

PANKAJ BHAGWAT<sup>a</sup> & ÉRIC MARCHAND

*a Université de Sherbrooke, Département de mathématiques, Sherbrooke Qc, CANADA, J1K 2R1  
(e-mails: pankaj.uttam.bhagwat@usherbrooke.ca; eric.marchand@usherbrooke.ca)*

## ABSTRACT

We study frequentist risk properties of predictive density estimators for mean mixtures of multivariate normal distributions, involving an unknown location parameter  $\theta \in \mathbb{R}^d$ , and which include multivariate skew normal distributions. We provide explicit representations for Bayesian posterior and predictive densities, including the benchmark minimum risk equivariant (MRE) density, which is minimax and generalized Bayes with respect to an improper uniform density for  $\theta$ . For four dimensions or more, we obtain Bayesian densities that improve uniformly on the MRE density under Kullback-Leibler loss. We also provide plug-in type improvements, investigate implications for certain type of parametric restrictions on  $\theta$ , and illustrate and comment the findings based on numerical evaluations.

*Keywords and phrases:* Bayes predictive density; Dominance; Kullback-Leibler loss; Minimax; Minimum risk equivariant; Mean mixtures; Multivariate Normal, Skew-normal distribution.

## 1. Introduction

The findings of this paper relate to predictive density estimation for mean mixture of normal distributions. The modelling of data via mixing multivariate normal distributions has found many applications and lead to methodological challenges for statistical inference. These include finite mixtures, as well as continuous mixing on the mean and/or the variance. Whereas, scale or variance mixtures of multivariate normal distributions compose a quite interesting subclass of spherically symmetric distributions, modelling asymmetry requires mixing on the mean and prominent examples are generated via variance-mean mixtures (e.g., [7]), as well as mean mixtures of multivariate normal distributions (e.g., [1, 5]) and references therein). Moreover, such mean mixtures, which are the subject of study here, generate or are connected to multivariate skew-normal distributions (e.g., [6]) which have garnered much interest over the years.

The development of shrinkage estimation techniques, namely since Stein's inadmissibility finding ([26]) concerning the maximum likelihood or best location equivariant estimator under squared error loss in three dimensions or more, has had a profound impact on statistical theory, thinking, methods, and practice. Examples include developments on sparsity and regularization methods, empirical Bayes estimation, adaptive inference, small area estimation, non-parametric function estimation, and predictive density estimation. Cast in a decision-theoretic framework, Stein's original result has been expanded in many diverse ways, namely to other distributions or probability models, and namely for spherically symmetric and elliptically symmetric distributions (see for instance, [11]). There have been fewer findings for multivariate skew-normal or mean mixtures of normal distributions, but the recent work of Kubokawa et al. [18] establishes point estimation minimax

---

<sup>1</sup>January 25, 2022

improvements of the best location equivariant estimator under quadratic loss, when the dimension of the location parameter is greater than or equal to four, and with underlying known perturbation parameter.

Predictive density estimation has garnered much interest over the past twenty years or so, and addresses fundamental issues in statistical predictive analysis. Decision-theoretic links between shrinkage point estimation and shrinkage predictive density estimation for normal models have surfaced (e.g., [14], [13]) and stimulated much activity (see for instance [12]), including findings for restricted parameter spaces (e.g., [10]). The main objective of this work is thus to explore the problem of predictive density estimation for mean mixtures of normal (MMN) distributions. A secondary objective is to provide novel representations for Bayesian posterior distributions and predictive densities, which have been found to be lacking in the literature.

Following early findings of Komaki (e.g., [14]) on the predictive density estimation problem for multivariate normal models under Kullback-Leibler loss, George, Liang and Xu in [13] exhibited further parallels with the point estimation problem for normal distribution under quadratic loss. They provide sufficient conditions on marginal distributions and prior distributions to get improved shrinkage predictive density estimators when the dimension is greater than or equal to three. Thus, motivated by these connections, it is interesting to investigate whether such shrinkage plays any role in the predictive density estimation problem for mean-mixture of multivariate normal models and we focus on frequentist risk efficiency of predictive density estimators under Kullback-Leibler loss. Our contribution here consists in identifying classes of plug-in type predictive densities and of Bayes predictive densities which are minimax and dominate the benchmark minimum equivariant estimator (MRE) for the case when the dimension of the location parameter is greater than or equal to four.

The organization of this manuscript is as follows. Section 2.1 contains several introductory definitions, properties and examples of MMN models, including a useful canonical form which subdivides the MMN distributed vector into  $d$  independent components, one of which a univariate MMN distribution and the others as normal distributions. Section 2.2 focuses on the predictive estimation framework with a KL loss decomposition, and an initial representation for the MRE density accompanied by various examples. Section 2.3 expands on the calculation of minimax risk and a representation in terms of the entropy of a univariate distribution. Section 3 is devoted to Bayesian posterior and predictive analysis with several novel representations. Sections 4.1 and 4.2, namely Theorems 4.4, Theorem 4.5 and Corollary 4.1, contain the main dominance findings, with plug-in type and Bayesian improvements. In both cases, the main technique employed rests upon the canonical transformation presented in Section 2.1 and permits to split up the KL risk as the addition of two parts, one of which can be operated on using known normal model prediction analysis findings. Section 4.3 deals with parametric restrictions and further applications of Theorems 4.4 and 4.5. Finally, some further illustrations are provided in Section 5.

## 2. Preliminary results and definitions

Here are some details, properties and definitions on mean mixture of normal distributions, its canonical form, and predictive density estimation. In the following, we will denote  $\phi_d(z; \Sigma)$  the probability density function (pdf) of a  $N_d(0, \Sigma)$  distribution evaluated at  $z \in \mathbb{R}^d$  and for positive definite  $\Sigma$ . When  $\Sigma = I_d$ , we may simplify the writing to  $\phi_d(z)$ , and then for  $d = 1$  to  $\phi(z)$ . We

will denote  $\Phi$  the cdf of a  $N(0, 1)$  distribution.

## 2.1. The model

The distributions of interest are mean-mixtures of multivariate normal distributions, both for our observables and densities to be estimated by a predictive density estimator. Such distributions connect to multivariate skew-normal distributions and have been the object of interest in recent work with studies of stochastic properties (e.g., [1], [5]), and shrinkage estimation about its location parameter ([18]).

**Definition 2.1.** *A random vector  $X \in \mathbb{R}^d$  is said to have a mean-mixture of normal distributions (MMN), denoted as  $X \sim MMN_d(\theta, a, \Sigma, \mathcal{L})$ , if it admits the representation*

$$X|V = v \sim N_d(\theta + va, \Sigma), V \sim \mathcal{L}, \quad (2.1)$$

where  $\theta \in \mathbb{R}^d$  is a location parameter,  $a \in \mathbb{R}^d - \{0\}$  is a known perturbation vector,  $\Sigma$  is a known positive definite covariance matrix, and  $V$  is a scalar random variable with cdf  $\mathcal{L}$ .

Alternatively, the random vector  $X$  has stochastic representation

$$X = \theta + \Sigma^{1/2}Z + Va, \quad (2.2)$$

where  $Z \sim N_d(0, I_d)$  and  $V \sim \mathcal{L}$  on  $\mathbb{R}$ , and its probability density function can be expressed as:

$$\begin{aligned} p(x|\theta) &= \mathbb{E}^V \{ \phi_d(x - \theta - Va, \Sigma) \} \\ &= \phi_d(x - \theta, \Sigma) \mathbb{E}^V \left( e^{-\frac{V^2}{2} a^T \Sigma^{-1} a} e^{V(x-\theta)^T \Sigma^{-1} a} \right). \end{aligned} \quad (2.3)$$

Thus, we note that the density function of MMN random vector can be decomposed in two parts: one symmetrical density  $\phi_d(\cdot)$  and the other part which is a function of the projection of  $(x - \theta)$  in the direction of  $\Sigma^{-1}a$ . Moreover, this construction isolates the asymmetry in the direction  $\Sigma^{-1}a$  and the scale is controlled by the random variable  $V$ .

**Remark 2.1.** *It is easy to see that the family of MMN distributions is closed under linear combinations of independent components. Specifically, if  $X_i|\theta \sim MMN_d(\theta, a, \Sigma_i, \mathcal{L}_i)$ ,  $i = 1, \dots, n$ , are independently distributed, then  $\sum_{i=1}^n b_i X_i|\theta \sim MMN_d((\sum_{i=1}^n b_i)\theta, a, \sum_{i=1}^n b_i^2 \Sigma_i, \mathcal{L}_0)$  with  $\mathcal{L}_0$  the cdf of the mixing variable  $V_0 = \sum_{i=1}^n b_i V_i$ . Namely, for the identically distributed case with  $\Sigma_i = \Sigma$  and the sample mean with  $b_i = 1/n$ , we obtain that*

$$\bar{X}|\theta \sim MMN_d(\theta, a, \Sigma/n, \mathcal{L}_0), \text{ with } \mathcal{L}_0 \text{ the cdf of } \bar{V}.$$

*It thus follows, as observed in [18], that findings applicable for a single MMN distributed observable  $X$  can be extended to the random sample case.*

We now turn our attention to a fundamental decomposition, or canonical form, ([1]) for MMN distributions which will be most useful.

