}{\left(1-\xi_{1,r}\cos(2\phi_i)-\xi_{2,r}\sin(2\phi_i)\right)}}$$

Using relations for trigonometric functions and $\vartheta_i = 2\phi_i$, we obtain

$$\zeta_{1,r+1} = \frac{\sigma_{11,r+1}-\sigma_{22,r+1}}{\sigma_{11,r+1}+\sigma_{22,r+1}} = \frac{\sum\limits_{i=1}^{n} w_i \frac{\cos(\vartheta_i)}{1-\xi_{1,r}\cos(\vartheta_i)-\xi_{2,r}\sin(\vartheta_i)}}{\sum\limits_{i=1}^{n} w_i \frac{1}{1-\xi_{1,r}\cos(\vartheta_i)-\xi_{2,r}\sin(\vartheta_i)}},$$

$$\zeta_{2,r+1} = \frac{2\sigma_{12,r+1}}{\sigma_{11,r+1}+\sigma_{22,r+1}} = \frac{\sum\limits_{i=1}^{n} w_i \frac{\sin(\vartheta_i)}{1-\xi_{1,r}\cos(\vartheta_i)-\xi_{2,r}\sin(\vartheta_i)}}{\sum\limits_{i=1}^{n} w_i \frac{1}{1-\xi_{1,r}\cos(\vartheta_i)-\xi_{2,r}\sin(\vartheta_i)}}.$$

This results in Algorithm 2. From $\zeta_1$ and $\zeta_2$ we obtain the desired estimation of parameters $(a,\gamma)$ of the wrapped Cauchy distribution via (3) and (4).

---

**Algorithm 2** ML estimation for the wrapped Cauchy distribution

---

**Input:** $\vartheta_1,\ldots,\vartheta_n \in [-\pi,\pi)$, $n \geq 3$, $w \in \mathring{\Delta}_n$ fulfilling Assumption 4.9
**Initialization:** $\zeta_{1,0} = \zeta_{2,0} = 0$
**for** $r = 0,\ldots$ **do**

$$\zeta_{1,r+1} = \frac{\sum\limits_{i=1}^{n} w_i \frac{\cos(\vartheta_i)}{1-\zeta_{1,r}\cos(\vartheta_i)-\zeta_{2,r}\sin(\vartheta_i)}}{\sum\limits_{i=1}^{n} w_i \frac{1}{1-\zeta_{1,r}\cos(\vartheta_i)-\zeta_{2,r}\sin(\vartheta_i)}}$$

$$\zeta_{2,r+1} = \frac{\sum\limits_{i=1}^{n} w_i \frac{\sin(\vartheta_i)}{1-\zeta_{1,r}\cos(\vartheta_i)-\zeta_{2,r}\sin(\vartheta_i)}}{\sum\limits_{i=1}^{n} w_i \frac{1}{1-\zeta_{1,r}\cos(\vartheta_i)-\zeta_{2,r}\sin(\vartheta_i)}}$$

---

# 6. Applications in Image Analysis

In this section, we describe how our GMMF can be used to denoise images corrupted by different kinds of (additive) noise, in particular Cauchy, Gaussian and wrapped Cauchy noise.

## 6.1. Nonlocal Denoising Approach

Let $f \colon \mathcal{G} \to \mathbb{R}$ be a noisy image, where $\mathcal{G} = \{1, \ldots, n_1\} \times \{1, \ldots, n_2\}$ denotes the image domain. Here and in all subsequent cases we extend the image by mirroring at the boundary. We assume that each pixel $f_i$, $i = (i_1, i_2) \in \mathcal{G}$ is affected by noise in an independent and identical way. Based on Theorem 2.1, this can be modeled as

$$f_i = u_i + \sigma \frac{\eta}{\sqrt{y}}, \qquad i \in \mathcal{G},$$

where $u$ is the noise-free image we wish to reconstruct, $\eta$ is a realization of $Z \sim \mathcal{N}(0, 1)$, and $y$ is a realization of $Y \sim \Gamma\left(\frac{\nu}{2}, \frac{\nu}{2}\right)$, where $Z$ and $Y$ are independent. The given parameter $\nu \geq 1$ determines the amount of outliers and $\sigma > 0$ their strength. This results in independent realizations $f_i$ of $T_\nu(u_i, \sigma)$ distributed random variables.

**Selection of Samples.** The estimation of the noise-free image requires to select for each $i \in \mathcal{G}$ a set of indices of samples $\mathcal{S}(i)$ that are interpreted as i.i.d. realizations of $T_\nu(u_i, \sigma)$. We focus here on a nonlocal approach, which is based on an image self-similarity assumption stating that small patches of an image can be found several times in the image. Then, the set $\mathcal{S}(i)$ constitutes of the indices of the centers of $K$ patches that are most similar to the patch centered at $i \in \mathcal{G}$. This requires the selection of the patch size and an appropriate similarity measure, which need to be adapted to the noise statistic and the noise level and is described in the next paragraph. In order to avoid a computational overload one typically restricts the search zone for similar patches to a $w \times w$ search window around $i \in \mathcal{G}$. Having defined for each $i \in \mathcal{G}$ a set of indices of samples $\mathcal{S}(i)$, the noise-free image can be estimated as

$$\hat{u}_i \in \underset{u_i}{\operatorname{argmin}} \left\{ L\left(u_i, \sigma | \{f_j\}_{j \in \mathcal{S}(i)}\right) \right\}, \quad i \in \mathcal{G}.$$

However, by using such a pixelwise treatment we implicitly assume that the pixels of an image are independent of each other, which is in practice a rather unrealistic assumption; in fact, in natural images they are locally usually highly correlated. Taking the local dependence structure into account may improve the results of image restoration methods, which motivates to take whole patches (and not only their centers) and estimate their parameters using Algorithm 1 as

$$(\hat{\mu}_i, \hat{\Sigma}_i) \in \underset{\mu_i, \Sigma_i}{\operatorname{argmin}} \left\{ L\left(\mu_i, \Sigma_i | \{p_j\}_{j \in \mathcal{S}(i)}\right) \right\}, \quad i \in \mathcal{G},$$

where $p_j$ denotes a patch centered at $j \in \mathcal{S}(i)$. If $\nu > 2$, such that the covariance matrix of the Student-$t$ distribution exists, we take the estimated correlation into account and restore

29

a patch $\hat{p}_i$ as

$$\hat{p}_i = \hat{\mu}_i + \left( \hat{\Sigma}_i - \frac{\nu}{\nu - 2} \sigma^2 I \right) \hat{\Sigma}_i^{-1} (p_i - \hat{\mu}_i),$$

which is in this case the best linear unbiased estimator (BLUE) of $p_i$, see [25]. If $\nu \leq 2$, we set $\hat{p}_i := \hat{\mu}_i$. Proceeding as above gives multiple estimates for each single image pixel that are averaged in the end to obtain the final image.

**Patch Similarity.** The selection of similar patches constitutes a fundamental step in our nonlocal denoising approach. At this point, the question arises how to compare noisy patches and numerical examples show that an adaptation of the similarity measure to the noise distribution is essential for a robust similarity evaluation. In [9], the authors formulated the similarity between patches as a statistical hypothesis testing problem and proposed among other criteria a similarity measure based on a generalized likelihood test, which we use in the following. Further details on this approach can also be found in [26]. We briefly summarize the main ideas.

