



Robust finite mixture modeling of multivariate unrestricted skew-normal generalized hyperbolic distributions

Mohsen Maleki¹ · Darren Wraith² · Reinaldo B. Arellano-Valle³

Received: 18 June 2017 / Accepted: 9 May 2018 / Published online: 19 May 2018
© Springer Science+Business Media, LLC, part of Springer Nature 2018

Abstract

In this paper, we introduce an unrestricted skew-normal generalized hyperbolic (SUNGH) distribution for use in finite mixture modeling or clustering problems. The SUNGH is a broad class of flexible distributions that includes various other well-known asymmetric and symmetric families such as the scale mixtures of skew-normal, the skew-normal generalized hyperbolic and its corresponding symmetric versions. The class of distributions provides a much needed unified framework where the choice of the best fitting distribution can proceed quite naturally through either parameter estimation or by placing constraints on specific parameters and assessing through model choice criteria. The class has several desirable properties, including an analytically tractable density and ease of computation for simulation and estimation of parameters. We illustrate the flexibility of the proposed class of distributions in a mixture modeling context using a Bayesian framework and assess the performance using simulated and real data.

Keywords Bayesian analysis · Finite mixtures · MCMC · Unrestricted skew-normal generalized hyperbolic family · Skew-normal · Generalized hyperbolic distribution

1 Introduction

Statistical models based on finite mixtures of distributions have been extensively used in a wide variety of applications. Applying finite mixture models to real datasets allows fitting different characteristics of the empirical distribution, such as multimodality, skewness, kurtosis and heterogeneity, across observations. For general reviews of mixture models and applications, see Hogan and Laird (1997), Böhning (2000), McLachlan and Peel (2000), Frühwirth-Schnatter (2006), Lin (2010) and Mengersen et al. (2011).

While the vast majority of work on mixture models has focused on Gaussian mixture models, in many applications the tails of Gaussian distributions are shorter than appropriate and the Gaussian shape is not suitable for highly

asymmetric data. Recent research has thus focused on fitting finite mixture models with more flexible distributional forms. The Student- t and the contaminated Gaussian distributions are two symmetric members of the scale mixtures of normal (SMN) family of distributions due to Andrews and Mallows (1974), which provide attractive heavy-tailed alternatives to the Gaussian distribution. Building upon this work is the class of scale mixtures of skew-normal (SMSN) distributions proposed by Branco and Dey (2001). The class of SMSN distributions provides location-scale density functions which depend on additional parameters of shape and kurtosis, and includes as special cases the normal and skew-normal (SN) densities, as well as the full SMN class of symmetric densities. Special symmetric and skew-symmetric heavier tail members of the SMSN family are, e.g., the Student- t , Cauchy, skew- t (ST), skew-Cauchy (SC), skew-contaminated normal (SCN) and skew-slash (SSL) distributions. Comprehensive coverage of the fundamental theory and new developments for SN and related distributions is given by Azzalini and Capitanio 2014; see also Genton (2004), Arellano-Valle and Genton (2005, 2010), Arellano-Valle and Azzalini (2006), Arellano-Valle and Azzalini (2006).

✉ Darren Wraith
d.wraith@qut.edu.au

¹ Department of Statistics, Shiraz University, Shiraz, Iran

² Institute of Health and Biomedical Innovation (IHBI), Queensland University of Technology (QUT), Brisbane, QLD, Australia

³ Department of Statistics, Universidad Católica de Chile, Santiago, Chile

Many of the different distributions within the class of SMSN distributions have been developed, and their performance assessed in the context of mixture models. Lin et al. (2007), Lin et al. (2009) and Pyne et al. (2009) studied mixtures of skew-normal distributions. Frühwirth-Schnatter and Pyne (2010) considered Bayesian inference for finite mixtures of univariate and multivariate SN and ST distributions. Basso et al. (2010) considered the robust mixture modeling based on the SMSN family. Wang et al. (2009), Lin (2010), Lee and McLachlan (2014), Vrbik and McNicholas (2012) and Forbes and Wraith (2014) considered mixtures of multivariate ST distribution. Maleki and Arellano-Valle (2017) proposed a time series model based on finite mixtures of SMSN distributions. For a review of mixtures of SN and ST distributions, see Lee and McLachlan (2013a, b).

Other distributional forms within the SMSN family of distributions have also been examined. Karlis and Santourian (2009) developed mixtures of multivariate normal inverse Gaussian distributions. Franczak et al. (2014) examined mixtures of shifted asymmetric Laplace (SAL) distributions. Morris et al. (2014) proposed mixtures of contaminated SAL distributions. Browne and McNicholas (2015) and Wraith and Forbes (2015) examined mixtures of generalized hyperbolic distributions.

In the SMSN class of distributions, although the mixing distribution typically controls the tail behavior, it can also affect the behavior of skewness (Branco and Dey 2001). A recent work on this theme providing flexibility in both skewness and heavy tails has been considered by Vilca et al. (2014) in a class of distributions referred to as multivariate SN generalized hyperbolic (SNGH) distributions. In this setting, the mixing distribution follows a generalized inverse Gaussian distribution (GIG), which has previously been demonstrated to provide considerable flexibility in modeling heavy-tailed data (Wraith and Forbes 2015).

In other recent works, a broad class of skewed distributions has been explored by Lee and McLachlan (2016) in a mixture model context focusing on the unified ST (SUT) distribution (Arellano-Valle and Azzalini 2006) and the fundamental ST distribution (Arellano-Valle and Genton 2005), including as a special case the location-scale variant of the canonical fundamental or unrestricted ST (or skew unified t ; SUT) distribution. A particular feature and advantage of the SUT distribution is that it encompasses as special cases the canonical fundamental or unrestricted SN (or skew unified normal; SUN) distribution (Arellano-Valle and Genton 2005) and other SN or ST variants (e.g., Sahu et al. 2003; Arellano-Valle et al. 2007; Lachos et al. 2007, 2010), thus providing considerable flexibility for modeling where the best fitting distribution can be chosen simply (automatically) through parameter estimation or use of model choice criteria.

In this paper, we propose a very general class of distributions which extends the previous work on the SUN and SUT

distributions by considering a mixing distribution for this class of models which follows a generalized inverse Gaussian (GIG). We refer to this new family of distributions as an unrestricted skew-normal generalized hyperbolic distribution (SUNGH). The new family provides a very general framework for a large class of distributions and has several desirable properties, including an analytically tractable density and ease of computation for simulation and estimation of parameters. The family also provides a high degree of flexibility for the modeling of complex multivariate data with different degrees of asymmetry, kurtosis and heavy tails. A particular attractiveness of this family of distributions is that it encompasses as special cases all of the distributions previously considered in the SMSN family and extensions to the unrestricted classes (e.g., SUT and SUN). Thus, this class of distributions provides a much needed unified framework where the choice of the best fitting distribution can proceed quite naturally through either parameter estimation or by placing constraints on specific parameters and assessing through model choice criteria. We illustrate the advantages of this new family in the finite mixture modeling context using a Bayesian framework.

There are some computational advantages to using a Bayesian framework in a mixture model setting. First, allowing for the influence or effect of missing data on parameter estimates is quite natural in a Bayesian setting as various patterns of missing data (e.g., class-dependent missingness) can be imputed at each MCMC iteration from the posterior predictive distribution (e.g., using a mixture model defined using open-source software such as JAGS or NIMBLE). In contrast, often quite separate and complex methods are needed for maximum likelihood estimation in these settings (Lin et al. 2009 and Wang et al. 2004). Further, for the complex distributions we consider in this paper, previous work using the EM algorithm has at times relied on approximations (Lee and McLachlan 2016) or calculations of derivatives involving complex functions (Browne and McNicholas 2015) for the estimation of parameters. This difficulty also extends to the estimation of the standard errors for parameters (if they are available) using asymptotic approximations to the observed information matrix if the sample size is large or resorting to a bootstrap method which is computationally demanding (Basso et al. 2010). At times, the standard errors for the parameters are also unavailable (particularly for the GH distribution) or not mentioned (e.g., Browne and McNicholas 2015). This is not to say that estimation in the Bayesian setting is devoid of potential computational issues, in particular the issue of label switching is a more prominent issue compared to methods using ML estimation (Mengersen et al. (2011)).

The paper is organized as follows. In Sect. 2, we provide some background to the SUN and GIG distributions. Section 3 outlines the details and properties of the new SUNGH

family. In Sect. 4, we present a Bayesian analysis of a finite mixture model following a SUNGH distribution. In Sect. 5, we illustrate the performance of the proposed approach on real and simulated data. Finally, in Sect. 6, we present our main conclusions and discuss some areas of further research.

2 SUN and GIG distributions

2.1 Preliminaries

Following Arellano-Valle and Genton (2005), Arellano-Valle and Azzalini (2006) and Arellano-Valle et al. (2007), we say that a $p \times 1$ random vector X follows an unrestricted skew-normal (SUN) with $p \times 1$ location vector μ , $p \times p$ positive definite dispersion matrix Σ and $p \times q$ skewness parameter matrix Λ , denoted by $X \sim \text{SUN}_{p,q}(\mu, \Sigma, \Lambda)$, if its probability density function (pdf) is

$$f(x|\mu, \Sigma, \Lambda) = 2^q \phi_p(x|\mu, \Psi) \Phi_q(\Lambda^\top \Psi^{-1}(x - \mu)|\Upsilon), \quad x \in \mathbb{R}^p, \tag{1}$$

where $\Psi = \Sigma + \Lambda\Lambda^\top$, $\Upsilon = I_q - \Lambda^\top \Psi^{-1} \Lambda = (I_q + \Lambda^\top \Sigma^{-1} \Lambda)^{-1}$, and $\phi_k(\cdot|\mu, \Psi)$ and $\Phi_q(\cdot|\Upsilon)$ are, respectively, the pdf and cumulative distribution function (cdf) of the multivariate normal distributions given by $N_p(\mu, \Psi)$ and $N_q(0, \Upsilon)$. The SUN class of multivariate distributions defined by (1) contains various special cases. For instance, we recover the multivariate normal when $\Lambda = 0$, the multivariate SN which called here restricted SN (rMSN) when $q = 1$, and the multivariate SN of Sahu et al. (2003) when $p = q$ and Λ being a diagonal matrix. In fact, the SUN distribution becomes an important special case of the unified SN distribution (SUN) studied by Arellano-Valle and Azzalini (2006).