**Lemma 2.1.** *For a random vector  $X \sim MMN_d(\theta, a, \Sigma, \mathcal{L})$  as in (2.1), there exists an orthogonal matrix  $H$  such that the first row of  $H$  is proportional to  $a^\top \Sigma^{-1/2}$  and  $Z = H\Sigma^{-1/2}X$  has a  $MMN_d(H\Sigma^{-1/2}\theta, a_0, I_d, \mathcal{L})$  distribution with  $a_0 = (\sqrt{a^\top \Sigma^{-1}a}, 0, \dots, 0)^\top$ .*

Such a  $Z$  may be referred to as a canonical form and is comprised of  $d$  independent components. Moreover  $Z - H\Sigma^{-1/2}\theta$  has  $d - 1$  components which are  $N(0, 1)$  distributed and another distributed as  $MMN_1(0, a_0, 1, \mathcal{L})$ . Such a canonical form construction is not unique and depends on the choice of  $H$ .

As already mentioned, the family of MMN distributions contains many interesting distributions and we refer to the above-mentioned references for various properties. We expand here with illustrations, which will also inform us for our predictive density problem and related Bayesian posterior analysis. A prominent example is the multivariate skew normal distribution due to Azzalini and Dalla Valle [6]. If we consider  $V \sim TRN(0, 1)$ , the standard truncated normal distribution on  $R_+$  in (2.1), we get the multivariate skew-normal family of distributions with densities

$$p(x|\theta) = 2\phi_d(x - \theta; \Sigma + aa^T) \Phi\left(\frac{(x - \theta)^\top \Sigma^{-1}a}{\sqrt{1 + a^\top \Sigma^{-1}a}}\right). \quad (2.4)$$

We denote this as  $X \sim SN_d(\theta, a, \Sigma)$ . Here, we note that  $V \sim \sqrt{\chi_1^2}$ , i.e. the square root of a Chi-square distribution with  $k = 1$  degrees of freedom. Various other choices of the mixing density have appeared in the literature (e.g., [5]), namely cases where  $V \sim \sqrt{\chi_k^2}$  or  $V$  is Gamma distributed. Here is a general result containing such cases as well as many others.

**Theorem 2.1.** *For a mixing density of the form*

$$\ell(v) = h(v) e^{-vc_2} e^{-\frac{v^2}{2}c_1} \mathbb{I}_{(0, \infty)}(v), \quad (2.5)$$

with  $c_1 > 0, c_2 \in \mathbb{R}$  or  $c_1 = 0, c_2 \geq 0$ , the corresponding pdf of  $X$  in (2.1) is given by

$$p(x|\theta) = \frac{1}{c_1} \phi_d(x - \theta, \Sigma) \frac{\mathbb{E}\left[h\left\{\frac{1}{c_1}\left(Z + \frac{c_2'}{c_1}\right)\right\} \middle| Z + \frac{c_2'}{c_1} \geq 0\right]}{R\left(\frac{c_2'}{c_1}\right)}, \quad (2.6)$$

with  $Z \sim N(0, 1)$ ,  $c_1' = (c_1 + a^\top \Sigma^{-1}a)^{1/2}$ ,  $c_2' = (x - \theta)^\top \Sigma^{-1}a - c_2$ , and  $R(\cdot)$  the reverse Mill's ratio given by  $R(t) = \phi(t)/\Phi(t), t \in \mathbb{R}$ .

**Proof.** *The result follows from (2.3) as*

$$\begin{aligned} \mathbb{E}^V\left(e^{-\frac{v^2}{2}a^\top \Sigma^{-1}a} e^{V(x-\mu)^\top \Sigma^{-1}a}\right) &= \int_0^\infty e^{-\frac{v^2}{2}(c_1')^2} e^{vc_2'} h(v) dv \\ &= \frac{\sqrt{2\pi}}{c_1'} e^{\frac{c_2'^2}{2c_1'^2}} \int_0^\infty h(v) \frac{c_1'}{\sqrt{2\pi}} e^{-\frac{c_1'^2}{2}\left(v - \frac{c_2'}{c_1'}\right)^2} dv \\ &= \frac{\sqrt{2\pi}}{c_1'} e^{\frac{c_2'^2}{2c_1'^2}} \mathbb{E}\left\{h\left(\frac{Z}{c_1'} + \frac{c_2'}{c_1'}\right) \middle| Z + \frac{c_2'}{c_1'} \geq 0\right\} \Phi\left(\frac{c_2'}{c_1'}\right). \quad \square \end{aligned}$$

We point out that the above Theorem applies for  $c_1 = c_2 = 0$  and thus covers all absolutely continuous distributions on  $\mathbb{R}_+$ . Here are nevertheless specific examples of Theorem 2.1 and model density (2.6).

**Example 2.1.** (A) *Gamma mixing with  $h(v) = \frac{v^{\alpha-1}}{\Gamma(\alpha)\beta^\alpha}$ . Theorem 2.1 applies with  $c_1 = 0$  and*

$c_2 = 1/\beta$ , and the model density is given by (2.6) with the above  $h$ ,  $c'_1 = (a^\top \Sigma^{-1} a)^{1/2}$  and  $c'_2 = (x - \theta)^\top \Sigma^{-1} a - (1/\beta)$ . The density was studied in [1, 2]. The exponential case with  $\alpha = 1$  simplifies with

$$p(x|\theta) = \frac{1}{\beta c'_1} \frac{\phi_d(x - \theta; \Sigma)}{R\left(\frac{c'_2}{c'_1}\right)}. \quad (2.7)$$

More generally for positive integer  $\alpha$ , the density's expression brings into play the  $(\alpha - 1)^{th}$  lower-truncated moment of a normal distribution. For instance, with  $\mathbb{E}\{(Z + \Delta)|Z + \Delta \geq 0\} = \Delta + R(\Delta)$ , we obtain for the case  $\alpha = 2$  the model density:

$$p(x|\theta) = \frac{\phi_d(x - \theta, \Sigma)}{(c'_1 \beta)^2} \left\{ \frac{c'_2/c'_1}{R(c'_2/c'_1)} + 1 \right\},$$

with the above  $c'_1$  and  $c'_2$ .

**(B)**  $\sqrt{\chi_k^2}$  mixing with  $h(v) = \frac{(\frac{1}{2})^{k/2-1}}{\Gamma(k/2)} v^{k-1}$ ,  $c_1 = 1$ ,  $c_2 = 0$ , and  $k > 0$ . The corresponding model density is given by (2.6) with the above  $h$ ,  $c'_1 = (1 + a^\top \Sigma^{-1} a)^{1/2}$ , and  $c'_2 = (x - \theta)^\top \Sigma^{-1} a$ .

The density was given in [5] and, as previously noted, the case  $k = 1$  reduces to the skew-normal case in (2.4). As in Example **(A)** for positive integer  $k$ , the density's expression involves a lower-truncated moment of a normal distribution.

**(C)** Kummer type II mixing with  $c_2 = c/\sigma$ ,  $c_1 = 0$ ,  $h(v) = \frac{\sigma^b}{\Gamma(a)\psi(a, 1-b, c)} \frac{v^{a-1}}{(v+\sigma)^{a+b}}$  with  $a, c, \sigma > 0$ ,  $b \in \mathbb{R}$ , and  $\psi$  the confluent hypergeometric function of type II defined for  $\gamma_1, \gamma_3 > 0$  and  $\gamma_2 \in \mathbb{R}$  as  $\psi(\gamma_1, \gamma_2, \gamma_3) = \frac{1}{\Gamma(\gamma_1)} \int_{\mathbb{R}_+} t^{\gamma_1-1} (1+t)^{\gamma_2-\gamma_1-1} e^{-\gamma_3 t} dt$ . This class of densities includes for  $b = -a$  the Gamma densities in **(A)**, as well as Beta type II densities for  $c = 0$  and  $b > 0$ . The resulting mean-mixture density is given by (2.6) and involves interesting expectations of the form  $\mathbb{E}\left(\frac{W^{a-1}}{(W+\sigma)^{a+b}} | W \geq 0\right)$  where  $W \sim N(\Delta, 1)$  with  $\Delta = c'_2/c'_1$ .

## 2.2. The prediction problem

Consider  $X|\theta \sim MMN_d(\theta, a, \Sigma_X, \mathcal{L}_1)$  and  $Y|\theta \sim MMN_d(\theta, a, \Sigma_Y, \mathcal{L}_2)$ , independently distributed as in Definition 2.1, i.e.

$$X|\theta, V_1 \sim N_d(\theta + V_1 a, \Sigma_X), Y|\theta, V_2 \sim N_d(\theta + V_2 a, \Sigma_Y), \text{ with } V_1 \sim \mathcal{L}_1, V_2 \sim \mathcal{L}_2. \quad (2.8)$$

Let  $p(x|\theta)$  and  $q(y|\theta)$  denote the conditional densities of  $X$  and  $Y$  given  $\theta$ , respectively. Based on observing  $X = x$ , we consider the problem of finding a suitable predictive density estimator  $\hat{q}(y; x)$  for  $q(y|\theta)$ ,  $y \in \mathbb{R}^d$ .