Modeling noisy images in a stochastic way allows to formulate the question whether two patches $p$ and $q$ are similar as a hypothesis test. Two noisy patches $p, q$ are considered to be similar if they are realizations of independent random variables $X \sim p_{\vartheta_1}$ and $Y \sim p_{\vartheta_2}$ that follow the same parametric distribution $p_\vartheta$, $\vartheta \in \Theta$ with a common parameter $\vartheta$ (corresponding to the underlying noise-free patch), i.e. $\vartheta_1 = \vartheta_2 \equiv \vartheta$. Therewith, the evaluation of the similarity between noisy patches can be formulated as the following hypothesis test:

$$\mathcal{H}_0 \colon \vartheta_1 = \vartheta_2 \qquad \text{vs.} \qquad \mathcal{H}_1 \colon \vartheta_1 \neq \vartheta_2. \tag{27}$$

In this context, a similarity measure $S$ maps a pair of noisy patches $p, q$ to a real value $c \in \mathbb{R}$. The larger this value $c$ is, the more the patches are considered to be similar.

In the following, we describe how a similarity measure $S$ can be obtained based on a suitable test statistic for the hypothesis testing problem. In general, according to the Neyman-Pearson Theorem, see, e.g. [6], the optimal test statistic (i.e. the one that maximizes the power for any given size $\alpha$) for single-valued hypotheses of the form

$$\mathcal{H}_0 \colon \vartheta = \vartheta_0 \qquad \text{vs.} \qquad \mathcal{H}_1 \colon \vartheta = \vartheta_1$$

is given by a likelihood ratio test. Note that single-valued testing problems correspond to a disjoint partition of the parameter space of the form $\Theta = \Theta_0 \dot{\cup} \Theta_1$, where $\Theta_i = \{\vartheta_i\}$, $i = 0, 1$. Despite being a very strong theoretical result, the practical relevance of the Neyman-Pearson Theorem is limited due to the fact that $\Theta_0$ and $\Theta_1$ are in most applications not single-valued. Instead, the testing problem is a so called *composite testing* problem, meaning that $\Theta_0$ and/or $\Theta_1$ contain more than one element. It can be shown that for composite testing

problems there does not exist a uniformly most powerful test. Now, the idea to generalize the Neyman-Pearson test to composite testing problems is to obtain first two candidates (or representatives) $\hat{\vartheta}_i$ of $\Theta_i$, and $i = 0, 1$ respectively, e.g. by maximum-likelihood estimation, and then to perform a Neyman-Pearson test using the computed candidates $\hat{\vartheta}_0$ and $\hat{\vartheta}_1$ in the definition of the test statistic. In case that an ML-estimation is used to determine $\hat{\vartheta}_0$ and $\hat{\vartheta}_1$, the resulting test is called *Likelihood Ratio Test* (LR test). Although there are in general no theoretical guarantees concerning the power of LR tests, they usually perform very well in practice if the sample size used to estimate $\hat{\vartheta}_0$ and $\hat{\vartheta}_1$ is large enough. This is due to the fact that ML estimators are asymptotically efficient. Several classical tests, e.g. one and two-sided $t$-tests, are either direct LR tests or equivalent to them.

In the sequel, we show how the above framework can be applied to our testing problem (27). First, let $x_1$ and $y_1$ be two single pixels for which we want to test whether they are realizations of the same distribution with unknown common parameter. The LR statistic reads as

$$\lambda(x_1, y_1) = \frac{\sup\limits_{\vartheta \in \mathcal{H}_0} \left\{ \mathcal{L}(\vartheta | x_1, y_1) \right\}}{\sup\limits_{\vartheta} \left\{ \mathcal{L}(\vartheta | x_1, y_1) \right\}} = \frac{\sup\limits_{\vartheta} \left\{ \mathcal{L}(\vartheta | x_1) \mathcal{L}(\vartheta | y_1) \right\}}{\sup\limits_{\vartheta} \left\{ \mathcal{L}(\vartheta | x_1) \right\} \sup\limits_{\vartheta} \left\{ \mathcal{L}(\vartheta | y_1) \right\}},$$

where in our situation $\mathcal{L}(\vartheta | x_1, y_1)$ denotes the likelihood function with respect to $\mu$ while $\nu$ and $\sigma$ are assumed to be known, and the notation $\vartheta \in \mathcal{H}_0$ means that the supremum is taken over those parameters $\vartheta$ fulfilling $\mathcal{H}_0$. We use this statistic as similarity measure, i.e.,

$$S(x_1, y_1) := \lambda(x_1, y_1).$$

More generally, since we assume the noise to affect each pixel in an independent and identical way, the similarity of two patches $p = (x_1, \ldots, x_t)$ and $q = (y_1, \ldots, y_t)$ is obtained as the product of the similarity of its pixels

$$S(p, q) := \prod_{i=1}^{t} S(x_i, y_i).$$

In case of the Student-$t$ distribution with $\nu > 0$ degrees of freedom, the similarity measure between two patches $p = (p_1, \ldots, p_t)$ and $q = (q_1, \ldots, q_t)$ can be computed as

$$S(p, q) = \prod_{i=1}^{t} \left( 1 + \frac{1}{\nu} \left( \frac{p_i - q_i}{2\sigma} \right)^2 \right)^{-(\nu+1)}.$$

In practice, we take a scaled logarithm of $S$ in order to avoid numerical instabilities, resulting
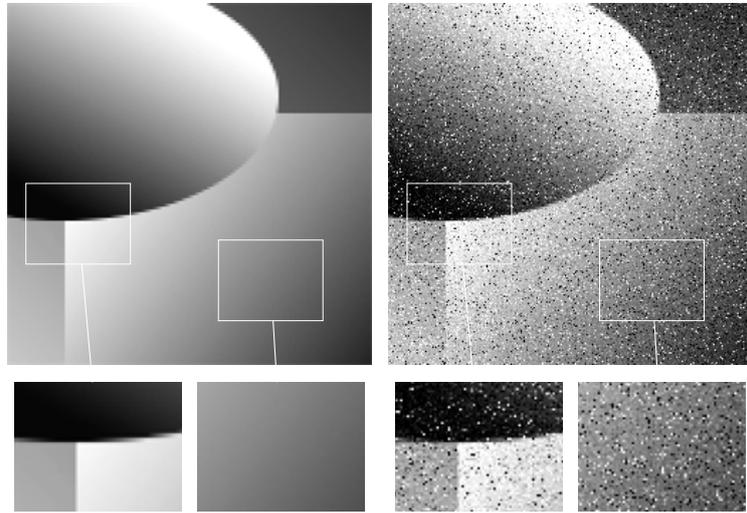
in the distance measure

$$d(p, q) = \sum_{i=1}^{t} \log \left( \nu + \left( \frac{p_i - q_i}{2\sigma} \right)^2 \right).$$
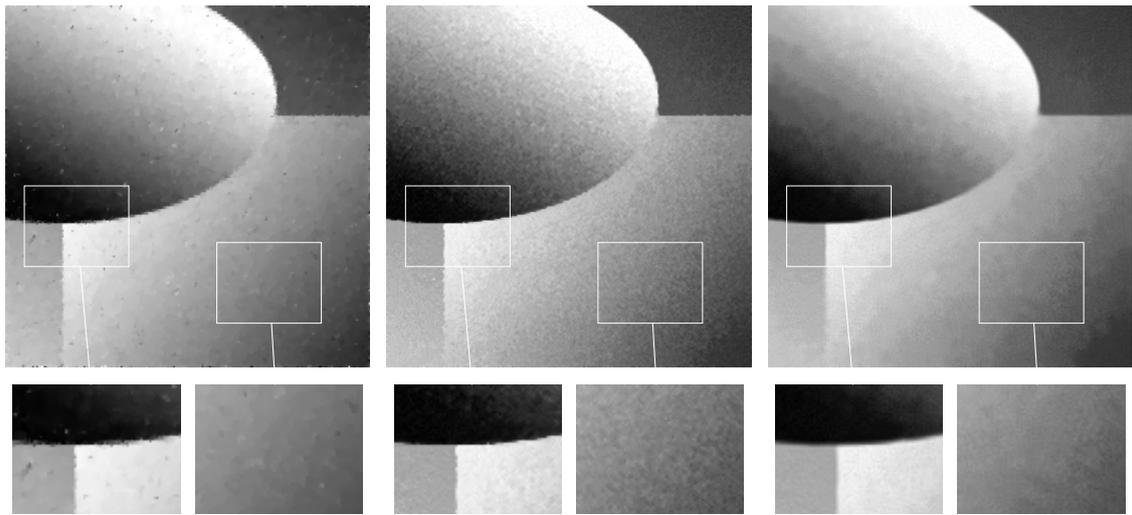
## 6.2. Numerical Examples