The random vector $X \sim \text{SUN}_{p,q}(\mu, \Sigma, \Lambda)$ can be stochastically represented from different ways. According to Arellano-Valle et al. (2006), the SUN random vector X has selection representation given by

$$X = \mu + (V_1|V_0 > 0), \tag{2}$$

where the condition $V_0 > 0$ means that each element of V_0 is positive, and

$$\begin{pmatrix} V_0 \\ V_1 \end{pmatrix} \sim N_{p+q} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} I_q & \Lambda^\top \\ \Lambda & \Sigma \end{pmatrix} \right).$$

The representation in (2) becomes a selection representation of the rMSN distribution when $q = 1$, i.e., when the latent vector V_0 is replaced by a one-dimensional normal random variable V_0 . Also, if we let $V_0 = W_0$ and $V_1 = W_1 + \Lambda W_0$, where $W_0 \sim N_q(0, I_q)$ and $W_1 \sim N_p(0, I_p)$ are independent,

it follows from (2) that the stochastic representation of X is given by

$$X = \mu + \Lambda |W_0| + \Sigma^{1/2} W_1, \tag{3}$$

where $|W_0|$ is the vector formed with the absolute value of each component of W_0 . For more details, see Arellano-Valle et al. (2006), Arellano-Valle and Azzalini (2006) and Arellano-Valle et al. (2007). In particular, the mean vector and covariance matrix of X are given by $E[X] = \mu + \sqrt{2/\pi} \Lambda \mathbf{1}_p$ and $\text{Cov}[X] = \Psi - \frac{2}{\pi} \Lambda \mathbf{1}_q \mathbf{1}_q^\top \Lambda^\top$, respectively, where $\mathbf{1}_q$ denotes the vector of ones with length q .

In this work, we consider the extension of the scale mixtures of rMSN (SMRSN or SMSN) distributions to the scale mixtures of SUN (SMSUN) distributions. Specifically, we consider the family of random vectors defined by

$$Y = \mu + \kappa(U)^{1/2} X, \tag{4}$$

where $X \sim \text{SUN}_{p,q}(0, \Sigma, \Lambda)$, $\kappa(\cdot)$ is a positive scale function and U is a mixing random variable which is independent of X . For our proposed SUNGH distribution, we consider the SMSUN class of distributions defined by (4) when the mixing random variable U follows a GIG distribution.

2.2 The family of GIG distribution

The GIG class is a rich family of flexible distributions with positive support that has been studied by several authors. For instance, see Good (1953), Barndorff-Nielsen and Halgreen (1977), Jørgensen (1982), among others. Thus, the choice of a GIG distribution for the scale mixing variable U in (4) is a natural candidate and provides a highly flexible unified class of multivariate distributions for multivariate statistical analysis.

The GIG distribution has several but equivalent representations in terms of its parameterization. In this paper, and in order to simplify and have closed-form posterior distributions in the Bayesian framework adopted here, we consider (without loss of generality) the following two representations:

First representation $\text{GIG}^*(v, \gamma, \rho)$: A random variable U has a GIG distribution, denoted by $U \sim \text{GIG}^*(v, \gamma, \rho)$, if its pdf is given by

$$\begin{aligned} \mathcal{GIG}^*(u|v, \gamma, \rho) &= \left(\frac{\gamma}{\rho}\right)^v \frac{u^{v-1}}{2K_v(\rho\gamma)} \\ &\exp\left(-\frac{1}{2}\left(\frac{\rho^2}{u} + \gamma^2 u\right)\right), \quad u > 0, \end{aligned} \tag{5}$$

where $K_r(x)$ is the modified Bessel function of the third kind of order r evaluated at x , and the parameter spaces are given by $\gamma > 0$, $\rho > 0$ and $-\infty < v < +\infty$.

Second representation $GIG_*(v, \psi, \eta)$: A random variable U follows a GIG distribution denoted by $U \sim GIG_*(v, \psi, \eta)$, if its pdf is given by

$$\mathcal{GIG}_*(u | v, \psi, \eta) = \frac{(u/\eta)^{v-1}}{2\eta K_v(\psi)} \exp\left(-\frac{\psi}{2}\left(\frac{u}{\eta} + \frac{\eta}{u}\right)\right), u > 0, \tag{6}$$

where $K_r(x)$ is defined previously and the parameter spaces are $\psi > 0, \eta > 0$ and $-\infty < v < +\infty$. This representation will be used to simplify the posterior representation of the GIG parameters. In this case, the m th moment of the random variable $U^{1/2}$ is given by

$$E\left(U^{m/2}\right) = \frac{K_{v+m/2}(\psi)}{K_v(\psi)} \eta^{m/2}, m = 1, 2, \dots$$

The equivalence between both representations of the GIG distribution considered in (5) and (6) is obtained by observing the one-to-one relationship between their parameters given by $\psi = \rho\gamma$ and $\eta = \rho/\gamma$. Particular members of the GIG class lead to a variety of skewed distributions belonging to the proposed family. The inverse Gaussian is one member of this class which has been extensively studied by Chhikara and Folks (1989), Seshadri (1993) and Johnson et al. (1994, chap. 15). Two additional members of the GIG class are the hyperbola and the positive hyperbolic distributions, both of which have been studied by Barndorff-Nielsen (1978) and Barndorff-Nielsen and Blaesild (1980). The exponential, gamma and inverse gamma distributions are also special members of the GIG family. For a recent study on these distributions, see Vilca et al. (2014) and references therein.

In this paper, we define the multivariate random variable Y via (4), and by considering a multivariate SUN random variable X according to (3) and a GIG scale random variable U distributed according to the second representation in (6). As mentioned previously, we refer to this proposed family as SUNGH distributions.

3 The family of SUNGH distributions

An alternative way to define SUNGH distribution follows by replacing Eq. (3) in Eq. (4). From this, we can say that a $p \times 1$ random vector Y follows a SUNGH distribution if

$$Y = \mu + \Lambda W + \kappa(U)^{1/2} \Sigma^{1/2} W_1, \tag{7}$$

where μ is a $p \times 1$ location vector, Σ is a $p \times p$ scale matrix, Λ is a $p \times q$ shape matrix, $W = \kappa^{1/2}(U) |W_0|$, $W_0 \sim N_q(\mathbf{0}, I_q)$, $W_1 \sim N_p(\mathbf{0}, I_p)$ and $U \sim GIG_*(v, \psi, \eta)$, with W_0, W_1 and U being independent random quantities. These assumptions also imply that W is also independent of

W_1 . Note that if we set $W = U, \kappa(u) = u$ and $q = 1$ we obtain the GH distribution proposed by McNeil et al. (2005) and considered in the mixture model context by Browne and McNicholas (2015). For this reason, the GH distribution is more restrictive (less flexible) compared to the SUNGH distribution. Since the conditional distribution of Y given $U = u$ is given by $Y | U = u \sim SUN_{p,q}(\mu, \kappa(u) \Sigma, \kappa(u)^{1/2} \Lambda)$, the marginal pdf of Y becomes the infinite mixture of the SUN pdf in (1) given by

$$f(y | \mu, \Sigma, \Lambda, \varpi) = 2^q \int_0^\infty \phi_p(y | \mu, \kappa(u) \Sigma) \Phi_q\left(\kappa(u)^{-1/2} \Lambda^\top \Psi^{-1}(y - \mu) | \Upsilon\right) \mathcal{GIG}_*(u | \varpi) du, \tag{8}$$

$y \in R^p$, where $\varpi = (v, \psi, \eta)^\top$, and Ψ and Υ defined as in (1). In what follows, we will refer to the SUNGH random vector in (7) as $Y \sim SUNGH_{p,q}(\mu, \Sigma, \Lambda, \varpi)$.

Note that there are some identifiability issues concerning the GIG parameters ϖ and skewness matrix Λ . Using (8) the density is not identifiable as for any parameter $c > 0$, the parameters $(\mu, \Sigma, \Lambda, v, \psi, \eta)$ and $(\mu, c\Sigma, c\Lambda, v, \psi/c, c\eta)$ yield the same density. A simple fix which results in an identifiable density is to set $\eta = 1$ and so $\varpi = (v, \psi)^\top$. An alternative parameterization which can provide for greater flexibility is discussed in Wraith and Forbes (2015). Further, any permutation matrix can be multiplied by W from the stochastic representations (7) without any changes to the distribution of Y , so sorting Λ by the norm of the columns or some other sorting method is also needed to ensure identifiability of the proposed model.

Varying the scale mixing function $\kappa(U)$ for a given distribution of U belonging to the $GIG_*(\varpi)$ class leads to a variety of members in the SUNGH family. Alternatively, we can fix the scale function and vary the distribution of U within the $GIG_*(\varpi)$ class. In the latter case, a convenient choice for the scale function is $\kappa(u) = u$, for which the pdf in (8) becomes

$$f(y | \mu, \Sigma, \Lambda, \varpi) = 2^q \mathcal{GH}_p(y | \mu, \Psi, \mathbf{0}, v') GH_q(\mathbf{B} | \mathbf{0}, \Upsilon, \mathbf{0}, v''), y \in R^p, \tag{9}$$

where $v' = (v, \sqrt{\psi/\eta}, \sqrt{\psi\eta})^\top$, $v'' = (v - p/2, \sqrt{\psi/\eta}, q'(y))^\top$, $q'(y)^2 = (y - \mu)^\top \Psi^{-1}(y - \mu) + \psi\eta$, $\Psi = \Sigma + \Lambda\Lambda^\top$, $\Upsilon = I_q - \Lambda^\top \Psi^{-1} \Lambda$ and $\mathbf{B} = \Lambda^\top \Psi^{-1}(y - \mu)$, \mathcal{GH}_p and GH_q denote the p -variate pdf and q -variate cdf of the generalized Hyperbolic distribution, respectively (Wraith and Forbes 2015).