The ubiquitous Kullack-Leibler (KL) divergence between two Lebesgue densities  $f$  and  $g$  on  $\mathbb{R}^m$ , defined as

$$\rho(f, g) = \int_{\mathbb{R}^m} f(t) \log \frac{f(t)}{g(t)} dt,$$

is the basis of Kullback-Leibler loss given by

$$L(\theta, \hat{q}) = \rho(q_\theta, \hat{q}). \quad (2.9)$$

We will make use of Lemma 2.1's canonical form as in (2.1) to transform a mean mixture of normal distributions vector into two independent components and to capitalize on the corresponding simplification for KL divergence which is as follows.

**Lemma 2.2.** *Let  $T = (T_{(1)}, T_{(2)}) \in \mathbb{R}^m$  and  $U = (U_{(1)}, U_{(2)}) \in \mathbb{R}^m$  be random vectors subdivided into independently distributed components  $T_{(i)}$  and  $U_{(i)}$  of dimensions  $m_i$  for  $i = 1, 2$  with  $m_1 + m_2 = m$ . Denote  $f$  and  $g$  the densities of  $T$  and  $U$ , respectively, and  $f_1, f_2, g_1, g_2$  the densities of  $T_{(1)}, T_{(2)}, U_{(1)}, U_{(2)}$ , respectively. Then, we have*

$$\rho(f, g) = \rho(f_1, g_1) + \rho(f_2, g_2). \quad (2.10)$$

**Proof.** *By independence, we have*

$$\rho(f, g) = \mathbb{E}^T \left\{ \log \left( \frac{f_1(T_{(1)})f_2(T_{(2)})}{g_1(T_{(1)})g_2(T_{(2)})} \right) \right\} = \mathbb{E}^T \left\{ \log \left( \frac{f_1(T_{(1)})}{g_1(T_{(1)})} \right) \right\} + \mathbb{E}^T \left\{ \log \left( \frac{f_2(T_{(2)})}{g_2(T_{(2)})} \right) \right\},$$

which is (2.10). □

We evaluate the performance of the density estimators using KL loss in (2.9), and the associated KL risk function

$$R_{KL}(\theta, \hat{q}) = \int_{\mathbb{R}^d} \left\{ \int_{\mathbb{R}^d} q(y|\theta) \log \frac{q(y|\theta)}{\hat{q}(y; x)} dy \right\} p(x|\theta) dx. \quad (2.11)$$

For a prior density  $\pi$  for  $\theta$  with respect to a  $\sigma$ -finite measure  $\nu$ , it is known (e.g., [3, 4]) that the Bayes predictive density is given by

$$\hat{q}_\pi(y; x) = \int_{\mathbb{R}^d} q(y|\theta) p(x|\theta) \pi(\theta) d\nu(\theta). \quad (2.12)$$

A benchmark predictive density estimator for  $q(y|\theta)$ ,  $y \in \mathbb{R}^d$ , is given by the Bayes predictive density estimator  $\hat{q}_U(y; X)$ ,  $y \in \mathbb{R}^d$ , with respect to the uniform prior density on  $\mathbb{R}^d$ . It is known to be the minimum risk equivariant (MRE) predictive density estimator under changes of location, as well as minimax. In [16], a representation, which applies to both integrated squared-error loss and KL loss, for  $\hat{q}_U$  is provided. For our prediction problem, the following result makes use of this representation and summarizes the above optimality properties.

**Lemma 2.3.** *The MRE predictive density estimator of the density of  $Y$  relative to model (2.8) under KL loss, is given by the Bayes predictive density  $\hat{q}_U$  under prior  $\pi_U(\theta) = 1$ . Furthermore, we have*

$$\hat{q}_U(\cdot; X) \sim MMN_d(X, a, \Sigma_X + \Sigma_Y, \mathcal{L}_3), \quad (2.13)$$

where  $\mathcal{L}_3$  is the cdf of  $V_3 = V_2 - V_1$ . Finally,  $\hat{q}_U(y; X)$  is minimax under KL loss.

**Proof.** The MRE and minimax properties are given in [24] and [20]. For a location family prediction problem with  $X \sim p(x - \theta)$  and  $Y \sim q(y - \theta)$  independently distributed, it is shown in [16] that

$$\hat{q}_U(y; X) = q * \bar{p}(y - x), \text{ with } \bar{p}(t) = p(-t),$$

i.e., the convolution of  $q$  and the additive inverse of  $p$  followed by a change of location equal to  $x$ . For model (2.1), the above convolution  $q * \bar{p}$  is given by the density of  $Y - X$  in model (2.1) with  $\theta = 0$ , and the result follows since

$$Y - X|V_1, V_2 \sim N_d((V_2 - V_1)a, \Sigma_X + \Sigma_Y). \quad \square$$

Here, we can see that the MRE density estimator also belongs to the class of MMN distributions with same perturbation parameter  $a$  and location parameter  $x$ . As well, the distribution of the difference  $V_2 - V_1$  plays a key role in Theorem 2.3's representation of the MRE predictive density, and as illustrated in the next subsection of examples.

### 2.3. Minimax risk and entropy

The Kullback-Leibler risk expressions brings into play the entropy associated with MMN distributions. Such a measure is not easily manipulated into a closed form (see for instance [9] for the study of entropy for skewed-normal distributions), but they can be expressed in terms of the entropy of a univariate MMN distribution, as illustrated with the following expansion of the constant and minimax risk of the MRE density  $\hat{q}_U$  in the context of model (2.8). For a Lebesgue density on  $\mathbb{R}^d$ , defined as

$$H(f) = - \int_{\mathbb{R}^d} f(t) \log f(t) dt,$$

we will make use of the following well-known and easily established properties.

**Lemma 2.4.** (a) For  $T \in \mathbb{R}^d$  with density  $f$  and  $U = \psi(T) \sim g$  with  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^d$  invertible with inverse Jacobian  $J_\psi$ , we have  $H(g) = -\mathbb{E} \log |J_\psi| + H(f)$ ;

(b) Let  $T = (T_{(1)}, T_{(2)}) \sim f$  be a random vector with independently distributed components  $T_{(1)} \sim f_1$  on  $\mathbb{R}^{m_1}$  and  $T_{(2)} \sim f_2$  on  $\mathbb{R}^{m_2}$ . Then (as in Lemma 2.2), we have  $H(f) = H(f_1) + H(f_2)$ .

As implied by part (a) of the above lemma, the entropy  $H(f_\mu)$  is constant as a function of  $\mu$  for location family densities  $f_\mu(t) = f_0(t - \mu)$ , as is the case for  $MMN_d(\mu, b, \Sigma, \mathcal{L})$  densities. Now, we have the following dimension reduction decomposition for the entropy  $H_d(b, \Sigma, \mathcal{L})$  of a  $MMN_d(0, b, \Sigma, \mathcal{L})$  density.

**Lemma 2.5.** We have for  $d \geq 2$ :

$$H_d(b, \Sigma, \mathcal{L}) = H_1(\sqrt{b^\top \Sigma^{-1} b}, 1, \mathcal{L}) + \frac{d-1}{2} \{1 + \log(2\pi)\} + \frac{1}{2} \log |\Sigma|.$$

**Proof.** Let  $X \sim MMN_d(0, b, \Sigma, \mathcal{L})$ , which has entropy  $H_d(b, \Sigma, \mathcal{L})$ , and set  $Z = H \Sigma^{-1/2} X \sim f_Z$  with  $H$  orthogonal having first row  $\frac{b^\top \Sigma^{-1/2}}{\sqrt{(b^\top \Sigma^{-1} b)}}$ . It follows from part (a) of Lemma 2.4 that  $H(f_Z) = -\frac{1}{2} \log |\Sigma| + H_d(b, \Sigma, \mathcal{L})$ . From Lemma 2.1, we have  $Z = (Z_1, Z_{(2)})^\top$  with  $Z_1 \sim MMN_1(0, \sqrt{(b^\top \Sigma^{-1} b)}, 1, \mathcal{L})$  and  $Z_2 \sim N_{d-1}(0, I_{d-1})$  independently distributed, and the result follows from part (b) of Lemma 2.4 and a straightforward evaluation of the entropy  $H(\phi_{d-1})$ .  $\square$

With the above, we conclude with an expression for the constant and minimax risk.