**Cauchy Noise**   As mentioned in the introduction, the initial motivation for this work was the consideration of Cauchy noise in [26, 35] and thus we tested our approach on images corrupted by additive Cauchy noise ($\nu = 1$) with noise level $\sigma = 10$. Since the noise level is very high, we chose $n = 50$ patches of size $5 \times 5$ for the denoising. It turns out that the differences in terms of PSNR or SSIM compared to the current state of the art method [26] are small and nearly not visible in images with much textured regions. However, the improvement is large in images with many constant or smoothly varying areas. This becomes in particular apparent in case of the test image given in Figure 2. Here, the top row displays the original image (left) together with its noisy version (right). The bottom row shows from left to right the results obtained using the variational method presented in [35], the pixelwise [26] and the patchwise nonlocal myriad filter. While in case of the variational method some of the outliers remain, the result of [26] is rather grainy, which is much improved by our new approach. This is also reflected in the corresponding PSNR and SSIM values stated in the captions of the figure.

As to be expected, the larger the patch size, the smoother is the resulting image. While this is desirable in constant or smoothly varying regions, fine structures and details are lost when the patch size is too large. Thus, a natural idea is to adapt the filter to the local structure of the image, that is, using a multivariate myriad filter in constant areas and a pixelwise myriad filter in regions with many details. We tried the following naive approach: First, we computed the one-dimensional as well as the multivariate myriad filtered images. Then, for each pixel of the image we decide whether it belongs to a homogeneous area or not and restore the pixel accordingly. In order to detect constant regions we used the variance of the pixels in similar patches, which we expect to be low in regions without much details. Results of this approach are shown in Figures 3 and 4, the corresponding PSNR and SSIM values are given in Table 2, where we used $5 \times 5$ patches a sample size of $n = 40$ for all images. In all examples one nicely sees that with the adaptive approach both fine details as well as constant areas are very well reconstructed, which is not the case when applying only the one-dimensional or only the multivariate myriad filter.

**Gaussian Noise**   For $\nu \to \infty$, the Student-$t$ distribution converges to the normal distribution. Thus, for large $\nu$ we expect the nonlocal GMMF to be able to denoise images corrupted by Gaussian noise as well. This is illustrated in Figure 5, which shows the *barbara* image

(a) PSNR: $+\infty$, SSIM: 1     (b) PSNR: 16.3270, SSIM:0.0870

(c) PSNR: 33.5452, SSIM: 0.8651 (d) PSNR: 33.9144, SSIM:0.8410 (e) PSNR: 37.7263, SSIM: 0.9398

Figure 2: Denoising of the test image (top left) corrupted with additive Cauchy noise ($\nu = 1$, $\sigma = 10$) (top right) using the methods proposed in [35] (bottom left), [26] (bottom middle) and our nonlocal GMMF (bottom right).

Figure 3: Denoising of images (first column) corrupted with additive Cauchy noise ($\nu = 1$, $\sigma = 10$) using a one-dimensional myriad filter [26] (second column), the multivariate myriad filter (third column) and a combination of both (fourth column).

Figure 4: Denoising of images (first column) corrupted with additive Cauchy noise ($\nu = 1$, $\sigma = 10$) using a one-dimensional myriad filter [26] (second column), the multivariate myriad filter (third column) and a combination of both (fourth column).

| Image | PSNR | | | SSIM | | |
|---|---|---|---|---|---|---|
| | 1d | multivariate | adaptive | 1d | multivariate | adaptive |
| | | | $\sigma = 10$ | | | |
| cameraman | 26.2721 | 25.5515 | 27.3726 | 0.7759 | 0.7376 | 0.7650 |
| boat | 26.2230 | 25.8486 | 26.6643 | 0.7349 | 0.7209 | 0.7418 |
| house | 25.1464 | 24.6241 | 25.8974 | 0.7467 | 0.7203 | 0.7689 |
| parrot | 26.4647 | 26.1817 | 27.4714 | 0.7842 | 0.7727 | 0.7858 |
| plane | 25.8675 | 25.6023 | 27.1730 | 0.7694 | 0.7501 | 0.7799 |
| leopard | 24.7906 | 24.1551 | 26.2640 | 0.7535 | 0.7348 | 0.7692 |

Table 2: Comparison of PSNR and SSIM values of the one-dimensional, the multivariate and an adaptive myriad filter.

(top left) corrupted by additive white Gaussian noise of standard deviation $\sigma = 10$ (top right) together with the denoising result using the minimum means squared error estimator proposed in the state-of-the-art nonlocal Bayes algorithm [27] (bottom left) and our nonlocal GMMF with $\nu = 1000$ (bottom right), this time applied with the BLUE estimator instead of the mean. Again, we chose $n = 40$ samples, but since the noise is not very strong we used $3 \times 3$ patches in this example. Both methods reconstruct the image very well and there are no visible differences. Note that we show the results obtained after one iteration of the respective algorithms in order to see the influence of the different estimators and no other effects. The complete denoising procedure proposed in [27] uses this first-step denoised image as an oracle image for patch selection in a second iteration and applies further fine-tuning steps such as a special treatment of homogeneous areas and patch aggregation to obtain the final image.

**Remark 6.1** (Robustness of Parameters). *Our algorithm depends on several parameters, namely the patch width $s$, the size of the search window $w$, and the number of similar patches $n$. An extensive grid search revealed that the results are rather robust towards small changes in the parameters, and patch sizes between $3 \times 3$ and $5 \times 5$ as well as $n = 40$ to $n = 50$ samples yield results that are visually indistinguishable and have similar PSNR values. The size of the search window needs to be adapted to the patch size and the number of samples; it has to be large enough to guarantee that enough similar patches can be found. As a rule of thumb, the stronger the noise, the larger we choose the patch size.*

**Wrapped Cauchy Noise.** Next, we apply Algorithm 2 to denoise $\mathbb{S}^1$-valued images $f \colon \mathcal{G} \to \mathbb{S}^1$ corrupted by wrapped Cauchy noise,

$$f_i = (u_i + \gamma \eta) \operatorname{mod} 2\pi, \qquad \eta \sim C(0, \gamma), \gamma > 0, \quad i \in \mathcal{G},$$

(a) PSNR: $+\infty$, SSIM: 1  (b) PSNR: 28.1241, SSIM: 0.7140

(c) PSNR: 31.2445, SSIM: 0.7873  (d) PSNR: 31.3024, SSIM: 0.7898

Figure 5: Denoising of the test image (top left) corrupted with additive Gaussian noise ($\nu = \infty$, $\sigma = 10$) (top right) using the nonlocal MMSE algorithm [27] (bottom left) and our nonlocal GMMF (bottom right).