The flexibility of the SUNGH family proposed in (8) can also be observed by varying the value of the dimension q . In fact, for $q = 1$ (the restricted case) we obtain as a special

case of (8) the SN generalized hyperbolic (SNGH) distributions considered in Vilca et al. (2014), and thus some known SMSN (or SMRSN) distributions, as well the corresponding symmetrical variant for $\Lambda = \mathbf{0}$.

A special case of the GIG distribution is the gamma distribution, so the proposed family of distributions covers the canonical fundamental unrestricted skew-normal (CFUSN) distribution of Arellano-Valle and Genton (2005), and the canonical fundamental unrestricted skew-t (CFUST) distribution of Lee and McLachlan (2016). Subsequently, a mixture model approach covering these distributions contains the finite mixtures of CFUSN and CFUST. By considering (9) in the symmetric case, the SUNGH and GH studied by Wraith and Forbes (2015) and Browne and McNicholas (2015) are similar, but in the asymmetric case these families are different. In particular, a greater degree of flexibility is available for the SUNGH family by allowing the skewness parameter to be multivariate ($p \times q$) rather than $p \times 1$. The SUNGH family also has several desirable properties outlined in Propositions 2 to 6 below which will allow the family to be used in a variety of statistical models (e.g., mixed models and regression).

Known members of the SMSN family contained in the SNGH family are the SN, ST, SSL and skew-Laplace (SLP), and their respective symmetric versions. In the unrestricted case ($q > 1$), the proposed family contains several subfamilies of distributions (symmetric and asymmetric) considered in the literature. For instance, if in (9) we let $q = p$, $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$ and $\kappa(u) = 1$, then the multivariate skew-normal distribution of Sahu et al. (2003) is obtained. Finally, if $\Lambda = \mathbf{0}$ (symmetric case) and $\kappa(u) = u$, then (9) becomes the symmetric generalized hyperbolic (GH) distribution introduced by Barndorff-Nielsen and Halgreen (1977).

In the following propositions, we present some necessary and useful properties of the SUNGH family for the next sections. The proofs of these results are presented in ‘‘Appendix.’’

Proposition 1 *Let $Y \sim \text{SUNGH}_{p,q}(\mu, \Sigma, \Lambda, \varpi)$. Then, the following results hold:*

- (a) if $k_1 = E[\kappa(U)^{1/2}] < \infty$, then $E[Y] = \mu + \sqrt{\frac{2}{\pi}}k_1\Lambda\mathbf{1}_q$,
- (b) if $k_2 = E[\kappa(U)] < \infty$, then $\text{Var}[Y] = k_2\psi - \frac{2}{\pi}\Lambda[k_2I_q - (k_2 - k_1^2)\mathbf{1}_q\mathbf{1}_q^T]\Lambda^T$.

Proposition 2 *Let $Y \sim \text{SUNGH}_{p,q}(\mu, \Sigma, \Lambda, \varpi)$. Then, $Y \sim \text{SUNGH}_{p,q+m}(\mu, \Sigma, \Lambda^*, \varpi)$ for each $m = 1, 2, \dots$, where $\Lambda^* = (\Lambda_{p \times q} \ \mathbf{0}_{p \times m})$ or $\Lambda^* = (\mathbf{0}_{p \times m} \ \Lambda_{p \times q})$.*

Proposition 3 *Let $Y \sim \text{SUNGH}_{p,q}(\mu, \Sigma, \Lambda, \varpi)$. Then, for each $\mathbf{b} \in \mathbf{R}^n$ and full row rank matrix $\mathbf{B} \in \mathbf{R}^{n \times p}$ we have*

$$X = \mathbf{b} + \mathbf{B}Y \sim \text{SUNGH}_{n,q}(\mathbf{b} + \mathbf{B}\mu, \mathbf{B}\Sigma\mathbf{B}^T, \mathbf{B}\Lambda, \varpi).$$

Proposition 4 *Let $Y \sim \text{SUNGH}_{p,q}(\mu, \Sigma, \Lambda, \varpi)$. Partition $Y = (Y_1^T, Y_2^T)^T$, where the first and second sub-vectors are of dimensions $p_1 \times 1$ and $p_2 \times 1$, respectively, with $p_1 + p_2 = p$. The corresponding partition of the parameters (μ, Σ, Λ) is*

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}, \Lambda = \begin{pmatrix} \Lambda_1 \\ \Lambda_2 \end{pmatrix},$$

where μ_i, Σ_{ii} and Λ_i have dimensions $p_i \times 1, p_i \times p_i$ and $p_i \times q$, respectively, for $i = 1, 2$. Then, the marginal distribution of Y_i is $\text{SUNGH}_{p_i,q}(\mu_i, \Sigma_{ii}, \Lambda_i, \varpi), i = 1, 2$.

Proposition 5 *If under the same conditions of Proposition 4, we have $\Sigma_{12} = \Sigma_{21} = \mathbf{0}$ then a necessary and sufficient condition to have null correlation between Y_1 and Y_2 is that $\Lambda_1 = \mathbf{0}$ or $\Lambda_2 = \mathbf{0}$.*

Proposition 6 *Consider the same conditions of Proposition 4 with the partition of shape matrix $\Lambda = (\Lambda_{ij})_{i,j=1,2}$, where Λ_{ij} has dimension $p_i \times q_j$, with $q_1 + q_2 = q$. If $\Sigma_{12} = \Sigma_{21}^T = \mathbf{0}$ and $\Lambda_{12} = \mathbf{0}$ or $\Lambda_{21} = \mathbf{0}$, then $Y_i \sim \text{SUNGH}_{p_i,q_i}(\mu_i, \Sigma_{ii}, \Lambda_{ii}, \varpi), i = 1, 2$, and $\text{Cov}(Y_1, Y_2) = -\frac{2}{\pi}k_1^2\Lambda_{12}\mathbf{1}_{q_1}\mathbf{1}_{q_2}^T\Lambda_{21}^T$.*

4 Finite mixtures of SUNGH family

4.1 FM-SUNGH model

In this section, we consider finite mixtures of the proposed SUNGH family of distributions (hereafter FM-SUNGH). To establish notation, we consider the usual mixture model defined as

$$f(y; \Theta, \mathbf{p}) = \sum_{k=1}^K p_k f(y; \Theta_k), \tag{10}$$

where $\Theta = (\Theta_1, \dots, \Theta_K)$, with $\Theta_k = (\mu_k, \Sigma_k, \Lambda_k, \nu_k, \psi_k, \eta_k), k = 1, \dots, K, \mathbf{p} = (p_1, \dots, p_K)^T$ (for which $p_k > 0, k = 1, \dots, K$ and $\sum_{k=1}^K p_k = 1$), $\nu_k = (\nu_{k1}, \dots, \nu_{kp})^T, \psi_k = (\psi_{k1}, \dots, \psi_{kp})^T, \eta_k = (\eta_{k1}, \dots, \eta_{kp})^T$ and $f(y; \Theta_k)$ given by (8). This model hereafter will be called FM-SUNGH. The identifiability of mixtures of distributions has been studied by Teicher (1963) and Holzmann et al. (2006) to ensure that the FM-SUNGH is identifiable.

The SUNGH family is a rich class of distributions, and various particular forms from this family have been considered over the last few years in the case of mixture models. In Table 1, we outline details of some of the distributions and the corresponding parameters within the SUNGH family.

Using the mixture model representation in (10), for each i.i.d. sample in the form of Y_1, \dots, Y_n , we can utilize

an (latent) indicator (allocation) variables Z_1, \dots, Z_n , to assign observations to belong to different components of the mixture ($k = 1, \dots, K$). The standard assumption for the allocation random variables Z_1, \dots, Z_K is that they follow a multinomial joint distribution: $Z_i = (Z_{i1}, \dots, Z_{iK}) \sim \text{Multinomial}(K, p_1, \dots, p_K)$, so that $P(Z_i = k) = p_k$; $i = 1, \dots, n, k = 1, \dots, K$. In terms of Z_i , we can conclude that

$$Y_i | Z_i = k \stackrel{\text{ind.}}{\sim} \text{SUNGH}(\Theta_k), P(Z_i = k) = p_k.$$

Let $C = \{Y, U, W, Z\}$ denote the complete data, where $Y = (Y_1^\top, \dots, Y_n^\top)^\top$ is the observed variable and $U = (U_{11}, \dots, U_{1K}, \dots, U_{n1}, \dots, U_{nK})^\top$, $W = (W_{11}^\top, \dots, W_{1K}^\top, \dots, W_{n1}^\top, \dots, W_{nK}^\top)^\top$ and $Z = (Z_1, \dots, Z_n)^\top$ are the latent or unobserved variables. If we consider the SUNGH stochastic representation (7) in terms of a finite mixture model for $i = 1, \dots, n$ and $k = 1, \dots, K$, a hierarchical representation is

$$Y_i | W_{ik} = w_{ik}, U_{ik} = u_{ik}, Z_i = k \stackrel{\text{ind.}}{\sim} N_p(\mu_k + \Lambda_k w_{ik}, \kappa(u_{ik}) \Sigma_k), \tag{11}$$

$$W_{ik} | U_{ik} = u_{ik}, Z_i = k \stackrel{\text{ind.}}{\sim} HN_q(\mathbf{0}, \kappa(u_{ik}) I_q), \tag{12}$$

$$U_{ik} | Z_i = k \stackrel{\text{ind.}}{\sim} \text{GIG}_*(v_k, \psi_k, \eta_k), \tag{13}$$

$$Z \sim \text{Multinomial}(K, p_1, \dots, p_K), \tag{14}$$

where HN_q denotes the q -variate right half-normal distribution.