**Theorem 2.2.** *In the context of model (2.8), the Kullback-Leibler risk of the MRE density  $\hat{q}_U$  is given by*

$$R_{KL}(\theta, \hat{q}_U) = H_1(\sqrt{a^\top \Sigma_S^{-1} a}, 1, \mathcal{L}_3) - H_1(\sqrt{a^\top \Sigma_Y^{-1} a}, 1, \mathcal{L}_2) + \frac{1}{2} \log \frac{\Sigma_S}{\Sigma_Y}, \quad (2.14)$$

with  $\Sigma_S = \Sigma_X + \Sigma_Y$ .

**Proof.** We have for  $\theta \in \mathbb{R}^d$

$$\begin{aligned} R_{KL}(\theta, \hat{q}_U) &= \mathbb{E}_\theta \{ \log q_\theta(Y) - \log \hat{q}_U(Y; X) \} \\ &= H(\hat{q}_U) - H(q_\theta) \\ &= H_d(a, \Sigma_S, \mathcal{L}_3) - H_d(a, \Sigma_Y, \mathcal{L}_2), \end{aligned}$$

by the independence of  $X$  and  $Y$ , the constancy of location family density  $q_\theta$ , and since  $Y - X | \theta \sim MMN_d(0, a, \Sigma_S, \mathcal{L}_3)$ . The result then follows from Lemma 2.5.  $\square$

The particular case with  $\Sigma_X = \sigma_X^2 I_d$  and  $\Sigma_Y = \sigma_Y^2 I_d$  follows directly from (2.14) and yields

$$R_{KL}(\theta, \hat{q}_U) = H_1\left(\frac{\|a\|}{\sigma_S}, 1, \mathcal{L}_3\right) - H_1\left(\frac{\|a\|}{\sigma_Y}, 1, \mathcal{L}_2\right) + \frac{d}{2} \log \frac{\sigma_S^2}{\sigma_Y^2}, \quad (2.15)$$

## 2.4. Minimum risk predictive density: Examples

Theorem 2.1 tells us that the minimum risk predictive density is given by  $\hat{q}_U(\cdot; X) \sim MMN_d(X, a, \Sigma_X + \Sigma_Y, \mathcal{L}_3)$  with  $\mathcal{L}_3$  the cdf of  $V_2 - V_1$ . The result is quite general and can be viewed as an extension of the multivariate normal case with  $a = 0$  and  $\hat{q}_U(\cdot; X) \sim N_d(X, \Sigma_X + \Sigma_Y)$ . Here are some interesting examples. When continuous, the mixing distributions can be taken to have a scale parameter equal to one without loss of generality, since a multiple can be integrated into the shape vector  $a$ .

- (A) For the case of degenerate  $V_2$  with  $\mathbb{P}(V_2 = v_2) = 1$ , i.e., when the distribution of  $Y | \theta$  is normal with  $Y \sim N_d(\theta + av_2, \Sigma_Y)$ , the minimum risk equivariant predictive density reduces to  $\hat{q}_U(\cdot; X) \sim MMN_d(X + av_2, -a, \Sigma_X + \Sigma_Y, \mathcal{L}_1)$ .
- (B) For the case of degenerate  $V_1$  with  $\mathbb{P}(V_1 = v_1) = 1$ , i.e., when the distribution of  $X$  is normal with  $X | \theta \sim N_d(\theta + av_1, \Sigma_X)$ , the minimum risk equivariant predictive density reduces to  $\hat{q}_U(\cdot; X) \sim MMN_d(X - av_1, a, \Sigma_X + \Sigma_Y, \mathcal{L}_2)$ .
- (C) We consider in this example  $V_1, V_2$  i.i.d. exponentially distributed with densities  $f(t) = e^{-t} \mathbb{I}_{(0, \infty)}(t)$ , as well as  $\Sigma_X = \sigma_X^2 I_d$  and  $\Sigma_Y = \sigma_Y^2 I_d$ . Here the distribution of  $V_3$  is Laplace or double-exponential with density  $\frac{1}{2} e^{-|v_3|}$  on  $\mathbb{R}$ . Therefore, from Theorem 2.1, we have

$$\begin{aligned} \hat{q}_U(y; x) &= \int_{\mathbb{R}} \frac{1}{2} e^{-|v_3|} \frac{1}{\sigma_S^d} \phi_d\left(\frac{y - x - av_3}{\sigma_S}\right) dv_3, \\ &= \phi_d(y - x; \sigma_S^2 I_d) \int_{\mathbb{R}_+} e^{-(v_3^2 \frac{\|a\|^2}{2\sigma_S^2} + v_3)} \cosh\left(v_3 \left(\frac{(y - x)^\top a}{\sigma_S^2}\right)\right) dv_3 \end{aligned}$$

with  $\sigma_S = (\sigma_X^2 + \sigma_Y^2)^{1/2}$ . By making use of Lemma 5.9 in the Appendix with  $A = \frac{\|a\|^2}{\sigma_S^2}$ ,



$B = -1 \pm \frac{(y-x)^\top a}{\sigma_S^2}$ , and  $c = 0$ , we obtain (for  $a \neq 0$ )

$$\hat{q}_U(y; x) = \sqrt{\frac{\pi\sigma_S^2}{2\|a\|^2}} \phi_d(y-x; \sigma_S^2 I_d) e^{\frac{\sigma_S^2}{2\|a\|^2} + \frac{\{(y-x)^\top a\}^2}{2\sigma_S^2\|a\|^2}} \times \left[ \left\{ e^{-\frac{(y-x)^\top a}{\|a\|^2}} \Phi\left(\frac{\sigma_S}{\|a\|}\left(\frac{(y-x)^\top a}{\sigma_S^2} - 1\right)\right) \right\} + \left\{ e^{\frac{(y-x)^\top a}{\|a\|^2}} \Phi\left(-\frac{\sigma_S}{\|a\|}\left(\frac{(y-x)^\top a}{\sigma_S^2} + 1\right)\right) \right\} \right].$$

(D) Consider  $V_1, V_2$  i.i.d. truncated normal distributed  $\text{TN}(0, 1)$  (or equivalently as  $\sqrt{\chi_1^2}$ ) for which  $X$  and  $Y$  are i.i.d. as multivariate skew-normal as in (2.4). A straightforward calculation yields the density

$$g_{V_3}(t) = 2\sqrt{2} \phi\left(\frac{t}{\sqrt{2}}\right) \Phi\left(-\frac{|t|}{\sqrt{2}}\right) \mathbb{I}_{\mathbb{R}}(t),$$

for  $V_3 =^d V_1 - V_2$ . It follows from Theorem 2.1, for  $\Sigma_X = \sigma_X^2 I_d$  and  $\Sigma_Y = \sigma_Y^2 I_d$ , denoting  $\sigma_S = (\sigma_X^2 + \sigma_Y^2)^{1/2}$ , that

$$\begin{aligned} \hat{q}_U(y; x) &= \int_{\mathbb{R}} 2\sqrt{2} \phi\left(\frac{t}{\sqrt{2}}\right) \Phi\left(-\frac{t}{\sqrt{2}}\right) \phi_d(y-x-at; \sigma_S^2 I_d) dt, \\ &= \frac{2}{\sqrt{\pi}} \phi_d(y-x; \sigma_S^2 I_d) \int_{\mathbb{R}_+} \Phi\left(-\frac{t}{\sqrt{2}}\right) e^{-\frac{t^2}{2}\left(\frac{1}{2} + \frac{a^\top a}{\sigma_S^2}\right)} \left\{ e^{\frac{(y-x)^\top at}{\sigma_S^2}} + e^{-\frac{(y-x)^\top at}{\sigma_S^2}} \right\} dt. \end{aligned}$$

Now, by making use of Lemma 5.9 with  $c = -\frac{\sqrt{2}}{2}$ ,  $A = \frac{2\sigma_S^2}{\sigma_S^2 + 2a^\top a}$ , and  $B = \pm \frac{(y-x)^\top a}{\sigma_S^2}$ , collecting terms, and setting  $f_k = \sqrt{\sigma_S^2 + ka^\top a}$ , we obtain the minimum risk equivariant predictive density

$$\begin{aligned} \hat{q}_U(y; x) &= \frac{4\sigma_S}{f_1} \phi_d\left(y-x; \sigma_S^2\left(I_d + \frac{aa^\top}{f_1^2}\right)\right) \\ &\times \left\{ \Phi_2\left(-\frac{(y-x)^\top a}{f_1 f_2}, \frac{\sqrt{2}(y-x)^\top a}{\sigma_S f_2}; \frac{-\sigma_S}{\sqrt{2} f_2}\right) + \Phi_2\left(\frac{(y-x)^\top a}{f_1 f_2}, -\frac{\sqrt{2}(y-x)^\top a}{\sigma_S f_2}; \frac{-\sigma_S}{\sqrt{2} f_2}\right) \right\}, \end{aligned}$$

where  $\Phi_2(z_1, z_2; \rho)$  the cdf evaluated at  $z_1, z_2 \in \mathbb{R}$  of a bivariate normal distributions with means equal to 0, variances equal to 1 and covariance equal to  $\rho$ . In the evaluation above, we made use of the identities  $(I - \frac{aa^\top}{f_2^2})^{-1} = I + \frac{aa^\top}{f_2^2}$  and  $|I + \frac{aa^\top}{f_1^2}| = 1 + \frac{a^\top a}{f_1^2}$ , which is a special case of the Sherman-Morrison formula for the matrix inversion of  $A + b_1 b_2^\top$  with  $A$  being a square matrix and  $b_1$  and  $b_2$  vectors of the same dimension.