where we chose $\gamma = 0.1$ corresponding to a moderate noise level, which yields $\rho = \mathrm{e}^{-\gamma} \approx 0.9048$. The original image as well as the noisy image are given in the top row of Figure 6. The similarity measure to find similar patches is in this case derived from the density of the wrapped Cauchy distribution and it is given by

$$S(p,q) = \frac{(1-\rho)^4}{\left(1 + \rho^2 - 2\rho \cos\left(\frac{p_i - q_i}{2}\right)\right)^2},$$

which leads to the distance

$$d(p,q) = \sum_{i=1}^{t} \log\left(1 + \rho^2 - 2\rho \cos\left(\frac{p_i - q_i}{2}\right)\right).$$

Here, we chose $n = 50$ patches of size $5 \times 5$ as samples, but this time extracted only their centers, estimated the parameters and restored the image pixelwise. We compare our approach with the variational method using an $L_2$ data term and a first and second order TV-regularizer in [3] and the nonlocal denoising algorithm based on second order statistic [25] (NL-MMSE). The results together with the mean-squared reconstruction error are given in the bottom row of Figure 6. Both the variational as well as the second order statistical method cannot cope with the impulsiveness of the wrapped Cauchy noise such that several wrong pixels remain, which is in particular visible in the background. Furthermore, the edges of the color squares and the transitions in the ellipse and in the circle are rather fringy. On the contrary, our method restores the image very well, if at all a slight grain can be observed in the color squares which is due to the pixelwise denoising that does not regard neighboring pixels appropriately.

**Denoising of InSAR Data**    Finally, we apply our denoising approach to a real-world data set given in [41][1]. It consists of an interferometric synthetic aperture radar (InSAR) image recorded in 1991 by the ERS-1 satellite capturing topographical information from the Mount Vesuvius. InSAR is a radar technique used in geodesy and remote sensing. Based on two or more synthetic aperture radar (SAR) images, maps of surface deformation or digital elevation are generated, using differences in the phase of the waves returning to an aircraft or satellite. The technique can potentially measure millimeter-scale changes in deformation over spans of days to years. It has applications in the geophysical monitoring of natural hazards, for example earthquakes, volcanoes and landslides, and in structural engineering, in particular monitoring of subsidence and structural stability. InSAR produces phase-valued images, i.e. in each pixel the measurement lies on the circle $\mathbb{S}^1$. The recorded image is shown
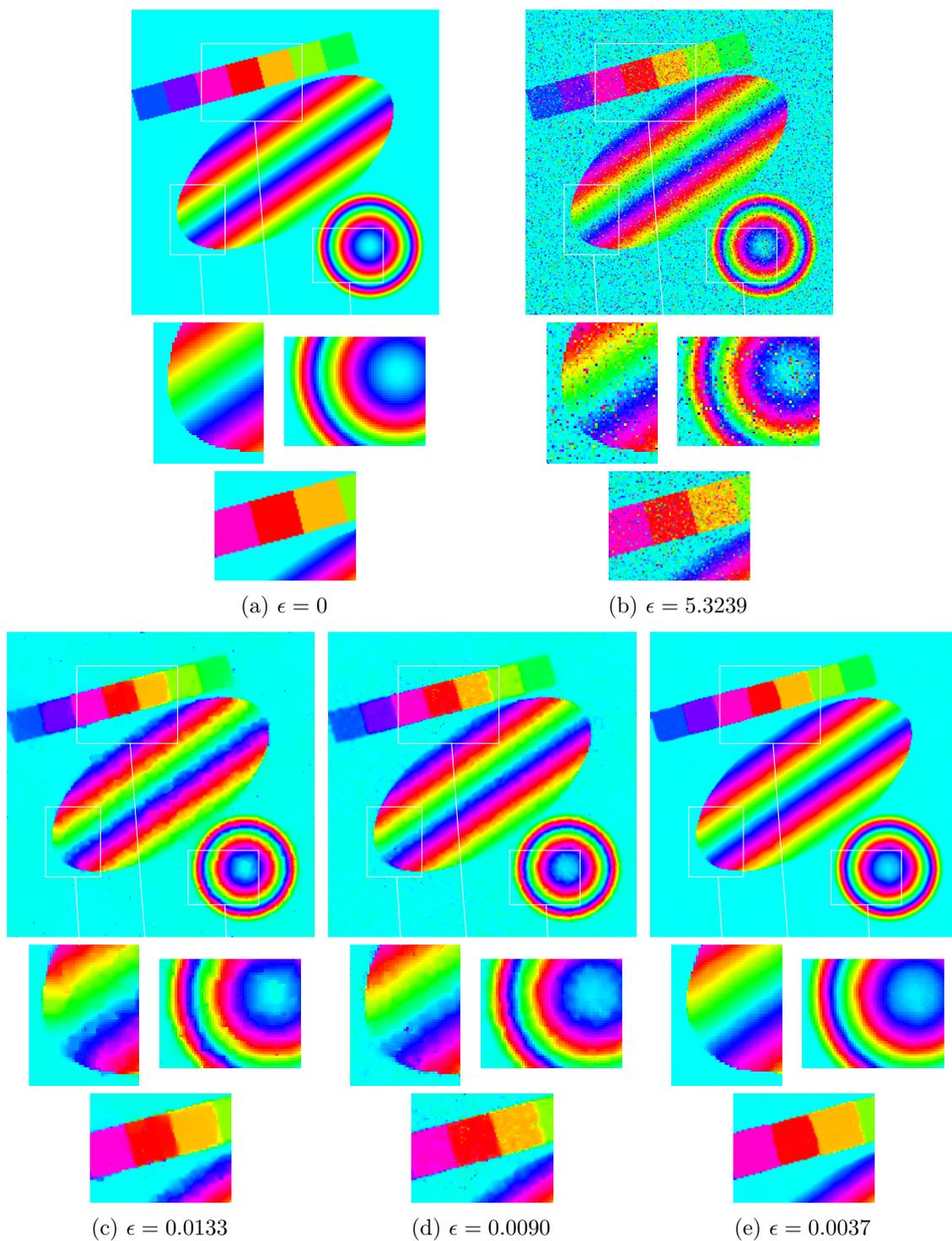
---

[1]The data is available online at https://earth.esa.int/workshops/ers97/program-details/speeches/rocca-et-al/

(a) $\epsilon = 0$        (b) $\epsilon = 5.3239$

(c) $\epsilon = 0.0133$    (d) $\epsilon = 0.0090$    (e) $\epsilon = 0.0037$

Figure 6: Denoising of an $\mathbb{S}^1$-valued image (top left) corrupted with additive wrapped Cauchy noise ($a = 0$, $\rho = 0.1$) (top right) using the variational method [3] (bottom left), the nonlocal MMSE method [25] (bottom middle) and our nonlocal GMMF (bottom right).

39

in Figure 7 top left. While the contour lines of the Vesuvius are clearly visible, the image suffers in particular in the middle and at the boundaries from strong noise. The top right image in Figure 7 depicts the denoising result using the variational approach in [3], while in the bottom row we show the results obtained with the nonlocal MMSE approach [25] (left) and with our method (right), where we used the parameters $\gamma = 0.1$, $n = 50$ and $5 \times 5$ patches. While both methods restore the outer contour lines very well, our approach yields a sharper image and preserves much more details, for instance the fine contour lines in the middle of the image.