The model’s complete data likelihood function is then given by

$$L(\Theta | C) = \prod_{i=1}^n \prod_{k=1}^K (\phi_p(y_i | \mu_k + \Lambda_k w_{ik}, \kappa(u_{ik}) \Sigma_k) H\phi_q(w_{ik} | \mathbf{0}, \kappa(u_{ik}) I_q) \mathcal{GIG}_*(u_{ik} | v_k, \psi_k, \eta_k) P(Z_i = k)), \tag{15}$$

where $H\phi_q(w | \mathbf{0}, \cdot) = \phi_q(w | \mathbf{0}, \cdot) I(w > \mathbf{0})$ is the q -variate right half-normal pdf.

4.2 Bayesian analysis

4.3 Priors

In this section, we choose priors for the parameters Θ which will be used in Applications section. By assuming independency between the different types of parameters in Θ and that the skewness matrix of each mixture component be in the form of $\Lambda_k = (\lambda_{k1} | \dots | \lambda_{kq})$, prior distributions for some of the FM-SUNGH model parameters are given by

$$p \sim \text{Dir}(\delta_1, \dots, \delta_K), \mu_k \sim N_p(m_k, M_k), \Sigma_k \sim IW_{t_k}(T_k), \lambda_{kt} \sim N_p(\ell_{kt}, L_{kt}), \quad t = 1, \dots, p$$

for $k = 1, \dots, K$, and where Dir and IW denote the Dirichlet and inverse Wishart distributions, respectively. An alternative representation of the skewness matrices priors and posteriors in the Gibbs updates is provided in ‘‘Appendix.’’ Prior distributions of the scaled factor variables for $k = 1, \dots, K$ are:

$$v_k \sim N(\mu_k, \sigma_k^2), \eta_k \sim \text{GIG}^*(\alpha_k, \chi_k, \varepsilon_k), \psi_k \sim \exp(\zeta_k).$$

4.3.1 Posteriors

By considering the likelihood function (15) and the priors specified previously, the joint posterior of Θ is given by

$$\pi(\Theta, u, w, z | y) \propto L(\Theta | y, u, w, z) \pi(\Theta).$$

The above joint posterior is intractable, but we can use an MCMC method such as Gibbs sampling and Metropolis–Hastings to draw samples using the conditional posterior distributions. To establish notation, let $B_k = \{i, Z_i = k\}$ be the set of observation indices for those y_i classified into the k th cluster and n_k is equal to the number of observations allocated to the k th component (cluster). Apart from the parameters for the scaled factor variables, all conditional posterior distributions have closed form and are as follows: (Note that $\Theta_{(-\theta)}$ denotes the set of parameters without its element θ .)

$$p | \Theta_{(-p)}, y, u, w, z_i = k \sim \text{Dir}(\delta_{pos.1}, \dots, \delta_{pos.K}),$$

where

$$\delta_{pos.k} = \delta_k + n_k; k = 1, \dots, K. \tag{16}$$

$$\mu_k | \Theta_{(-\mu_k)}, y, u, w, z_i = k \sim N_p(\mu, \Sigma), k = 1, \dots, K,$$

where

$$\mu = \Sigma \left(M_k^{-1} m_k + \sum_{B_k} \kappa(u_{ik})^{-1} \Sigma_k^{-1} (y_i - \Lambda_k w_{ik}) \right), \Sigma = \left[M_k^{-1} + \sum_{B_k} \kappa(u_{ik})^{-1} \Sigma_k^{-1} \right]^{-1}, \tag{17}$$

$$\Sigma_k | \Theta_{(-\Sigma_k)}, y, u, w, z_i = k \sim IW_{t_k+n}(T), k = 1, \dots, K,$$

where

$$T = \sum_{B_k} \kappa(u_{ik})^{-1} (y_i - \mu_k - \Lambda_k w_{ik}) (y_i - \mu_k - \Lambda_k w_{ik})^\top + T_k, \tag{18}$$

$$\lambda_{kt} | \Theta_{(-\lambda_{kt})}, y, u, w, z_i = k \sim N_p(\mu, \Sigma); k = 1, \dots, K; t = 1, \dots, p,$$

where

$$\mu = \Sigma \left(L_{kt}^{-1} \ell_{kt} + \sum_{B_k} \kappa(u_{ik})^{-1} w_{ik(t)} \Sigma_k^{-1} \right)$$

Table 1 Summary of non-normal finite mixture models belonging to the FM-SUNGH model

Finite mixture models	Component density	FM-SUNGH	References
FM of restricted skew-normal (FM-rMSN)	$f(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\delta}) = 2\phi_p(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma}) \Phi_1(\boldsymbol{\delta}^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \boldsymbol{\Upsilon}),$ $\boldsymbol{\Upsilon} = 1 - \boldsymbol{\delta}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}$	$\kappa(u) = 1, \quad q = 1$	Pyne et al. (2009)
FM of unrestricted skew-normal (FM-CFUSN)	$f(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\Lambda}) =$ $2^q \phi_p(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Psi}) \Phi_q(\boldsymbol{\Lambda}^\top \boldsymbol{\Psi}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \boldsymbol{\Upsilon}),$ $\boldsymbol{\Psi} = \boldsymbol{\Sigma} + \boldsymbol{\Lambda} \boldsymbol{\Lambda}^\top, \quad \boldsymbol{\Upsilon} = \mathbf{I}_q - \boldsymbol{\Lambda}^\top \boldsymbol{\Psi}^{-1} \boldsymbol{\Lambda} =$ $(\mathbf{I}_q + \boldsymbol{\Lambda}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\Lambda})^{-1}$	$\kappa(u) = 1$	Lin (2009)
FM of restricted skew-t (FM-rMST)	$f(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\delta}, n) =$ $2t_p(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma}, n) T_1(\boldsymbol{\delta}^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \sqrt{\frac{n+p}{n+d}} \mathbf{0}, \boldsymbol{\Upsilon}, n + p),$ $d = (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}), \boldsymbol{\Upsilon} = 1 - \boldsymbol{\delta}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}$	$\eta \downarrow 0, \quad v = \psi = n/2, \quad q = 1$	Pyne et al. (2009) and Vrbik and McNicholas (2012)
FM of unrestricted skew-t (FM-CFUST)	$f(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\delta}, n) =$ $2t_p(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma}, n) T_q(\boldsymbol{\Lambda}^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \sqrt{\frac{n+p}{n+d}} \mathbf{0}, \boldsymbol{\Lambda}, n + p),$ $\boldsymbol{\Lambda} = \text{diag}(\boldsymbol{\delta}), d = (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}), \boldsymbol{\Lambda} =$ $\mathbf{I}_p - \boldsymbol{\delta}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}$	$\eta \downarrow 0, \quad v = \psi = n/2$	Lin (2010), Lee and McLachlan (2014) and Lee and McLachlan (2016)
FM of multivariate student-t (FM-T)	$f(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma}, n) =$ $\frac{\Gamma((n+p)/2)}{\Gamma(n/2) \pi^{p/2} \boldsymbol{\Sigma} ^{n/2}} \left(1 + \frac{1}{n} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)^{-\frac{n+p}{2}}$	$\eta \downarrow 0, \quad v = \psi = n/2, \quad \boldsymbol{\Lambda} = \mathbf{0}, \quad q = 1$	McLachlan and Peel (2000)
FM of scale mixtures of skew-normal (FM-SMSN)	$f(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\lambda}) =$ $2 \int_0^\infty \phi_p(\mathbf{y} \boldsymbol{\mu}, u^{-1} \boldsymbol{\Sigma}) \Phi_1(u^{1/2} \boldsymbol{\lambda}^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})) dH(u)$	See Section 3 in Vilca et al. (2014)	Basso et al. (2010) and Maleki and Arellano-Valle (2017)
FM of generalized hyperbolic (FM-GH)	$f(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\alpha}, \lambda, \omega) =$ $\frac{e^{-\boldsymbol{\alpha}^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})} \left(\frac{\omega+d}{\gamma}\right)^{\frac{1}{2} - \frac{p}{4}} K_{\lambda - \frac{p}{2}} \left(\sqrt{\gamma(\omega+d)}\right)}{(2\pi)^{p/2} \boldsymbol{\Sigma} ^{1/2}} \frac{K_\lambda(\omega)}{K_\lambda(\omega)},$ $\gamma = \omega + \boldsymbol{\alpha}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}, d = (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$	Let $\mathbf{W} = U$ in Eq. (7), $\kappa(u) = u$ and $q = 1$	McNeil et al. (2005) and Browne and McNicholas (2015)

$$\begin{aligned} & (y_i - \mu_k - \Lambda_k(-t)w_{ik(-t)}), \\ \Sigma &= \left(L_{kt}^{-1} + \sum_{B_k} \kappa(u_{ik})^{-1} w_{ik(t)}^2 \Sigma_k^{-1} \right)^{-1}, \end{aligned} \tag{19}$$

where $\Lambda_k(-t)$ denotes the $p \times (q - 1)$ skewness matrix Λ_k with the t th column eliminated, $w_{ik(-t)}$ denotes the $(q - 1) \times 1$ vector w_{ik} vector with the t th element eliminated, and $w_{ik(t)}$ denotes the t th element of the vector w_{ik} .