### 3. Bayes posterior analysis and predictive densities

In this section, we expand on and document representations for Bayesian posterior and predictive densities for mean mixture of normal distributions.

### 3.1. Posterior densities

Bayesian posterior analysis of MMN models relate to the general form

$$X|K, \theta \sim f_{\theta, K}, K \sim g, \text{ and } \theta \sim \pi, \quad (3.16)$$

with observable  $X \in \mathbb{R}^d$ , density  $g$  of  $K$  free of  $\theta$ , and  $\pi$  prior density for  $\theta \in \mathbb{R}^d$ . Such a set-up leads to the following intermediate result, taken from [21].

**Lemma 3.6.** *For model (3.16), the posterior distribution of  $U = {}^d \theta|x$  admits the representation*

$$U|K' \sim \pi_{k', x} \text{ with } K' \sim g_{\pi, x}, \quad (3.17)$$

$\pi_{k', x}$  being the posterior density of  $\theta$  as if  $K = k'$  had been observed, and  $g_{\pi, x}(k') \propto g(k') m_{\pi, k'}(x)$  with  $m_{\pi, k'}$  being the marginal density of  $X$  as if  $K = k'$  had been observed.

We now apply the above to MMN distributions as in Definition 2.1.

**Example 3.2.** *We apply Lemma 3.6 to  $X|\theta \sim \text{MMN}_d(\theta, a, \Sigma, \mathcal{L})$  and the prior  $\theta \sim N_d(\mu, \Delta)$  with  $\Sigma, \Delta > 0$ . The above fits into model (3.16) with  $g$  taken to be the density of the mixing parameter  $K = V \sim \mathcal{L}$ , and  $f_{\theta, k}$  the  $N_d(\theta + ka, \Sigma)$  density. Conditional on  $K = k'$ , standard Bayesian analysis for the normal model tells us that*

$$\theta|k', x \sim N_d((I - P)x + P\mu - k'a, (I - P)\Sigma), \text{ and } X|k' \sim N_d(\mu + k'a, \Sigma + \Delta), \quad (3.18)$$

with  $P = \Sigma(\Sigma + \Delta)^{-1}$ , which yields the densities  $\pi_{k', x}$  and  $m_{\pi, k'}$  of Lemma 3.6. Then from Lemma 3.6, we infer that

$$\theta|x \sim \text{MMN}_d((I - P)x + P\mu, a^* = -a, (I - P)\Sigma, \mathcal{L}^*), \quad (3.19)$$

where the distribution  $\mathcal{L}^*$  has density

$$g_{\pi, x}(k') \propto g(k') e^{-\frac{A}{2}k'^2 + Bk'}, \text{ with } A = a^\top(\Sigma + \Delta)^{-1}a \text{ and } B = (x - \mu)^\top(\Sigma + \Delta)^{-1}a. \quad (3.20)$$

Furthermore, it follows immediately that

$$\mathbb{E}(\theta|x) = (I - P)x + P\mu - Pa \mathbb{E}(K'), \text{ with } K' \sim g_{\pi, x}. \quad (3.21)$$

**Remark 3.2.** *For the improper prior density  $\pi(\theta) = 1$ , one obtains  $\theta|x \sim \text{MMN}_d(x, -a, \Sigma, \mathcal{L})$  by a direct calculation. It can also be inferred from the above Example with  $\Delta = \tau^2 I_d$  and  $\tau^2 \rightarrow \infty$ .*

**Example 3.3.** *It is interesting to further study the above posterior distributions for the particular cases where the mixing density (i.e.,  $V$  or  $K$ ) of the MMN model is of the form*

$$g(k) \propto e^{-c_1 k^2/2 - c_2 k} \mathbb{I}_{(0, \infty)}(k), \quad (3.22)$$

with  $c_1 > 0, c_2 \in \mathbb{R}$  or  $c_1 = 0, c_2 > 0$ . Several of these distributions were presented in Example 2.1, but we recall that the cases  $c_1 > 0$  for instance, which correspond to truncated normal distributions on  $(0, \infty)$ , lead to skew-normal densities (2.4) for  $c_2 = 0$ . In the following, denote  $TN(a, b; (0, \infty))$  as a truncated normal distribution on  $(0, \infty)$  with shape parameter  $a \in \mathbb{R}$ , scale parameter  $b > 0$ , density  $\frac{1}{b} \frac{\phi((y-a)/b)}{\Phi(a/b)} \mathbb{I}_{(0, \infty)}(y)$ , and expectation  $a + bR(a/b)$ , with the reverse Mill's ratio  $R(\cdot)$ .

Now, it is easily seen for cases where  $K \sim g$  as in (3.22) that

$$\begin{aligned} g_{\pi,x}(k') &\propto e^{-(c_1+A)k'^2/2+(B-c_2)k'} \mathbb{I}_{(0,\infty)}(k') \\ &\propto \phi\left(\sqrt{A+c_1}k' - \frac{(B-c_2)}{\sqrt{A+c_1}}\right) \mathbb{I}_{(0,\infty)}(k'), \end{aligned}$$

which is the density of a  $TN\left(\frac{B-c_2}{A+c_1}, \frac{1}{\sqrt{A+c_1}}; (0, \infty)\right)$  distribution. Hence, the above, which yields the density associated with  $\mathcal{L}$ , provides a complete description of the posterior distribution in (3.19) for all considered cases of mixing density (3.22). Analogously, the corresponding expectation  $\mathbb{E}(K') = \frac{B-c_2}{A+c_1} + \frac{1}{\sqrt{A+c_1}} R\left(\frac{B-c_2}{\sqrt{A+c_1}}\right)$  provides an explicit expression for the posterior expectation  $\mathbb{E}(\theta|x)$  in (3.21).

## 3.2. Predictive densities

We now continue the above posterior analysis by focussing on the Bayes predictive density (i.e., the conditional density of  $Y$  given  $X = x$ ) for MMN distributions and a normally distributed prior for the unknown location parameter. In doing so, the following extension come into play.

**Definition 3.2.** A random vector  $Z \in \mathbb{R}^d$  is said to have a mean mixture of normal distribution with two directions, denoted as  $Z \sim MMN_d(\theta, a_1, a_2, \Sigma, \mathcal{L})$ , if it admits the representation

$$Z|V_1, V_2 \sim N_d(\theta + a_1W_1 + a_2W_2, \Sigma) \text{ with } (W_1, W_2) \sim \mathcal{L},$$

where  $\theta \in \mathbb{R}^d$  is a location parameter,  $a_1, a_2 \in \mathbb{R}^d$  are known perturbation vectors,  $\Sigma$  is a known positive definite covariance matrix, and  $W_1, W_2$  are scalar random variable with joint cdf  $\mathcal{L}$ .

We make use of the following intermediate result provided in [21] and applicable to mixture models of the form:

$$X|K, \theta \sim f_{\theta,K} \text{ with } K \sim g; Y|J, \theta \sim f_{\theta,J} \text{ with } J \sim h, \text{ and } \theta \sim \pi. \quad (3.23)$$

In the above set-up,  $X \in \mathbb{R}^d$  is observable, the mixing variables  $K$  and  $J$  are independently distributed with distributions free of  $\theta$ , the variables  $X$  and  $Y$  are conditionally independent on  $\theta$ , and  $\pi$  is a prior density for  $\theta \in \mathbb{R}^d$  with respect to a  $\sigma$ -finite measure  $\nu$ .

**Lemma 3.7.** For model (3.23), setting  $\pi_{k',x}$  and  $g_{\pi,x}$  as in Lemma 3.6, the Bayes predictive density of  $Y$  admits the mixture representation

$$Y|J', K' \sim q_{\pi}(\cdot|J', K'), \text{ with } J' \sim h, K' \sim g_{\pi,x} \text{ independent,}$$

and  $q_{\pi}(y|j', k') = \int_{\mathbb{R}^d} q_{\theta,j'}(y) \pi_{k',x}(\theta) d\nu(\theta)$ , which can be interpreted as the Bayes predictive density for  $Y$  as if  $Y \sim q_{\theta,j'}$  and  $K = k'$  had been observed.

Applied to mean mixture of multivariate normal distributions with a normal distributed prior, we obtain the following presented as a theorem.