In a second experiment, we examine the influence of the denoising on the reconstruction of the absolute phase. In order to do so, we use the phase unwrapping method PUMA proposed in [4], which recasts the problem as a max flow-min cut problem and is among the current state-of-the-art algorithms for phase unwrapping. We used the code provided by the authors of [4] with the default parameters choices. The results of the corresponding images of Figure 7 are given in Figure 8. While the absolute phase of the original image is strongly affected by the noise, the results based on the TV-regularized image are in some parts oversmoothed, so that the reconstruction does not provide any details in those areas. On the other hand, the reconstructions based on the NL-MMSE denoised image and based on our method provide much better results, where our method leads to a slightly finer resolution.

# 7. Conclusion

We introduced a generalized multivariate myriad filter based on (weighted) ML estimation of the multivariate Student-$t$ distribution and illustrated its usage in a nonlocal robust denoising approach. Furthermore, we showed how a special case of our algorithm can be related to projected normal distributions on the sphere $\mathbb{S}^{d-1}$ and for $d = 2$ further to the wrapped Cauchy distribution on the circle $\mathbb{S}^1$, which gives rise to robust denoising strategies for $\mathbb{S}^1$-valued images.

There are different directions for future work: First, we would like to extend our analysis to the case that additionally the degrees of freedom parameter $\nu$ is unknown and to SMM. Although an EM algorithm has already been derived for this case in [23], to the best of our knowledge there do not exist results concerning existence and/or uniqueness of the joint ML estimator.

Concerning our denoising approach, fine tuning steps as discussed in [28] such as aggregations of patches [42], the use of an oracle image or a variable patch size to better cope with textured and homogeneous image regions may improve the denoising results. Further, in all our examples we used uniform weights, but weights based for instance on spatial distance or similarity would make sense as well. Another question is how to incorporate linear operators (blur, missing pixels) into the image restoration. To this end, SMM could be used as priors

(a) Original InSAR image.

(b) Variational method [3].
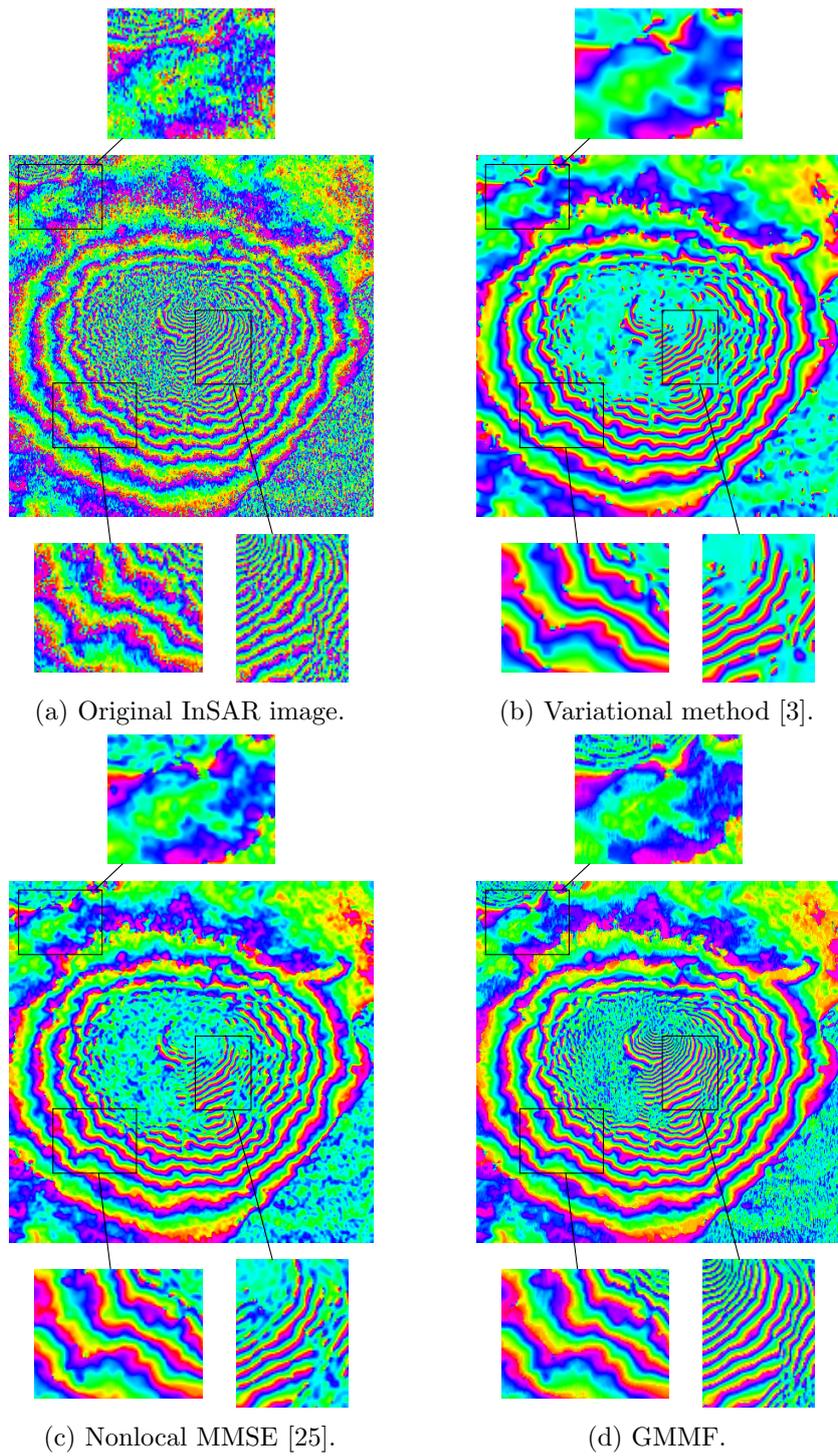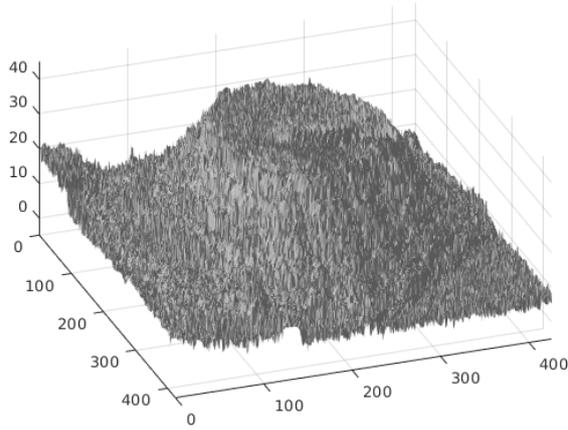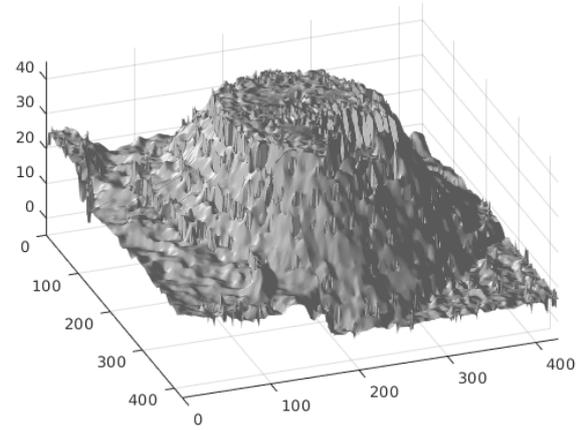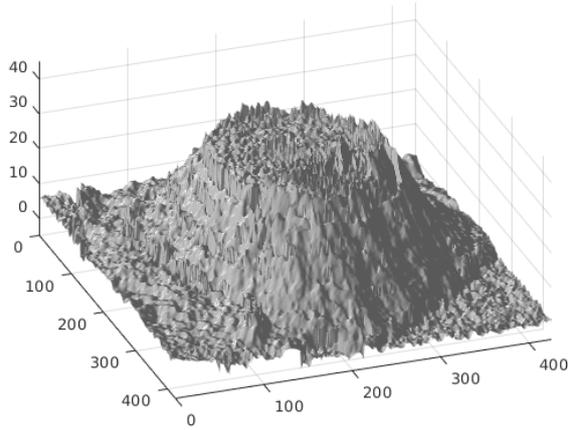
(c) Nonlocal MMSE [25].