The full conditional posterior distribution for the latent variables Z_i, U_{ik} and W_{ik} , for $i = 1, \dots, n; k = 1, \dots, K$, are given by:

$Z_i | \Theta, y, u, w \sim \text{Multinomial}(K, p_{p,1}, \dots, p_{p,K})$, where

$$p_{p,k} = \frac{p_k f(y_i; \Theta_k)}{\sum_{j=1}^K p_j f(y_i; \Theta_j)}, k = 1, \dots, K, \tag{20}$$

$U_{ik} | \Theta, y, w, z_i = k \sim \text{GIG}^*(a_u, b_u, c_u)$, where $\kappa(u) = u$ and

$$\begin{aligned} a_u &= v_{ik} - \frac{p + q}{2}, \\ b_u &= (\psi_{ik}/\eta_{ik})^{1/2}, \\ c_u &= \left(\left[w_{ik}^\top w_{ik} + (y_i - \mu_k - \Lambda_k w_{ik})^\top \right. \right. \\ & \quad \left. \left. \Sigma_k^{-1} (y_i - \mu_k - \Lambda_k w_{ik}) \right] + \psi_{ik} \eta_{ik} \right)^{1/2}. \end{aligned} \tag{21}$$

$W_{ik} | \Theta, y, u, z_i = k \sim \text{HN}_q(\mu, \Sigma)$, where

$$\begin{aligned} \mu &= \kappa(u_{ik})^{-1} \Sigma \Lambda_k^\top \Sigma_k^{-1} (y_i - \mu_k), \\ \Sigma &= \kappa(u_{ik}) \left(I_q + \Lambda_k^\top \Sigma_k^{-1} \Lambda_k \right)^{-1}. \end{aligned} \tag{22}$$

Finally, the full conditional posterior for the scaled factor variables $v_k, \psi_k, \eta_k, k = 1, \dots, K$, is as follows: $\eta_k | \Theta_{(-\eta_k)}, u, w, z_i = k \sim \text{GIG}^*(a_\eta, b_\eta, c_\eta)$, where

$$\begin{aligned} a_\eta &= \alpha_k - v_k n_k, \\ b_\eta &= \left(\chi_k^2 + \psi_k \sum_{B_k} 1/u_{ik} \right)^{1/2}, \\ c_\eta &= (\varepsilon_k^2 + \psi_k \sum_{B_k} u_{ik})^{1/2}. \end{aligned} \tag{23}$$

The full conditional posterior density of $v_k, k = 1, \dots, K$ is proportional to:

$$\pi(v_k | \Theta_{(-v_k)}, u, w, Z_i = k) \sim \pi_1(v_k) N(\mu_k, \sigma_k^2), \tag{24}$$

where $\pi_1(v_k) = (K v_k (\psi_k))^{-n_k} \prod_{B_k} (u_{ik}/\eta_k)^{v_k}$.

The full conditional posterior density of $\psi_k, k = 1, \dots, K$ is also proportional to:

$$\begin{aligned} \pi(\psi_k | \Theta_{(-\psi_k)}, u, w, z_i = k) &\sim \pi_2(\psi_k) \\ E\left(\zeta_k + \sum_{B_k} (u_{ik}/\eta_k + \eta_k/u_{ik})/2\right), \end{aligned} \tag{25}$$

where $\pi_2(\psi_k) = (K v_k (\psi_k))^{-n_k}$ and $E(\varphi)$ denotes the density of the exponential distribution with rate parameter φ .

Note that (24) and (25) do not have closed forms, but a Metropolis–Hastings or rejection sampling step can be embedded in the MCMC scheme to obtain draws from them.

5 Applications

In this section, we present a simulation study and applications on two real datasets to evaluate the performance of the proposed SUNGH model for clustering problems. For illustrative purposes, we choose K to be equal to two for all models presented.

5.1 Simulated data

To illustrate some of the differences between the SUNGH family of models, we consider the case of two clusters each sampled from a four-dimensional SUNGH distribution with known parameters which are slightly separated from each other. For the first and second cluster

$$\mu_1 = (4, 4, 4, 4)^\top, \quad \Sigma_1 = \Sigma_2 = \begin{pmatrix} 1 & \dots & 0.5 \\ \vdots & \ddots & \vdots \\ 0.5 & \dots & 1 \end{pmatrix} \text{ and}$$

$$\Lambda_1 = \begin{pmatrix} -4 & -4 \\ 1 & 3 \\ -4 & -4 \\ 1 & 3 \end{pmatrix},$$

$$\mu_2 = (-1, -1, -1, -1)^\top \text{ and } \Lambda_2 = \begin{pmatrix} 4 & -5 \\ 1 & 2 \\ 4 & -5 \\ 1 & 2 \end{pmatrix},$$

respectively. Both clusters shared the same parameters for the $\text{GIG}^*(v, \psi, \eta)$ distribution, where $v = -0.5, \psi = 1$ and $\eta = 1$. The sample size for each cluster is 300 and 450, respectively. A plot of the simulated data is shown in Fig. 1 with the observations belonging to each cluster labeled by different colors.

For estimation of the different models, largely non-informative prior distributions were used for each of the component parameters: $\mu = (\mu_1, \dots, \mu_4)^\top \sim N_4(\mathbf{0}, \Sigma)$, where $\Sigma = 10^3 I_4, \Sigma \sim IW_\tau(T)$, where $\tau = 4$ and $T = I_4$

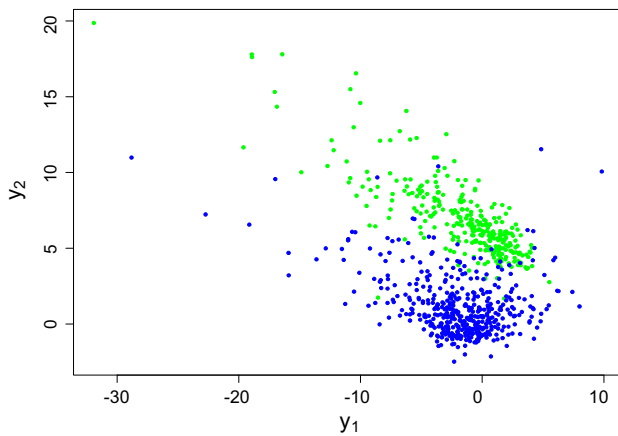


Fig. 1 Plot of results for simulated data. Colors indicate the groups to which observations belong in the first two dimensions. (Color figure online)

with skewness matrix $\Lambda_{4 \times 2}$ with priors of its columns as $\lambda_t \sim N_4(\ell_t, L_t)$ for which $\ell_t = \mathbf{0}$ and $L_t = 10^3 I_4$ and for $t = 1, 2$ (and in the matrix variate prior of the skewness matrix we can consider that $\Lambda_{4 \times 2} \sim MN_{4,2}(\mathbf{0}, 10^3 I_4, 10^3 I_4)$), $\nu \sim N(0, 10^3)$, $\eta \sim \text{GIG}^*(0.001, 2000, 0)$, $\psi \sim \exp(0.1)$ and $\mathbf{p} \sim \text{Dir}(1, \dots, 1)$. Also, we chose $\kappa(u) = u$ for the scale mixing function. All computations are implemented in the R software version 3.3.1 (R Core Team 2017) with a core i7 760 processor 2.8 GHz. The R and Nimble code for the models are available from the authors upon request. Gibbs sampling runs of 60,000 iterations with burn-in of 30,000 were used, and convergence criteria were established using the Gelman–Rubin statistic (Gelman and Rubin 1992) and by visual inspection. To address the issue of label switching over the MCMC iterations (Mengersen et al. (2011)), we used the *maximum a posteriori* estimate (MAP) to select one of the $k!$ modal regions and a distance based measure on the space of parameters to relabel parameters in proximity to this region (Celeux et al. 2000).

Model performance was assessed by comparing the classification accuracy and model selection criteria for different distributions within the family of SUNGH models (see Table 2). For classification accuracy, we report the adjusted rand index (ARI) (Hubert and Arabie 1985) which ranges from 0 (no match) to 1 (perfect match). We also report the EAIC and EBIC which are variations of the classical AIC and BIC criteria for use in a Bayesian setting (Carlin and

Louis 2011) (lower values indicate a better fit). In a mixture setting, it is also possible to compare the DIC values using one of the measures suggested by Celeux et al. (2006).

As to be expected, from Table 2 we can see quite clearly that the classification performance of the true model (SUNGH ($q = 2$)) is very good with an ARI of 0.87 and model choice criteria all appear to favor this model. A higher log-likelihood was found for the model SUNGH ($q = 3$) with a similar ARI score to the SUNGH ($q = 2$) model, but on other criteria this model was not favored due to the extra parameters involved. In applied settings, and where the true labels are unknown, a similar trade-off will be made between choosing more complex models with extra flexibility in the skewness matrix (higher q values) and relative improvement in model choice or goodness-of-fit measures. The performance of the SN and SNGH models is also to be expected given the relative lack of flexibility for the skewness parameter to accommodate the degree of skewness in all dimensions in this application.

5.2 Real applications

5.2.1 AIS example

In this example, we consider a dataset from the Australian Institute Sports (AIS) containing measures of physical activity for 202 athletes (102 male and 100 female) based on sex, red cell count, white cell count, hematocrit, hemoglobin, plasma ferritin concentration, body mass index, sum of skin folds, body fat percentage, lean body mass and finally height and weight of the athletes (Cook and Weisberg 1994). The data are available in the R package “sn” (see Azzalini 2015).

To assess the performance of the proposed SUNGH model, we use BMI and body fat percentage (Bfat) to classify male and female athletes. Figure 2a shows the observations for male (in black) and female (in red) athletes according to these two measures, suggesting a reasonably skewed distribution for both males and females with a particularly strong skewed and heavy-tailed distribution for male athletes. Figure 2b–d also shows the fitted contours and assigned labels for each observation for several of the models examined (SNGH, SN and SUNGH).