**Theorem 3.3.** (a) For  $X|\theta \sim MMN_d(\theta, a_X, \sigma_X^2 I_d, \mathcal{L}_1)$  and  $Y|\theta \sim MMN_d(\theta, a_Y, \sigma_Y^2 I_d, \mathcal{L}_2)$  independent with prior  $\theta \sim N_d(\mu, \tau^2 I_d)$ , the Bayes predictive density for  $Y$  is that of a

$$MMN_d(\omega x + (1-\omega)\mu, -\omega a_X, a_Y, (\omega\sigma_X^2 + \sigma_Y^2)I_d, \mathcal{L})$$

distribution, with  $\mathcal{L}$  the joint cdf of  $(K', J')$  with independently distributed  $K' \sim g_{\pi, x}$  as in (3.20) and  $J' \sim \mathcal{L}_2$ , with  $\omega = \tau^2/(\tau^2 + \sigma_X^2)$ ,  $A = \|a_X\|^2/(\sigma_X^2 + \tau^2)$ , and  $B = \{(x - \mu)^\top a_X\}/(\sigma_X^2 + \tau^2)$ .

- (b) Moreover, whenever  $a_Y = ca_X$  for  $a_X \neq 0$  and a fixed  $c \in \mathbb{R}$ , the above predictive distribution is  $MMN_d(\omega x + (1 - \omega)\mu, a_X, (\omega\sigma_X^2 + \sigma_Y^2)I_d, \mathcal{L}_3)$ , with  $\mathcal{L}_3$  the cdf of  $cJ' - \omega K'$ , and  $(J', K')$  distributed as above. Finally, for  $a_X = 0$ , i.e., for  $X|\theta \sim N_d(\theta, \sigma_X^2 I_d)$ , the predictive distribution is  $MMN_d(\omega x + (1 - \omega)\mu, a_Y, (\omega\sigma_X^2 + \sigma_Y^2)I_d, \mathcal{L}_2)$

**Proof.** Part (b) follows immediately from part (a). For part (a), consider  $X' = X - K'a_X$  and  $Y' = Y - J'a_Y$ . The result then follows from Lemma 3.7 with the familiar predictive density estimation result:

$$Y'|J', K', X' \sim N_d(\omega X' + (1 - \omega)\mu, (\omega\sigma_X^2 + \sigma_Y^2)I_d),$$

implying

$$q_\pi(\cdot|J', K') \sim N_d(\omega x + (1 - \omega)\mu - \omega a_X K' + a_Y J', (\omega\sigma_X^2 + \sigma_Y^2)I_d),$$

matching Definition 3.2 with  $(W_1, W_2) =^d (K', J')$ . □

**Remark 3.3.** We point out that the minimum risk predictive density matches the density in (b) with  $\tau^2 = \infty$ , i.e.,  $\omega = 1$ .

## 4. Dominance Results

In this section, we first provide KL risk improvements on the MRE predictive density  $\hat{q}_U$  for estimating the density of  $Y|\theta \sim MMN_d(\theta, a, \sigma_Y^2 I_d, \mathcal{L}_2)$  based on  $X|\theta \sim MMN_d(\theta, a, \sigma_X^2 I_d, \mathcal{L}_1)$  with  $d \geq 4$ . Such improvements are necessarily minimax as a consequence of Theorem 2.3. Our findings cover two types of improvements: (i) plug-in type (Section 4.1), and (ii) Bayesian improvements (Section 4.2). Furthermore, we provide analogue results for certain type of restricted parameter spaces which are also applicable for  $d = 2, 3$ . Examples will be provided in Section 5.

The restriction to covariance matrices that are multiple of identity is justified by convenience and the fact that there is no loss of generality in doing so.

**Remark 4.4.** Predictive density estimates are intrinsic by nature which implies that the developments of this section, presented for  $\Sigma_X = \sigma_X^2 I_d$  and  $\Sigma_Y = \sigma_Y^2 I_d$  in model (2.1) with known  $\sigma_X^2$  and  $\sigma_Y^2$ , apply as well for  $\Sigma_Y = c\Sigma_X$  with known  $\Sigma_X, \Sigma_Y$ , and  $c = \sigma_Y^2/\sigma_X^2$ . Indeed, one can consider  $X' = \Sigma_X^{-1/2}X$  for which  $X|\theta \sim MMN_d(\Sigma_X^{-1/2}\theta, \Sigma_X^{-1/2}a, I_d, \mathcal{L}_1)$  to estimate the density of  $Y' = \Sigma_X^{-1/2}Y$ , for which  $Y'|\theta \sim MMN_d(\Sigma_X^{-1/2}\theta, \Sigma_X^{-1/2}a, cI_d, \mathcal{L}_2)$ . In doing so, one produces a predictive density estimator  $q_1(y') = \hat{q}(y'; x'), y' \in \mathbb{R}^d$ , for the density  $q_{Y'}$  of  $Y'$ , which equates to  $q_2(y) = \hat{q}(\Sigma_X^{-1/2}y; \Sigma_X^{-1/2}x) |\Sigma_X^{-1/2}|; y \in \mathbb{R}^d$ ; as a predictive density estimator of the density  $q_Y$  of  $Y$ . Moreover, the Kullback-Leibler  $\rho(q_{Y'}, q_1)$  and  $\rho(q_Y, q_2)$  are equal, i.e.

$$\int_{\mathbb{R}^d} q_{Y'}(t) \log \frac{q_{Y'}(t)}{q_1(t)} dt = \int_{\mathbb{R}^d} q_Y(t) \log \frac{q_Y(t)}{q_2(t)} dt,$$

as seen with the change of variables  $t \rightarrow \Sigma_X^{-1/2}t$ .

## 4.1. Plug-in type improvements

In the normal case with  $X|\theta \sim N_d(\theta, \sigma_X^2 I_d)$  and  $Y|\theta \sim N_d(\theta, \sigma_Y^2 I_d)$  independently distributed, the MRE predictive density  $\hat{q}_U(\cdot; X) \sim N_d(X, (\sigma_X^2 + \sigma_Y^2)I_d)$  is inadmissible for  $d \geq 3$  and can be improved by plug-in type densities of the form  $q_{\hat{\theta}}(\cdot; X) \sim N_d(\hat{\theta}(X), (\sigma_X^2 + \sigma_Y^2)I_d)$ . Indeed, the KL risk performance of  $q_{\hat{\theta}}$  relates directly to the ‘‘dual’’ point estimation risk of  $\hat{\theta}(X)$  for estimating  $\theta$  under squared error loss  $\|\hat{\theta} - \theta\|^2$ , with  $q_{\hat{\theta}}(\cdot; X)$  dominating  $\hat{q}_U(\cdot; X)$  if and only if  $\hat{\theta}(X)$  dominates  $X$  ([10]). For MMN distributions, such a duality does not deploy itself in the same way, but does so after transformation of  $(X, Y)$  to a canonical form and through the intrinsic nature of predictive densities. The following result exhibits this and is applicable to  $d \geq 4$ .

**Theorem 4.4.** *Consider  $X, Y$  distributed as in model (2.8) with  $a \neq 0, d \geq 4, \theta \in \mathbb{R}^d, \Sigma_X = \sigma_X^2 I_d$ , and  $\Sigma_Y = \sigma_Y^2 I_d$ , and the problem of obtaining a predictive density estimator  $\hat{q}(y; X), y \in \mathbb{R}^d$ , for the density of  $Y$ . Let  $H = \begin{pmatrix} h_1^\top \\ H_2 \end{pmatrix}$  be an  $d \times d$  orthogonal matrix such that  $h_1 = \frac{a}{\|a\|}$ . Define the densities*

$$q_1(\cdot; X) \sim MMN_1(h_1^\top X, \|a\|, (\sigma_X^2 + \sigma_Y^2), \mathcal{L}_3) \quad \text{and} \quad q_{2, \hat{\zeta}_2}(\cdot; X) \sim N_{d-1}(\hat{\zeta}_2(H_2 X), (\sigma_X^2 + \sigma_Y^2)I_{d-1}).$$

*Then, the predictive density  $q_{H, \hat{\zeta}_2}(y; X) = q_1(h_1^\top y; X) \times q_{2, \hat{\zeta}_2}(H_2 y; X), y \in \mathbb{R}^d$ , dominates  $\hat{q}_U$  under KL loss if and only if  $\hat{\zeta}_2(Z_2)$  dominates  $Z_2$  as an estimator of  $\zeta_2 \in \mathbb{R}^{d-1}$  under squared error loss  $\|\hat{\zeta}_2 - \zeta_2\|^2$  and for the model  $Z_2|\zeta_2 \sim N_{d-1}(\zeta_2, \sigma_X^2 I_{d-1})$ .*

**Proof.** Set

$$X' = HX = \begin{pmatrix} X'_1 \\ X'_{(2)} \end{pmatrix}, \quad Y' = HY = \begin{pmatrix} Y'_1 \\ Y'_{(2)} \end{pmatrix}, \quad \text{and} \quad \zeta = H\theta = \begin{pmatrix} \zeta_1 \\ \zeta_{(2)} \end{pmatrix}, \quad (4.24)$$

with  $X'_1 = h_1^\top X, X'_{(2)} = H_2 X, Y'_1 = h_1^\top Y, Y'_{(2)} = H_2 Y, \zeta_1 = h_1^\top \theta$ , and  $\zeta_{(2)} = H_2 \theta$ . From Lemma 2.1, we have that  $X'_1, X'_{(2)}, Y'_1$ , and  $Y'_{(2)}$  are independently distributed with  $X'_1 \sim MMN_1(\zeta_1, \|a\|, \sigma_X^2, \mathcal{L}_1), Y'_1 \sim MMN_1(\zeta_1, \|a\|, \sigma_Y^2, \mathcal{L}_2), X'_{(2)} \sim N_{d-1}(\zeta_{(2)}, \sigma_X^2 I_{d-1})$ , and  $Y'_{(2)} \sim N_{d-1}(\zeta_{(2)}, \sigma_Y^2 I_{d-1})$ .