(d) GMMF.

Figure 7: Denoising of an InSAR image of the Mount Vesuvius (top left) using the variational method [3] (top right), the nonlocal MMSE approach [25] and our nonlocal GMMF (bottom right).
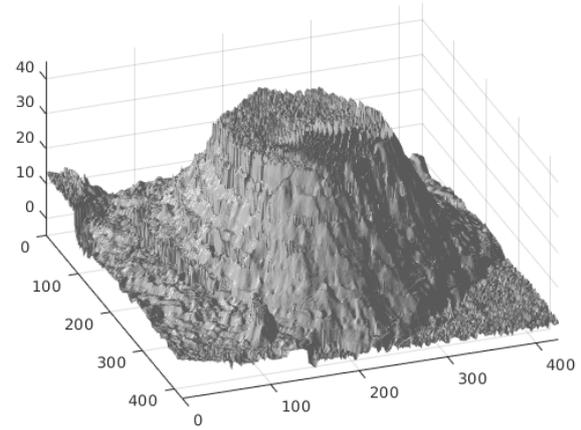
(a) Original InSAR image.

(b) Variational method [3].

(c) Nonlocal MMSE [25].

(d) GMMF.

Figure 8: Reconstruction of the absolute phase images from Figure 7 using the reconstruction algorithm PUMA [4].

within variational models.

## Acknowledgments

# A. Appendix

Proof of Lemma 2.2. (i) By definition of $\Theta$ in terms of $\Phi$ we have for the corresponding probability density functions

$$f_\Theta(\vartheta) = f_{2\Phi}(\vartheta) + f_{2\Phi}\left(\vartheta - \operatorname{sgn}(\vartheta)2\pi\right) = \frac{1}{2}\left[f_\Phi\left(\frac{\vartheta}{2}\right) + f_\Phi\left(\frac{\vartheta}{2} - \operatorname{sgn}(\vartheta)\pi\right)\right]$$

$$= \frac{1}{2}\left[f_\Phi\left(\frac{\vartheta}{2}\right) + f_\Phi\left(\left(\frac{\vartheta}{2} + \pi\right)\bmod 2\pi\right)\right] = f_\Phi\left(\frac{\vartheta}{2}\right).$$

Therefore we have to show that $f_\Phi\left(\frac{\vartheta}{2}\right) = g_w(\vartheta|a,\gamma)$. We parametrize $x = \begin{pmatrix}\cos(\phi)\\\sin(\phi)\end{pmatrix}$ and using relations of trigonometric functions, we obtain

$$x^\mathrm{T}\Sigma^{-1}x = \frac{1}{\sigma_{11}\sigma_{22} - \sigma_{12}^2}\left(\sigma_{22}x_1^2 - 2\sigma_{12}x_1x_2 + \sigma_{11}x_2^2\right)$$

$$= \frac{1}{\sigma_{11}\sigma_{22} - \sigma_{12}^2}\left(\sigma_{22}\cos^2(\phi) + \sigma_{11}\sin^2(\phi) - 2\sigma_{12}\cos(\phi)\sin(\phi)\right)$$

$$= \frac{1}{\sigma_{11}\sigma_{22} - \sigma_{12}^2}\left(\frac{1}{2}(\sigma_{11} + \sigma_{22}) - \frac{1}{2}(\sigma_{11} - \sigma_{22})\cos(2\phi) - \sigma_{12}\sin(2\phi)\right)$$

$$= \frac{\frac{1}{2}(\sigma_{11} + \sigma_{22})}{\sigma_{11}\sigma_{22} - \sigma_{12}^2}\left(1 - \frac{\sigma_{11} - \sigma_{22}}{\sigma_{11} + \sigma_{22}}\cos(2\phi) - \frac{2\sigma_{12}}{\sigma_{11} + \sigma_{22}}\sin(2\phi)\right)$$

$$= \frac{\frac{1}{2}\operatorname{tr}(\Sigma)}{|\Sigma|}\left(1 - \frac{\sigma_{11} - \sigma_{22}}{\sigma_{11} + \sigma_{22}}\cos(2\phi) - \frac{2\sigma_{12}}{\sigma_{11} + \sigma_{22}}\sin(2\phi)\right).$$

Setting $\zeta_1 = \frac{\sigma_{11} - \sigma_{22}}{\sigma_{11} + \sigma_{22}}$ and $\zeta_2 = \frac{2\sigma_{12}}{\sigma_{11} + \sigma_{22}}$ such that

$$\sqrt{1 - \zeta_1^2 - \zeta_2^2} = \frac{2\sqrt{\sigma_{11}\sigma_{22} - \sigma_{12}^2}}{\sigma_{11} + \sigma_{22}} = \frac{\sqrt{|\Sigma|}}{\frac{1}{2}\operatorname{tr}(\Sigma)}$$

we get

$$x^\mathrm{T}\Sigma^{-1}x = \frac{1 - \zeta_1\cos(2\phi) - \zeta_2\sin(2\phi)}{\sqrt{|\Sigma|}\sqrt{1 - \zeta_1^2 - \zeta_2^2}} \tag{28}$$

and

$$f_\Phi\left(\frac{\vartheta}{2}\right) = \frac{1}{2\pi}\frac{\sqrt{1 - \zeta_1^2 - \zeta_2^2}}{1 - \zeta_1\cos(\vartheta) - \xi_2\sin(\vartheta)}.$$

This has exactly the form (2) of $g_w$ if we identify

$$\xi_1 = \frac{2\rho}{1+\rho^2}\cos(a) = \zeta_1 = \frac{\sigma_{11}-\sigma_{22}}{\sigma_{11}+\sigma_{22}}, \quad \xi_2 = \frac{2\rho}{1+\rho^2}\sin(a) = \zeta_2 = \frac{2\sigma_{12}}{\sigma_{11}+\sigma_{22}}.$$

Squaring and adding these equations leads by observing that $\rho < 1$ to (3) and by devision of the equations to (4).

(ii) We first compute the density of $\frac{\Theta}{2} + \pi\Xi$ as

$$f_{\frac{\Theta}{2}+\pi\Xi}(\vartheta) = \frac{1}{4}f_{\frac{\Theta}{2}}(\vartheta+\pi)\mathbb{1}_{\left[-\frac{3\pi}{2},-\frac{\pi}{2}\right)}(\vartheta) + \frac{1}{2}f_{\frac{\Theta}{2}}(\vartheta)\mathbb{1}_{\left[-\frac{\pi}{2},\frac{\pi}{2}\right]}(\vartheta) + \frac{1}{4}f_{\frac{\Theta}{2}}(\vartheta-\pi)\mathbb{1}_{\left[\frac{\pi}{2},\frac{3\pi}{2}\right)}(\vartheta).$$