Table 3 presents the model choice criteria for the different models examined. The results suggest that the SUNGH has

Table 2 Results for simulated data example

Model	Log-likelihood (max)	EAIC	EBIC	DIC ₂	ARI
SN	−5936.9	12009.3	12198.8	11980.8	0.47
SNGH	−5692.1	11504.2	11693.6	11460.1	0.77
SUNGH ($q = 2$)	−5642.8	11458.9	11676.1	11444.3	0.87
SUNGH ($q = 3$)	− 5628.5	11487.5	11723.2	11514.1	0.86

The best values are indicated in bold

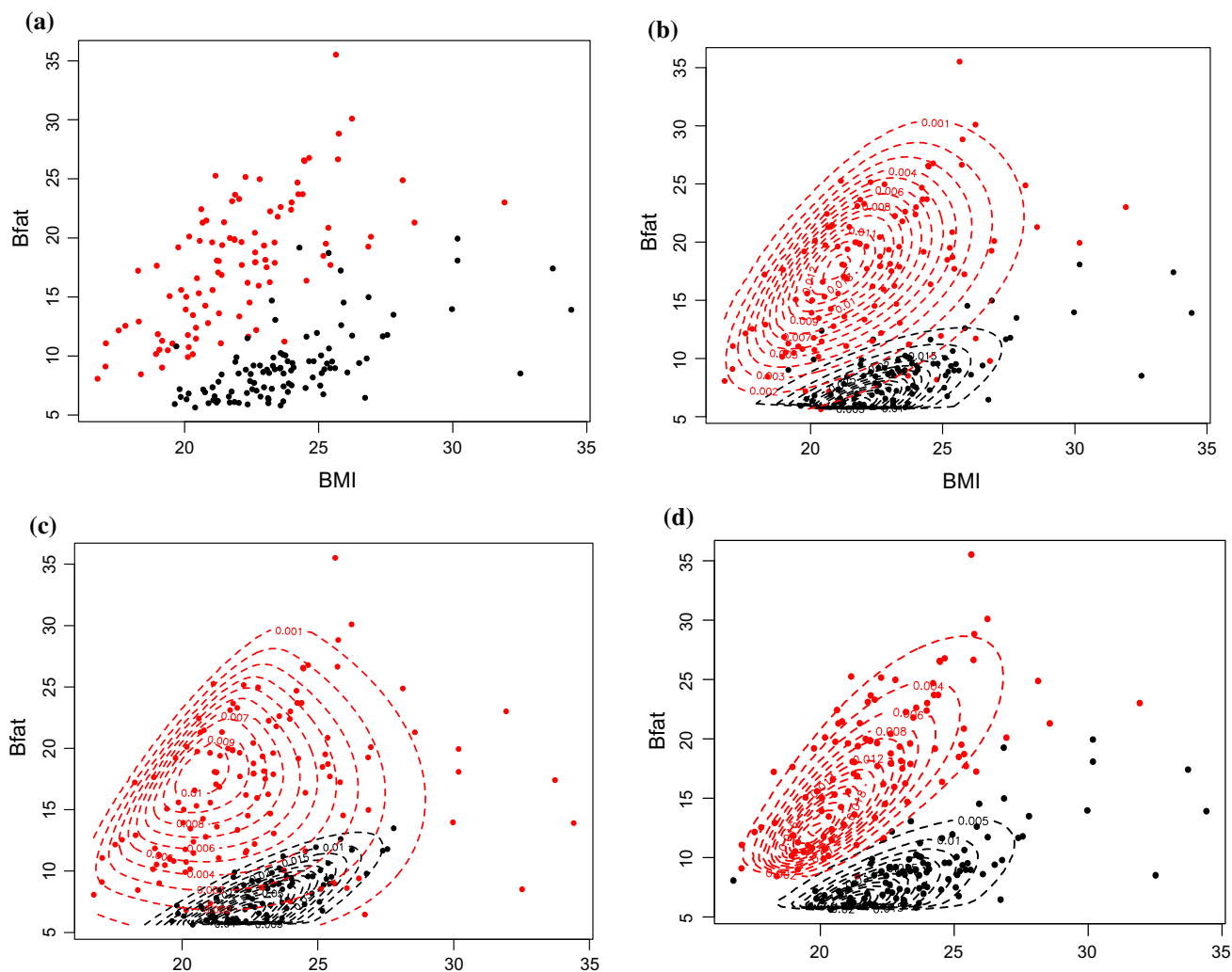


Fig. 2 Plots of results for AIS example: **a** true data, **b** SNGH, **c** SN, **d** SUNGH

Table 3 Model choice criteria for AIS data example

Model	Log-likelihood (max)	EAIC	EBIC	DIC ₂	ARI
SN	-1071.2	2200.6	2266.7	2176.3	0.52
SNGH	-1071.8	2205.2	2278.3	2186.8	0.64
(CFUST) SUT ($q = 2$)	-1074.3	2229.8	2309.2	2216.1	0.69
SUNGH ($q = 2$)	-1056.4	2197.8	2277.2	2186.8	0.79

The best values are indicated in bold; SN and SUT denote the restricted skew-normal and unrestricted skew-t distributions, respectively

the highest log-likelihood and the lowest EAIC, but the SN model has the lowest values for the EBIC and DIC₂ measures. However, the ARI for the SN model (= 0.52) is considerably lower than the SUNGH model (= 0.79) suggesting more support for the SUNGH model in terms of classification accuracy. We can also see these results reflected visually in Fig. 2 with the SUNGH (Fig. 2d) able to represent the skewed nature of the distribution for the two groups, particularly for the

male athletes. In contrast, the SN model (Fig. 2c) poorly represents the skewed distribution of the female athletes and the heavy-tailed nature of the distribution for the male athletes. As can be expected, the classification results for the SNGH model (Fig. 2b) are visually similar to the SUNGH; however, small differences (due to the reduced flexibility of the skewness parameter for the SNGH) can be observed which impact greatly on the classification accuracy (ARI = 0.64).

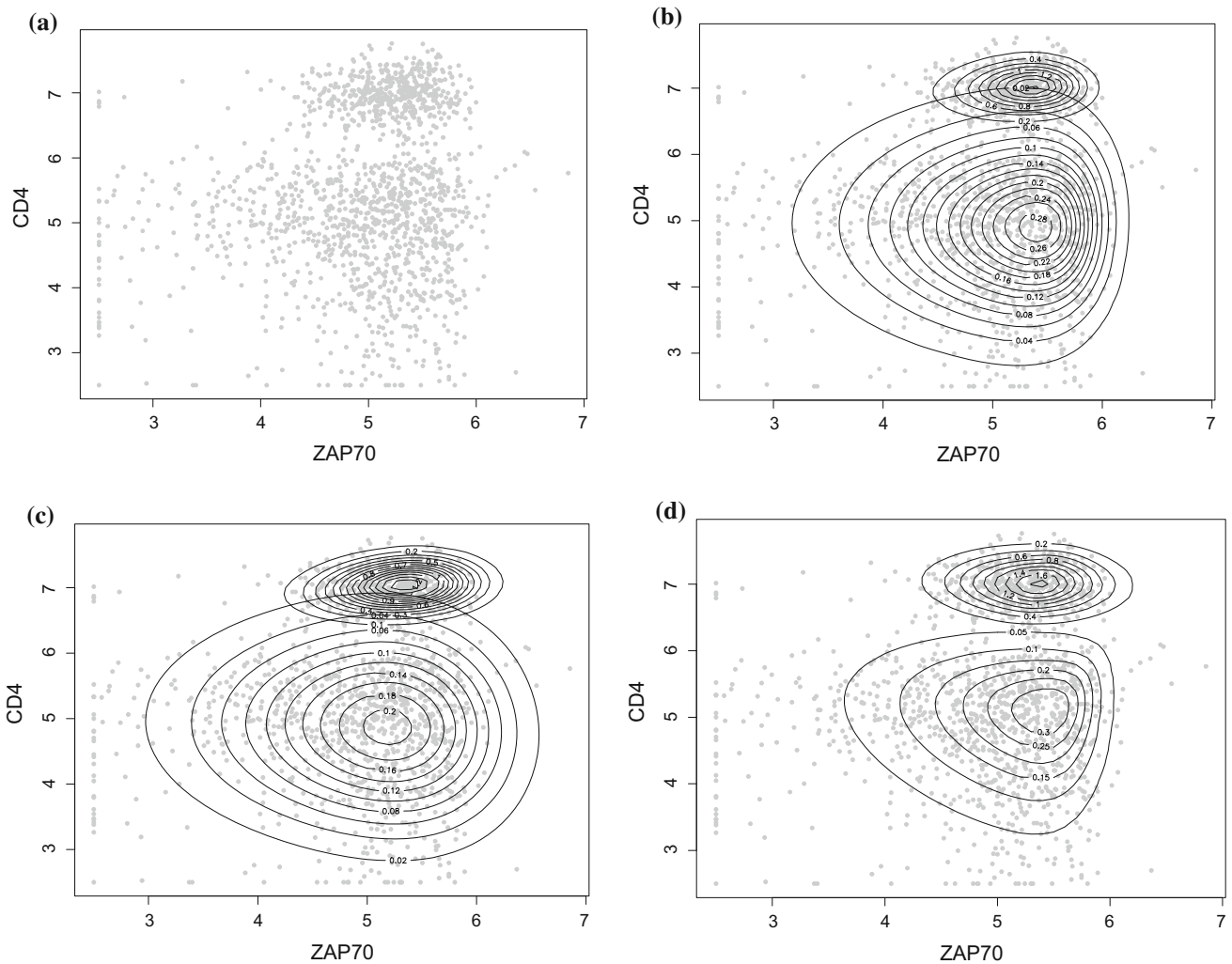


Fig. 3 Plots of observations and fitted contours for lymphoma example: **a** plot of observations for CD4 and ZAP70, **b** SNGH, **c** SN, **d** SUNGH

5.2.2 Lymphoma example

In another example, we examine a clustering problem for a lymphoma dataset analyzed by Lee and McLachlan (2013b). The data consist of a subset of data originally presented and collected by Maier et al. (2007). In Maier et al. (2007), blood samples from 30 subjects were stained with four fluorophore-labeled antibodies against CD4, CD45RA, SLP 76(pY 128) and ZAP 70 (pY 292) before and after an anti-CD3 stimulation. To illustrate the performance of different distributions within the SUNGH family, we will look at clustering a subset of the data containing the variables CD4 and ZAP70 (Fig. 3), which appear to be bimodal and display an asymmetric pattern. In particular, the largest mode appears to show both strong correlation between the two variables and substantial skewness in both dimensions.

From Fig. 3, we can see a clear difference between the SNGH and SUNGH models with the latter providing a closer fit to the two groups visible in the data. This is further supported by the model choice criteria with all three measures (EAIC, EBIC and DIC_2) favoring the SUNGH model (Table 4) compared to SNGH. The SUNGH model is also preferred over the SN model (Fig. 3c) with the latter model not appearing to fit or represent the larger component in the data. Overall, the SUNGH model appears to fit the two groups in these data quite well with the lowest values for two of the three model choice measures (EAIC and EBIC). Using the DIC_2 criteria, the lowest value appears to be for the SUT model, so there is some support for this in terms of model choice. However, the difference between DIC_2 values for this model and SUNGH is not great (SUNGH = 7201.3; SUT = 7198.5), suggesting little difference in terms of this measure.