Now consider the class of predictive densities of the form

$$q_{\hat{\zeta}_2}(y'; X') = q_1(y'_1; X'_1) \times q_{2, \hat{\zeta}_2}(y'_2; X'_2), \quad y' = (y'_1, y'_2) \in \mathbb{R}^d, \quad (4.25)$$

for estimating the density of  $Y'$ . As in Remark 4.4, the Kullback-Leibler risk performance of  $q_{H, \hat{\zeta}_2}(\cdot; X)$  for estimating the density of  $Y$  is equivalent to the Kullback-Leibler risk performance of  $q_{\hat{\zeta}_2}(\cdot; X')$  for estimating the density of  $Y'$ . Furthermore, observe that the MRE density estimator  $\hat{q}_U$  equates to density  $q_{\hat{\zeta}_{2,0}}(\cdot; X')$  with  $\hat{\zeta}_{2,0}(Y'_2) = Y'_2$ . It thus follows, with the independence of the

components of  $Y'$  and  $X'$ , Lemma 2.2, and setting  $Z_2 = X_{(2)}$  that

$$\begin{aligned}
R_{KL}(\theta, \hat{q}_U) - R_{KL}(\theta, q_{H, \hat{\zeta}_2}) &= R_{KL}(\theta, q_{\hat{\zeta}_2, 0}) - R_{KL}(\theta, q_{\hat{\zeta}_2}) \\
&= \mathbb{E} \log \left( \frac{q_1(Y'_1; X'_1)}{q_1(Y'_1; X'_1)} \right) + \mathbb{E} \log \left( \frac{q_{2, \hat{\zeta}_2}(Y'_2; X'_2)}{q_{2, \hat{\zeta}_2, 0}(Y'_2; X'_2)} \right) \\
&= \frac{1}{2(\sigma_X^2 + \sigma_Y^2)} \left( \mathbb{E} \|\hat{\zeta}_2(Z_2) - \zeta_2\|^2 - \mathbb{E} \|Z_2 - \zeta_2\|^2 \right), \quad (4.26)
\end{aligned}$$

which yields the result.  $\square$

The above dominance finding is quite general with respect to the specifications of  $a$ ,  $\mathcal{L}_1$ , and  $\mathcal{L}_2$  of model (2.8). Furthermore, observe by examining (4.26) that the risk difference depends on  $\theta$  only through  $\zeta_{(2)} = H_2\theta$  and this for any choice of  $H_2$ . More strikingly as seen with (4.26), the risk difference does not depend on the mixing distributions  $\mathcal{L}_1$  and  $\mathcal{L}_2$  and can be simply described by a quadratic risk difference of point estimators which arise in a  $(d-1)$  variate normal distribution problem. An illustration of Theorem 4.4 will be presented in Section 5.

## 4.2. Bayesian improvements

We now focus on Bayesian predictive densities that dominate  $\hat{q}_U$ . In doing so, we work with canonical forms as in Lemma 2.1, apply the partitioning argument of Lemma 2.2, and take advantage of known results for prediction in  $(d-1)$  multivariate normal models. We consider a class of improper priors on  $\theta$  which is the product measure of a (improper) uniform density over the linear subspace spanned by  $a$  and a second component of the prior ( $\pi_0$ ) supported on the subspace orthogonal to  $a$ . The measure of this nature splits resulting Bayes predictive densities into independent parts and leads to a decomposition the KL risk in two additive parts. Hence, the dominance result is obtained by dominating the part of the KL risk corresponding to the orthogonal space to  $a$ , where transformed variables are  $N_{d-1}$  distributed and where we can capitalize on known results. Namely, the superharmonicity of  $\pi_0$ , or its associated marginal density or its associated square root marginal density, will suffice for dominance and minimaxity.

**Theorem 4.5.** *Consider  $X, Y$  distributed as in model (2.8) with  $\Sigma_X = \sigma_X^2 I_d$ ,  $\Sigma_Y = \sigma_Y^2 I_d$ , and  $d \geq 2$ . Let  $H = \begin{pmatrix} h_1^\top \\ H_2 \end{pmatrix}$  be an  $d \times d$  orthogonal matrix such that  $h_1 = \frac{a}{\|a\|}$ . Let  $X', Y'$ , and  $\zeta$  be defined as in (4.24) and consider prior densities of the form*

$$\pi(\theta) = \pi_0(\zeta_{(2)}). \quad (4.27)$$

(a) *Then, the Bayes predictive density for  $Y$  is given by*

$$\hat{q}_\pi(y; X) = \hat{q}'_\pi(Hy; X'), y \in \mathbb{R}^d, \quad (4.28)$$

*with  $\hat{q}'_\pi(\cdot; x')$  the Bayes predictive density for  $Y'$  based on  $X'$ , given by*

$$\hat{q}'_\pi(y'; X') = \hat{q}_U(y'_1; X'_1) \times \hat{q}'_{\pi_0}(y'_{(2)}; X'_{(2)}), \quad (4.29)$$

*with: (i)  $\hat{q}_U(\cdot; X'_1)$  the MRE density, given in Theorem 2.3, of  $Y'_1 \sim MMN_1(\zeta_1, \|a\|, \sigma_Y^2, \mathcal{L}_2)$  based on  $X'_1 \sim MMN_1(\zeta_1, \|a\|, \sigma_X^2, \mathcal{L}_1)$ , and (ii)  $\hat{q}'_{\pi_0}(\cdot; X'_2)$  the Bayes predictive density for*

$Y'_{(2)} \sim N_{d-1}(\zeta_{(2)}, \sigma_Y^2 I_{d-1})$  based on  $X'_{(2)} \sim N_{d-1}(\zeta_{(2)}, \sigma_X^2 I_{d-1})$  and for prior density  $\pi_0(\zeta_{(2)})$  for  $\zeta_{(2)}$ ;

- (b) If  $d \geq 4$ , then  $\hat{q}_\pi$  given in (4.28) dominates the MRE  $\hat{q}_U$ , and is therefore minimax, if and only if  $\hat{q}'_{\pi_0}(\cdot; X'_2)$  dominates the MRE density for  $Y'_{(2)}$  based on  $X'_{(2)}$  given by a  $N_{d-1}(X'_{(2)}, (\sigma_X^2 + \sigma_Y^2)I_{d-1})$  density.

**Proof.**

- (a) Eq. (4.28) follows from the transformation of variables under the orthogonal matrix  $H$ . Note that the distribution of the transformed variables is

$$\begin{aligned} X' &\sim MMN_d(\zeta, a_0, \sigma_X^2 I_d, \mathcal{L}_1) \\ \text{and } Y' &\sim MMN_d(\zeta, a_0, \sigma_Y^2 I_d, \mathcal{L}_2), \end{aligned}$$

where  $a_0 = \left(\frac{a}{\|a\|}, 0, \dots, 0\right)^\top$ . The prior of the form (4.27) induces an improper uniform measure on  $\zeta_1$  and independent  $\pi_0(\zeta_{(2)})$  on  $\zeta_{(2)}$ . Along with the conditional independence of  $Y'_1$  and  $Y'_{(2)}$  given  $\zeta$ , we get the Bayes predictive density as (4.29).