Therewith, the density of $\left(\frac{\Theta}{2}+\pi\Xi\right)_{2\pi} := \left(\frac{\Theta}{2}+\pi\Xi\right)\bmod 2\pi$ is given by

$$f_{\left(\frac{\Theta}{2}+\pi\Xi\right)_{2\pi}}(\vartheta) = f_{\frac{\Theta}{2}+\pi\Xi}(\vartheta)\mathbb{1}_{[-\pi,\pi)}(\vartheta) + \underbrace{f_{\frac{\Theta}{2}+\pi\Xi}(\vartheta+2\pi)\mathbb{1}_{\left[-\pi,-\frac{\pi}{2}\right)}(\vartheta)}_{(*)} + \underbrace{f_{\frac{\Theta}{2}+\pi\Xi}(\vartheta-2\pi)\mathbb{1}_{\left[\frac{\pi}{2},\pi\right)}(\vartheta)}_{(**)}.$$

We calculate

$$(*) = \tfrac{1}{4}f_{\frac{\Theta}{2}}(\vartheta+2\pi+\pi)\mathbb{1}_{\left[-\frac{3\pi}{2},-\frac{\pi}{2}\right)}(\vartheta+2\pi)\mathbb{1}_{\left[-\pi,-\frac{\pi}{2}\right)}(\vartheta) + \tfrac{1}{2}f_{\frac{\Theta}{2}}(\vartheta+2\pi)\mathbb{1}_{\left[-\frac{\pi}{2},\frac{\pi}{2}\right)}(\vartheta+2\pi)\mathbb{1}_{\left[-\pi,-\frac{\pi}{2}\right)}(\vartheta)$$

$$+ \tfrac{1}{4}f_{\frac{\Theta}{2}}(\vartheta+2\pi-\pi)\mathbb{1}_{\left[\frac{\pi}{2},\frac{3\pi}{2}\right)}(\vartheta+2\pi)\mathbb{1}_{\left[-\pi,-\frac{\pi}{2}\right)}(\vartheta) = \tfrac{1}{4}f_{\frac{\Theta}{2}}(\vartheta+\pi)\mathbb{1}_{\left[-\pi,-\frac{\pi}{2}\right)}(\vartheta)$$

$$(**) = \tfrac{1}{4}f_{\frac{\Theta}{2}}(\vartheta-2\pi+\pi)\mathbb{1}_{\left[-\frac{3\pi}{2},-\frac{\pi}{2}\right)}(\vartheta-2\pi)\mathbb{1}_{\left[\frac{\pi}{2},\pi\right)}(\vartheta) + \tfrac{1}{2}f_{\frac{\Theta}{2}}(\vartheta-2\pi)\mathbb{1}_{\left[-\frac{\pi}{2},\frac{\pi}{2}\right)}(\vartheta-2\pi)\mathbb{1}_{\left[\frac{\pi}{2},\pi\right)}(\vartheta)$$

$$+ \tfrac{1}{4}f_{\frac{\Theta}{2}}(\vartheta-2\pi-\pi)\mathbb{1}_{\left[\frac{\pi}{2},\frac{3\pi}{2}\right)}(\vartheta-2\pi)\mathbb{1}_{\left[\frac{\pi}{2},\pi\right)}(\vartheta) = \tfrac{1}{4}f_{\frac{\Theta}{2}}(\vartheta-\pi)\mathbb{1}_{\left[\frac{\pi}{2},\pi\right)}(\vartheta),$$

which results in

$$f_\Phi(\vartheta) = f_{\left(\frac{\Theta}{2}+\pi\Xi\right)\bmod 2\pi}(\vartheta)$$

$$= \frac{1}{4}f_{\frac{\Theta}{2}}(\vartheta+\pi)\mathbb{1}_{\left[-\pi,-\frac{\pi}{2}\right)}(\vartheta) + \frac{1}{2}f_{\frac{\Theta}{2}}(\vartheta)\mathbb{1}_{\left[-\frac{\pi}{2},\frac{\pi}{2}\right]}(\vartheta) + \frac{1}{4}f_{\frac{\Theta}{2}}(\vartheta-\pi)\mathbb{1}_{\left[\frac{\pi}{2},\pi\right)}(\vartheta)$$

$$+ \frac{1}{4}f_{\frac{\Theta}{2}}(\vartheta+\pi)\mathbb{1}_{\left[-\pi,-\frac{\pi}{2}\right)}(\vartheta) + \frac{1}{4}f_{\frac{\Theta}{2}}(\vartheta-\pi)\mathbb{1}_{\left[\frac{\pi}{2},\pi\right)}(\vartheta)$$

$$= \frac{1}{2}f_{\frac{\Theta}{2}}(\vartheta+\pi)\mathbb{1}_{\left[-\pi,-\frac{\pi}{2}\right)}(\vartheta) + \frac{1}{2}f_{\frac{\Theta}{2}}(\vartheta)\mathbb{1}_{\left[-\frac{\pi}{2},\frac{\pi}{2}\right]}(\vartheta) + \frac{1}{2}f_{\frac{\Theta}{2}}(\vartheta-\pi)\mathbb{1}_{\left[\frac{\pi}{2},\pi\right)}(\vartheta)$$

$$= \frac{1}{2}f_{\frac{\Theta}{2}}(\vartheta)\mathbb{1}_{\left[-\frac{\pi}{2},\frac{\pi}{2}\right)}(\vartheta) + \frac{1}{2}f_{\frac{\Theta}{2}}\big((\vartheta+\pi)\bmod 2\pi\big)\mathbb{1}_{[-\pi,\pi)\setminus\left[-\frac{\pi}{2},\frac{\pi}{2}\right)}(\vartheta)$$

$$= f_\Theta(2\vartheta)\mathbb{1}_{\left[-\frac{\pi}{2},\frac{\pi}{2}\right)}(\vartheta) + f_\Theta\big((2\vartheta+\pi)\bmod 2\pi\big)\mathbb{1}_{[-\pi,\pi)\setminus\left[-\frac{\pi}{2},\frac{\pi}{2}\right)}(\vartheta).$$

With the same identifications as in part (i) we obtain the assertion. $\qquad\square$

# References

[1] A. Antoniadis, D. Leporini, and J.-C. Pesquet. Wavelet thresholding for some classes of non-Gaussian noise. *Statistica Neerlandica*, 56(4):434–453, 2002.

[2] A. Banerjee and P. Maji. Spatially constrained Student's $t$-distribution based mixture model for robust image segmentation. *Journal of Mathematical Imaging Vision*, 60(3):355–381, 2018.

[3] R. Bergmann, F. Laus, G. Steidl, and A. Weinmann. Second order differences of cyclic data and applications in variational denoising. *SIAM Journal on Imaging Sciences*, 7(4):2916–2953, 2014.

[4] J. M. Bioucas-Dias and G. Valadao. Phase unwrapping via graph cuts. *IEEE Transactions on Image processing*, (3):698–709, 2007.

[5] C. L. Byrne. *The EM Algorithm: Theory, Applications and Related Methods*. Lecture Notes, University of Massachusetts, 2017.

[6] G. Casella and R. L. Berger. *Statistical Inference*, volume 2. Duxbury Pacific Grove, CA, 2002.

[7] S. Chrétien and A. O. Hero. Kullback proximal algorithms for maximum-likelihood estimation. *IEEE Transactions on Information Theory*, 46(5):1800–1810, 2000.

[8] S. Chrétien and A. O. Hero. On EM algorithms and their proximal generalizations. *ESAIM: Probability and Statistics*, 12:308–326, 2008.

[9] C.-A. Deledalle, L. Denis, and F. Tupin. How to compare noisy patches? Patch similarity beyond Gaussian noise. *International Journal of Computer Vision*, 99(1):86–102, 2012.

[10] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38, 1977.

[11] M. Ding, T. Huang, S. Wang, J. Mei, and X. Zhao. Total variation with overlapping group sparsity for deblurring images under Cauchy noise. *Applied Mathematics and Computation*, 341:128–147, 2019.