Table 4 Model choice criteria for lymphoma example

Model	Log-likelihood (max)	EAIC	EBIC	DIC ₂
SN	−3657.0	7433.6	7539.8	7473.1
(CFUSN) SUN ($q = 2$)	−3567.8	7507.1	7634.6	7782.6
ST	−3577.6	7240.7	7347.0	7246.2
(CFUST) SUT ($q = 2$)	−3404.5	7051.8	7179.3	7198.5
SUNGH ($q = 2$)	−3396.4	7045.0	7172.5	7201.3
SNGH	−3578.9	7252.4	7358.7	7266.9

The best values are indicated in bold; SN and SUN denote the restricted and unrestricted skew-normal distributions, respectively; ST and SUT denote the restricted and unrestricted skew-t distributions, respectively

6 Conclusion

We have proposed a flexible family of unrestricted skew-normal generalized hyperbolic (SUNGH) distributions for application in clustering problems which are capable of representing distributions of asymmetric and heavy-tailed forms. The family contains several other well-known symmetric and asymmetric families of distributions such as scale mixtures of the skew-normal family (SMSN) as special cases. Various properties of the SUNGH family are well defined, and estimation of the parameters is relatively straightforward in a Bayesian framework with most of the Gibbs sampling updates available in closed form. Assessments of the performance of the proposed model on simulated and real data suggest that the family provides a considerable degree of freedom and flexibility in modeling data of varying tail behavior and directional shape. As this family of distributions and the parameterization we have adopted preserves several important propositions (e.g., closed under linear combinations), the SUNGH family can be used in a variety of other statistical models (e.g., linear multilevel/mixed models and regression).

Acknowledgements The authors would like to thank the coordinating editor and anonymous reviewers for their suggestions, corrections and encouragement, which helped us to improve earlier versions of the manuscript.

Appendix

A.1. Proof of Propositions 1 to 6

In this appendix, we prove Propositions 1 to 6.

Proof of Proposition 1 By considering (7),

$$\begin{aligned}
 \text{(a): } \mu_Y &= E[Y] = E_U[E_{Y|U}[Y]] \\
 &= E_U[E_{Y|U}[\mu + \kappa(u)^{1/2} \Lambda |W_0| + \kappa(u)^{1/2} \Sigma^{1/2} W_1]] \\
 &= E_U[\mu + \kappa(U)^{1/2} \Lambda E|W_0| + \kappa(U)^{1/2} \Sigma^{1/2} E(W_1)] \\
 &= \mu + k_1 \Lambda E|W_0|,
 \end{aligned}$$

$$\begin{aligned}
 \text{(b): } \text{Var}[Y] &= E[(Y - \mu_Y)(Y - \mu_Y)^T] \\
 &= E_U[E_{Y|U}[(Y - \mu_Y)(Y - \mu_Y)^T]] \\
 &= E_U[E_{Y|U}[(\Lambda[\kappa(u)^{1/2} |W_0| - k_1 E|W_0|] \\
 &\quad + \kappa(u)^{1/2} \Sigma^{1/2} W_1)([\kappa(u)^{1/2} |W_0| \\
 &\quad - k_1 E|W_0|] \Lambda^T + \kappa(u)^{1/2} W_1^T \Sigma^{1/2})]] \\
 &= \Lambda E_U[E_{Y|U}\{[\kappa(U)^{1/2} |W_0| \\
 &\quad - k_1 E|W_0|][\kappa(U)^{1/2} |W_0| \\
 &\quad - k_1 E|W_0|]\} \Lambda^T \\
 &\quad + E_U[\kappa(U) \Sigma^{1/2} E_{Y|U}[W_1 W_0^T] \Sigma^{1/2}]] \\
 &= \Lambda [k_2 E|W_0| |W_0|^T - k_1^2 E|W_0| E|W_0|^T] \Lambda^T + k_2 \Sigma \\
 &= \Lambda \left[k_2 \left(\frac{2}{\pi} \mathbf{1}_q \mathbf{1}_q^T + \left(1 - \frac{2}{\pi} \right) I_q \right) - \frac{2}{\pi} k_1^2 \mathbf{1}_q \mathbf{1}_q^T \right] \\
 &\quad \Lambda^T + k_2 \Sigma \\
 &= k_2 \Psi + \Lambda \left[(k_2 - k_1^2) \frac{2}{\pi} \mathbf{1}_q \mathbf{1}_q^T - \frac{2}{\pi} k_2 I_q \right] \Lambda^T \quad \square
 \end{aligned}$$

Proof of Proposition 2 By considering the stochastic representation (7) and the fact that W_0 (and so W) are uncorrelated, this subject proved. In the case of $\Lambda^* = (\Lambda_{p \times q} \mathbf{0}_{p \times m})$, relation (7) for $Y \sim \text{SUNGH}_{p,q+m}(\mu, \Sigma, \Lambda^*, \varpi)$ is equivalent to $Y = \mu + \Lambda^* W + \kappa(U)^{1/2} \Sigma^{1/2} W_1 = \mu + \Lambda W^{(1)} + \kappa(U)^{1/2} \Sigma^{1/2} W_1$, where $W^{(1)}$ is the first q components of W , and in the case of $\Lambda^* = (\mathbf{0}_{p \times m} \Lambda_{p \times q})$, relation (7) for $Y \sim \text{SUNGH}_{p,q+m}(\mu, \Sigma, \Lambda^*, \varpi)$ is equivalent to $Y = \mu + \Lambda^* W + \kappa(U)^{1/2} \Sigma^{1/2} W_1 = \mu + \Lambda W^{(2)} + \kappa(U)^{1/2} \Sigma^{1/2} W_1$, where $W^{(2)}$ is the last q components of W \square

Proof of Proposition 3 By considering the stochastic representation (7), we have that $b + BY = b + B\mu + B\Lambda W + \kappa(U)^{1/2} (B\Sigma B^T)^{1/2} W_1$ \square

Proof of Proposition 4 By considering Proposition 3, with $b = \mathbf{0}$ and the matrix B in the form of $(I_{p_1} \mathbf{0}_{p_1 \times p_2})$ or $(\mathbf{0}_{p_2 \times p_1} I_{p_2})$, respectively, this subject proved \square

Proof of Proposition 5 Since $Y = (Y_1^T, Y_2^T)^T$, from part (b) of the Proposition 1, we have $\text{Var}[Y] = (\text{Cov}(Y_i, Y_j))_{i,j=1,2}$

$= \left(\Sigma_{ij} + \Lambda_i \left[(k_2 - k_1^2) \frac{2}{\pi} \mathbf{1}_q \mathbf{1}_q^\top - \frac{2}{\pi} k_2 I_q \right] \Lambda_j^\top \right)$. Thus, if $\Sigma_{12} = \mathbf{0}$, then $\text{Cov}(Y_1, Y_2) = \Lambda_1 \left[(k_2 - k_1^2) \frac{2}{\pi} \mathbf{1}_q \mathbf{1}_q^\top - \frac{2}{\pi} k_2 I_q \right] \Lambda_2^\top$, thus following that each of the conditions $\Lambda_1 = \mathbf{0}$ or $\Lambda_2 = \mathbf{0}$ leads to $\text{Cov}(Y_1, Y_2) = \mathbf{0}$ \square

Proof of Proposition 6 The first part follows by applying Proposition 2 in the Proposition 4. For the proof of the second result, note from the proof of Proposition 5 that

$$\text{Cov}(Y_1, Y_2) = \left(\Lambda_{11} \mathbf{0}_{p_1 \times q_2} \right)_{p_1 \times q} \left[\left((k_2 - k_1^2) \frac{2}{\pi} \mathbf{1}_q \mathbf{1}_q^\top - \frac{2}{\pi} k_2 I_q \right) \right]_{q \times q} \left(\mathbf{0}_{p_2 \times q_1} \Lambda_{22} \right)_{q \times p_2}^\top.$$

Thus, using the partitions $I_q = \text{diag}(I_{q_1}, I_{q_2})$ and $\mathbf{1}_q = (\mathbf{1}_{q_1}^\top, \mathbf{1}_{q_2}^\top)^\top$ we obtain the proof \square

A.2. Matrix variate priors for skewness matrix

Considering the matrix variate priors in the form of $\Lambda_k \sim MN_{p,q}(N_k, S_k, F_k), k = 1, \dots, K$, where MN denotes the matrix normal distributions, this leads to the following posteriors instead of (19) as follows:

$\text{vec}(\Lambda_k) | \Theta_{(-\Lambda_k)}, y, u, w, z_i = k \sim N_{pq}(\mu, \Sigma); k = 1, \dots, K$, where

$$\mu = \Sigma \left[(S_k \otimes F_k)^{-1} \text{vec}(N_k) + \sum_{B_k} \kappa(u_{ik})^{-1} (M_{ik}^\top \otimes \Sigma_k^{-1}) \right], \tag{19a}$$

$$\Sigma = \left[(S_k \otimes F_k)^{-1} + \sum_{B_k} \kappa(u_{ik})^{-1} (\Sigma_k^{-1} \otimes L_{ik}) \right]^{-1},$$

where $L_{ik} = w_{ik} w_{ik}^\top$ and $M_{ik} = (y_i - \mu_k) w_{ik}^\top$, for which \otimes denotes the Kronecker product and vec denotes the vectorization of a matrix (a linear transformation which converts the matrix into a column vector).

Using these forms for the Gibbs updates may improve mixing and convergence to a stationary distribution. However, they involve the use of matrix variate distributions for which users may not be familiar; hence, a simpler (computational) update is provided in the main text.

References

Andrews, D.R., Mallows, C.L.: Scale mixture of normal distribution. *J. Roy. Stat. Soc. B* **36**, 99–102 (1974)
 Arellano-Valle, R.B., Azzalini, A.: On the unification of families of skew-normal distributions. *Scand. J. Stat.* **33**, 561–574 (2006)
 Arellano-Valle, R.B., Genton, M.G.: On fundamental skew distributions. *J. Multivar. Anal.* **96**, 93–116 (2005)

Arellano-Valle, R.B., Genton, M.G.: Multivariate unified skew-elliptical distributions. *Chil. J. Stat.* **2**, 17–34 (2010)
 Arellano-Valle, R.B., Branco, M.D., Genton, M.G.: A unified view on skewed distributions arising from selections. *Can. J. Stat.* **34**, 581–601 (2006)
 Arellano-Valle, R.B., Bolfarine, H., Lachos, G.H.: Bayesian inference for skew-normal linear mixed model. *J. Appl. Stat.* **33**, 561–574 (2007)
 Azzalini, A.: Package ‘sn’. <http://azzalini.stat.unipd.it/SN> (2015). Accessed 13 May 2017
 Azzalini, A., with the collaboration of Capitanio, A.: *The Skew-Normal and Related Families*. IMS Monographs Series. Cambridge University Press (2014)
 Barndorff-Nielsen, O.: Hyperbolic distributions and distributions on hyperbolae. *Scand. J. Stat.* **5**, 151–157 (1978)
 Barndorff-Nielsen, O., Blaesild, P.: Hyperbolic distributions. In: Kotz, S., Johnson, N.L., Read, C. (eds.) *Encyclopedia of Statistical Sciences*, vol. 3. Wiley, New York (1980)
 Barndorff-Nielsen, O., Halgreen, C.: Infinite divisibility of the hyperbolic and generalized inverse Gaussian distributions. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete* **38**, 309–311 (1977)
 Basso, R.M., Lachos, V.H., Cabral, C.R.B., Ghosh, P.: Robust mixture modeling based on the scale mixtures of skew-normal distributions. *Comput. Stat. Data Anal.* **54**, 2926–2941 (2010)
 Böhning, D.: *Computer-Assisted Analysis of Mixtures and Applications. Meta-Analysis, Disease Mapping and Others*. Chapman & Hall, Boca Raton (2000)
 Branco, M.D., Dey, D.K.: A general class of multivariate skew-elliptical distributions. *J. Multivar. Anal.* **79**, 99–113 (2001)
 Browne, R.P., McNicholas, P.D.: A mixture of generalized hyperbolic distributions. *Can. J. Stat.* **43**(2), 176–198 (2015)
 Carlin, B.P., Louis, T.A.: *Bayesian Methods for Data Analysis*. CRC Press, Boca Raton (2011)
 Celeux, G., Hurn, M., Robert, C.P.: Computational and inferential difficulties with mixture posterior distributions. *J. Am. Stat. Assoc.* **95**, 957–970 (2000)
 Celeux, G., Forbes, F., Robert, C.P., Titterton, D.M.: Deviance information criteria for missing data models. *Bayesian Anal.* **1**, 651–674 (2006)
 Chhikara, R.S., Folks, J.L.: *The Inverse Gaussian Distribution*. Marcel Dekker, New York (1989)
 Cook, R.D., Weisberg, S.: *An Introduction to Regression Graphics*. Wiley, New York (1994)
 Forbes, F., Wraith, D.: A new family of multivariate heavy-tailed distributions with variable marginal amounts of tail weight: application to robust clustering. *Stat. Comput.* **24**(6), 971–984 (2014)
 Franczak, B.C., Browne, R.P., McNicholas, P.D.: Mixtures of shifted asymmetric laplace distributions. *IEEE Trans. Pattern Anal. Mach. Intell.* **36**(6), 1149–1157 (2014)
 Frühwirth-Schnatter, S.: *Finite Mixture and Markov Switching Models*. Springer Series in Statistics. Springer, Berlin (2006)
 Frühwirth-Schnatter, S., Pyne, S.: Bayesian inference for finite mixtures of skew-normal and skew-t distributions. *Biostatistics* **11**(2), 317–336 (2010)
 Gelman, A., Rubin, D.B.: Inference from iterative simulation using multiple sequences. *Stat. Sci.* **7**, 457–511 (1992)
 Genton, M.G.: *Skew-Elliptical Distributions and Their Applications: A Journey Beyond Normality*. Chapman & Hall, Boca Raton (2004)
 Good, I.J.: The population frequencies of species and the estimation of population parameters. *Biometrika* **40**, 237–260 (1953)
 Hogan, J.W., Laird, N.M.: Mixture models for the joint distribution of repeated measures and event times. *Stat. Med.* **16**, 239–258 (1997)
 Holzmann, H., Munk, A., Gneiting, T.: Identifiability of finite mixtures of elliptical distributions. *Scand. J. Stat.* **33**(4), 753–763 (2006)

- Hubert, L., Arabie, P.: Comparing partitions. *J. Classif.* **2**, 193–218 (1985)
- Johnson, N.L., Kotz, S., Balakrishnan, N.: *Continuous Univariate Distributions*, vol. 1. Wiley, New York (1994)
- Jørgensen, B.: *Statistical Properties of the Generalized Inverse Gaussian distribution*. Springer, New York (1982)
- Karlis, D., Santourian, A.: Model-based clustering with non-elliptically contoured distributions. *Stat. Comput.* **19**(1), 73–83 (2009)
- Lachos, V.H., Bolfarine, H., Arellano-Valle, R.B.: Likelihood-based inference for multivariate skew-normal regression models. *Commun. Stat. Theory Methods* **36**(9), 1769–1786 (2007)
- Lachos, V.H., Ghosh, P., Arellano-Valle, R.B.: Likelihood based inference for skew-normal independent linear mixed models. *Stat. Sin.* **20**, 303–322 (2010)
- Lee, S.X., McLachlan, G.J.: Model-based clustering and classification with non-normal mixture distributions. *Stat. Methods Appl.* **22**(4), 427–454 (2013a)
- Lee, S.X., McLachlan, G.J.: On mixtures of skew normal and skew t distributions. *Adv. Data Anal. Classif.* **7**(3), 241–266 (2013b)
- Lee, S.X., McLachlan, G.J.: Finite mixtures of multivariate skew t distributions: some recent and new results. *Stat. Comput.* **24**, 181–202 (2014)
- Lee, S.X., McLachlan, G.J.: Finite mixtures of canonical fundamental skew t-distributions: the unification of the restricted and unrestricted skew t-mixture models. *Stat. Comput.* **26**, 573–589 (2016)
- Lin, T.I.: Maximum likelihood estimation for multivariate skew normal mixture models. *J. Multivar. Anal.* **100**(2), 257–265 (2009)
- Lin, T.I.: Robust mixture modeling using multivariate skew t distributions. *Stat. Comput.* **20**(3), 343–356 (2010)
- Lin, T.I., Lee, J.C., Yen, S.Y.: Finite mixture modeling using the skew-normal distribution. *Stat. Sin.* **17**(b), 909–927 (2007)
- Lin, T.I., Ho, H.J., Chen, C.L.: Analysis of multivariate skew normal models with incomplete data. *J. Multivar. Anal.* **100**(10), 2337–2351 (2009)
- Maier, L.M., Anderson, D.E., De Jager, P.L., Wicker, L.S., Hafler, D.A.: Allelic variant in CTLA4 alters t cell phosphorylation patterns. *Proc. Natl. Acad. Sci. USA* **104**, 18607–18612 (2007)
- Maleki, M., Arellano-Valle, R.B.: Maximum a-posteriori estimation of autoregressive processes based on finite mixtures of scale-mixtures of skew-normal distributions. *J. Stat. Comput. Simul.* **87**(6), 1061–1083 (2017)
- McLachlan, G.J., Peel, D.: *Finite Mixture Models*. Wiley, Chichester (2000)
- McNeil, A.J., Frey, R., Embrechts, P.: *Quantitative Risk Management: Concepts, Techniques and Tools*. Princeton University Press, Princeton (2005)
- Mengersen, K., Robert, C., Titterton, D.M.: *Mixtures: Estimation and Applications*. Wiley, Chichester (2011)
- Morris, K., McNicholas, P.D., Punzo, A., Browne, R.P.: Robust Asymmetric Clustering. ArXiv e-print [arxiv:1402.6744](https://arxiv.org/abs/1402.6744) (2014)
- Pyne, S., Hu, X., Wang, K., Rossin, E., Lin, T.I., Maier, L.M., Baecher-Allan, C., McLachlan, G.J., Tamayo, P., Hafler, D.A., De Jager, P.L., Mesirov, J.P.: Automated high-dimensional flow cytometric data analysis. *Proc. Natl. Acad. Sci.* **106**(21), 8519–8524 (2009)
- R Core Team.: *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/> (2017). Accessed 20 June 2017
- Sahu, S.K., Dey, D.K., Branco, M.D.: A new class of multivariate skew distributions with applications to Bayesian regression models. *Can. J. Stat.* **31**(2), 129–150 (2003)
- Seshadri, V.: *The Inverse Gaussian Distribution: A Case Study in Exponential Families*. Oxford University Press, New York (1993)
- Teicher, H.: Identifiability of finite mixtures. *Ann. Math. Stat.* **34**(4), 1265–1269 (1963)
- Vilca, F., Balakrishnan, N., Zeller, C.B.: Multivariate skew-normal generalized hyperbolic distribution and its properties. *J. Multivar. Anal.* **128**, 73–85 (2014)
- Vrbik, I., McNicholas, P.D.: Analytic calculations for the EM algorithm for multivariate skew-t mixture models. *Stat. Probab. Lett.* **82**(6), 1169–1174 (2012)
- Wang, H.X., Zhang, Q.B., Luo, B., Wei, S.: Robust mixture modelling using multivariate t-distribution with missing information. *Pattern Recogn. Lett.* **25**(6), 701–710 (2004)
- Wang, K., Ng, S.K., McLachlan, G.J.: Multivariate skew t mixture models: applications to fluorescence-activated cell sorting data. In: *Digital Image Computing: Techniques and Applications*, Los Alamitos, California, pp. 526–531. IEEE (2009)
- Wraith, D., Forbes, F.: Location and scale mixtures of Gaussians with flexible tail behaviour: properties, inference and application to multivariate clustering. *Comput. Stat. Data Anal.* **90**(Oct.), 61–73 (2015)