[12] L. Dümbgen. On Tyler's $M$-functional of scatter in high dimension. *Annals of the Institute of Statistical Mathematics*, 50:471–491, 1998.

[13] L. Dümbgen, M. Pauly, and T. Schweizer. $M$-functional of multivariate scatter. *Statistics Surveys*, 9:32—-105, 2015.

[14] L. Dümbgen and D. Tyler. On the breakdown properties of some multivariate $M$-functionals. *Scandinavian Journal of Statistics*, 32:247–264, 2005.

[15] G. Frahm. *Generalized elliptical distributions: theory and applications.* PhD Thesis, Universität Köln, 2004.

[16] D. Gerogiannis, C. Nikou, and A. Likas. The mixtures of Student's *t*-distributions as a robust framework for rigid registration. *Image and Vision Computing*, 27(9):1285–1294, 2009.

[17] D. Hernandez-Stumpfhauser, F. J. Breidt, and M. J. van der Woerd. The general projected normal distribution of arbitrary dimension: Modeling and Bayesian inference. *Bayesian Analysis*, 12(1):113 – 133, 2017.

[18] D. G. Kendall. Pole-seeking Brownian motion and bird navigation. *Journal of the Royal Statistical Society*, 36:261–294, 1974.

[19] J. T. Kent and D. E. Tyler. Maximum likelihood estimation for the wrapped Cauchy distribution. *Journal of Applied Statistics*, 15(2):247–254, 1988.

[20] J. T. Kent and D. E. Tyler. Redescending *M*-estimates of multivariate location and scatter. *The Annals of Statistics*, 19(4):2102–2119, 1991.

[21] J. T. Kent, D. E. Tyler, and Y. Vard. A curious likelihood identity for the multivariate *t*-distribution. *Communications in Statistics-Simulation and Computation*, 23(2):441–453, 1994.

[22] S. Kotz and S. Nadarajah. *Multivariate t-Distributions and Their Applications.* Cambridge University Press, 2004.

[23] K. L. Lange, R. J. Little, and J. M. Taylor. Robust statistical modeling using the *t* distribution. *Journal of the American Statistical Association*, 84(408):881–896, 1989.

[24] A. Lanza, S. Morigi, F. Sciacchitano, and F. Sgallari. Whiteness constraints in a unified variational framework for image restoration. *Journal of Mathematical Imaging and Vision*, 60(9):1503–1526, 2018.

[25] F. Laus, M. Nikolova, J. Persch, and G. Steidl. A nonlocal denoising algorithm for manifold-valued images using second order statistics. *SIAM Journal on Imaging Sciences*, 10(1):416–448, March 2017.

[26] F. Laus, F. Pierre, and G. Steidl. Nonlocal myriad filters for Cauchy noise removal. *Journal of Mathematical Imaging and Vision*, 60(8):1324–1354, 2018.

[27] M. Lebrun, A. Buades, and J.-M. Morel. A nonlocal Bayesian image denoising algorithm. *SIAM Journal on Imaging Sciences*, 6(3):1665–1688, 2013.

[28] M. Lebrun, M. Colom, A. Buades, and J. Morel. Secrets of image denoising cuisine. *Acta Numerica*, 21:475–576, 2012.

[29] C. Liu and D. B. Rubin. ML estimation of the $t$ distribution using EM and its extensions, ECM and ECME. *Statistica Sinica*, pages 19–39, 1995.

[30] J. M. Borwein and R. M. Corless. Gamma and factorial in the monthly. *The American Mathematical Monthly*, 125, 2017.

[31] K. V. Mardia. *Statistics of Directional Data*. Academic Press, 1972.

[32] K. V. Mardia and P. E. Jupp. *Directional Statistics*, volume 494. John Wiley & Sons, 2009.

[33] R. A. Maronna. Robust $M$-estimators of multivariate location and scatter. *Annals in Statistics*, 4(1):51–67, 01 1976.

[34] G. McLachlan and T. Krishnan. *The EM Algorithm and Extensions*. John Wiley and Sons, Inc., 1997.

[35] J.-J. Mei, Y. Dong, T.-Z. Huang, and W. Yin. Cauchy noise removal by nonconvex ADMM with convergence guarantees. *Journal of Scientific Computing*, 74(2):743–766, 2018.

[36] X.-L. Meng and D. Van Dyk. The EM algorithm - an old folk-song sung to a fast new tune. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 59(3):511–567, 1997.

[37] S. Nadarajah and S. Kotz. Estimation methods for the multivariate $t$-distribution. *Acta Applicandae Mathematicae*, 102(1):99–118, 2008.

[38] T. M. Nguyen and Q. J. Wu. Robust Student's-$t$ mixture model with spatial constraints and its application in medical image segmentation. *IEEE Transactions on Medical Imaging*, 31(1):103–116, 2012.

[39] D. Peel and G. J. McLachlan. Robust mixture modelling using the $t$ distribution. *Statistics and Computing*, 10(4):339–348, 2000.

[40] K. B. Petersen and M. S. Pedersen. *The matrix cookbook*. Lecture Notes, Technical University of Denmark, 2008.

[41] F. Rocca, C. Prati, and A. M. Guarnieri. Possibilities and limits of SAR interferometry. *ESA SP*, pages 15–26, 1997.

[42] A. Saint-Dizier, J. Delon, and C. Bouveyron. A unified view on patch aggregation. hal-preprint hal-01865340.

[43] F. Sciacchitano, Y. Dong, and T. Zeng. Variational approach for restoring blurred images with Cauchy noise. *SIAM Journal on Imaging Sciences*, 8(3):1894–1922, 2015.

[44] S. Setzer, G. Steidl, and T. Teuber. On vector and matrix median computation. *Journal of Computational and Applied Mathematics*, 236:2200–2222, 2012.

[45] G. Sfikas, C. Nikou, and N. Galatsanos. Robust image segmentation with mixtures of Student's $t$-distributions. In *2007 IEEE International Conference on Image Processing*, volume 1, pages I – 273–I – 276, 2007.

[46] D. E. Tyler. A distribution-free $M$-estimator of multivariate scatter. *Annals of Statistics*, 15:234–251.

[47] D. E. Tyler. Statistical analysis for the angular central Gaussian distribution on the sphere. *Biometrika*, 74:579–589.

[48] A. Van Den Oord and B. Schrauwen. The Student-$t$ mixture as a natural image patch prior with application to image compression. *Journal of Machine Learning Research*, 15(1):2061–2086, 2014.

[49] D. A. van Dyk. *Construction, implementation, and theory of algorithms based on data augmentation and model reduction.* PhD Thesis, The University of Chicago, 1995.

[50] F. Wang and G. A. E. Modeling space and space-time directional data using projected Gaussian processes. *Journal of the American Statistical Association*, 109:1565–1580, 2014.

[51] F. Wang and A. E. Gelfand. Directional data analysis under the general projected normal distribution. *Statistical Methodology*, 10:113 – 127, 2013.

[52] G. S. Watson. *Statistics on Spheres.* Wiley, 1983.

[53] C. J. Wu. On the convergence properties of the EM algorithm. *The Annals of Statistics*, 11(1):95–103, 1983.

[54] Z. Yang, Z. Yang, and G. Gui. A convex constraint variational method for restoring blurred images in the presence of alpha-stable noises. *Sensors*, 18(4):1175, 2018.

[55] Z. Zhou, J. Zheng, Y. Dai, Z. Zhou, and S. Chen. Robust non-rigid point set registration using Student's-$t$ mixture model. *PloS one*, 9(3):e91381, 2014.