- (b) Observe that the MRE density estimator  $\hat{q}_U(\cdot; X)$  corresponds to  $\pi_0(\theta) = 1$ , i.e., the improper uniform density on  $\zeta_{(2)} \in \mathbb{R}^{d-1}$ . By virtue of Lemma 2.2, the KL risk difference between  $\hat{q}_U(\cdot; X)$  and  $\hat{q}_\pi(\cdot; X)$  is then expressed as

$$\begin{aligned} R_{KL}(\theta, \hat{q}_U) - R_{KL}(\theta, \hat{q}_\pi) &= \mathbb{E} \log \hat{q}_\pi(Y; X) - \mathbb{E} \log \hat{q}_U(Y; X) \\ &= \mathbb{E} \log \hat{q}'_{\pi_0}(Y'_{(2)}; X'_{(2)}) - \mathbb{E} \log \hat{q}'_U(Y'_{(2)}; X'_{(2)}) \\ &= R_{KL}(\zeta_{(2)}, \hat{q}'_U) - R_{KL}(\zeta_{(2)}, \hat{q}'_{\pi_0}), \end{aligned}$$

and part (b) follows. □

**Remark 4.5.** *Theorem 4.5's dominance finding in part (b) is unified with respect to the model settings  $a$ ,  $\mathcal{L}_1$  and  $\mathcal{L}_2$ , as well as the dimension  $d \geq 4$ ,  $\sigma_X^2$ , and  $\sigma_Y^2$ . Furthermore, as seen in the lines of the proof, the difference in risks between the predictive densities  $\hat{q}_U$  and  $\hat{q}_\pi$ : (i) does not depend on the mixing  $\mathcal{L}_1$  and  $\mathcal{L}_2$ , and (ii) depends on  $\theta$  only through  $\zeta_{(2)} = H_2\theta$ .*

Starting with [14], continuing namely with [13], several Bayesian predictive densities  $\hat{q}'_{\pi_0}(\cdot; X'_2)$  have been shown to satisfy the dominance condition in part (b) of the above Theorem. Such choices lead to dominating predictive densities of  $\hat{q}_U$ . In [13], analogously to the quadratic risk estimation problem with multivariate normal observables (e.g., [27, 11]), sufficient conditions for minimaxity are conveniently expressed in terms of the marginal density of  $Z \sim N_{d-1}(\zeta_{(2)}, \sigma^2 I_{d-1})$  associated with density  $\pi_0$  and given by

$$m_{\pi_0}(z, \sigma^2) = \int_{\mathbb{R}^{d-1}} \phi_{d-1}(z - \zeta_{(2)}, \sigma^2 I_{d-1}) \pi_0(\zeta_{(2)}) d\zeta_{(2)}.$$

The superharmonicity of either  $\pi_0$ ,  $m_{\pi_0}(z, \sigma^2)$  for  $z \in \mathbb{R}^{d-1}$ , for various values of  $\sigma^2$ , or as well of  $\sqrt{m_{\pi_0}(z, \sigma^2)}$ , each lead to sufficient conditions for minimaxity. We recall here that the superharmonicity of  $h : \mathbb{R}^{d-1} \rightarrow \mathbb{R}$  holds whenever the Laplacian  $\Delta^2 h(t) = \sum_{i=1}^{d-1} \frac{\partial^2 h(t)}{\partial t_i^2}$  exists with  $\Delta^2 h(t) \leq 0$  for  $t \in \mathbb{R}^{d-1}$ .

**Corollary 4.1.** Consider the prediction context of Theorem 4.5 and a prior density  $\pi_0$  as in (4.27) other than the uniform density. Suppose that  $m_{\pi_0}(z, \sigma_X^2)$  is finite for all  $z \in \mathbb{R}^{d-1}$  and that  $d \geq 4$ . Then, the following conditions are each sufficient for  $\hat{q}_\pi(\cdot; X)$  given in (4.28) with prior density as in (4.27) to dominate the MRE density  $\hat{q}_U$ :

- (i)  $\Delta^2 m_{\pi_0}(z, \sigma^2) \leq 0$ ,  $z \in \mathbb{R}^{d-1}$ , for  $\frac{\sigma_X^2 \sigma_Y^2}{\sigma_X^2 + \sigma_Y^2} < \sigma^2 < \sigma_X^2$ , with strict inequality on a set of positive Lebesgue measure on  $\mathbb{R}^{d-1}$  for at least one  $\sigma^2$ ;
- (ii)  $\Delta^2 \sqrt{m_{\pi_0}(z, \sigma^2)} \leq 0$ ,  $z \in \mathbb{R}^{d-1}$ , for  $\frac{\sigma_X^2 \sigma_Y^2}{\sigma_X^2 + \sigma_Y^2} < \sigma^2 < \sigma_X^2$ , with strict inequality on a set of positive Lebesgue measure on  $\mathbb{R}^{d-1}$  for at least one  $\sigma^2$ ;
- (iii) The prior  $\pi_0$  is such that  $\Delta^2 \pi_0(\zeta_{(2)}) \leq 0$  a.e.

**Proof.** The results follow from part (b) of Theorem 4.5 and Theorem 1 - Corollary 2 in [13].  $\square$

Choices of the prior density  $\pi_0$  satisfying the conditions of Corollary 4.1 thus rest upon analyses for the normal case which are plentiful. In particular, several examples of  $\pi_0$ , and the resulting predictive density  $\hat{q}'_{\pi_0}$ , are provided in [13]. These provide explicit representations of minimax predictive densities  $\hat{q}_\pi$  given in (4.28). A detailed example is presented in Section 5.

The orthogonality decomposition used in this Section leads to a further interesting representation which generalizes the one obtained in the multivariate normal case, and for which we now expand upon. For the multivariate normal case, referring to Theorem 4.5's decomposition, with  $X'_{(2)} \sim N_{d-1}(\zeta_{(2)}, \sigma_X^2 I_{d-1})$  independent of  $Y'_{(2)} \sim N_{d-1}(\zeta_{(2)}, \sigma_Y^2 I_{d-1})$ , a well-known representation of the Bayes predictive density associated with prior density  $\pi_0$  for  $\zeta_{(2)}$ , given by [13], is

$$\hat{q}'_{\pi_0}(y'_{(2)}; x'_{(2)}) = \hat{q}'_U(y'_{(2)}; x'_{(2)}) \times \frac{m_{\pi_0}(w'_{(2)}; \sigma_W^2)}{m_{\pi_0}(x'_{(2)}, \sigma_X^2)}, \quad (4.30)$$

with  $w'_{(2)} = \frac{\sigma_X^2 y'_{(2)} + \sigma_Y^2 x'_{(2)}}{\sigma_X^2 + \sigma_Y^2}$  and  $\sigma_W^2 = \frac{\sigma_X^2 \sigma_Y^2}{\sigma_X^2 + \sigma_Y^2}$ , and where  $\hat{q}'_U(\cdot; X'_{(2)})$  is the MRE predictive density of the density of  $Y'_{(2)}$  based on  $X'_{(2)}$ , and given by a  $N_{d-1}(x'_{(2)}, (\sigma_X^2 + \sigma_Y^2) I_{d-1})$  density.

For the MMN case, we now have the following.

**Lemma 4.8.** For a prior  $\pi_0$  and  $H$  in Theorem 4.5, the corresponding Bayes predictive density  $\hat{q}_\pi$  admits the representation

$$\hat{q}_\pi(y; x) = \hat{q}_U(y; x) \times \frac{m_{\pi_0}(H_2 w, \sigma_W^2)}{m_{\pi_0}(H_2 x, \sigma_X^2)}, \quad (4.31)$$

with  $w = \frac{\sigma_X^2 y + \sigma_Y^2 x}{\sigma_X^2 + \sigma_Y^2}$ .

**Proof.** Using the set-up of Theorem 4.5, and expressions (4.28) and (4.29), the MRE predictive density is obtained as

$$\hat{q}_U(y; X) = \hat{q}_U(y_1, X_1) \times \hat{q}'_U(y'_{(2)}; X'_{(2)}), y \in \mathbb{R}^d.$$

Therefore, from (4.28) and (4.29) again, as well as from 4.30, we obtain

$$\hat{q}_\pi(y; X) = \hat{q}'_U(h_1^\top y; X_1) \times \hat{q}'_U(H_2 y; X_2) \times \frac{m_{\pi_0}(w'_{(2)}; \sigma_W^2)}{m_{\pi_0}(x'_{(2)}, \sigma_X^2)},$$



which yields the result.  $\square$

To conclude describing the dominance findings of this section and of Section 4.1, we point out that the plug-in type improvements of Theorem 4.5 and the Bayesian dominance results of Theorem 4.5 and Corollary 4.1 are applicable regardless of the choice of the orthogonal completion  $H_2$  of  $H$ , thus adding to choices of  $\pi_0$  leading to minimaxity. Furthermore, the above developments are unified and the findings are applicable for all MMN models (2.8) with  $\Sigma_X = \sigma_X^2 I_d$  and  $\Sigma_Y = \sigma_Y^2 I_d$ , as well as for  $\Sigma_Y = c\Sigma_X$  as justified in Remark 4.4.

**Remark 4.6.** *A particular appealing choice of  $H_2$ , which will be further explored below in Sections 4.3 and 5, is the such that  $H_2^\top H_2 = I_d - \frac{aa^\top}{a^\top a}$  in which case*

$$\|\zeta_{(2)}\|^2 = \theta^\top \left( I_d - \frac{aa^\top}{a^\top a} \right) \theta, \quad (4.32)$$

and spherically symmetric densities  $\pi_2(\zeta_2) = g(\|\zeta_2\|^2)$  lead to prior densities in (4.27) of the form

$$\pi(\theta) = g \left\{ \theta^\top \left( I_d - \frac{aa^\top}{a^\top a} \right) \theta \right\} = g \left( \left\| \theta - \frac{a^\top \theta}{a^\top a} a \right\|^2 \right). \quad (4.33)$$

Such densities do not depend on  $\|a\|$  and have contours given by hypersurfaces of cylinders with axis given by  $a$  (or  $h_1 = \frac{a}{\|a\|}$ ). Here is an example of three contours for  $d = 3$  and  $a = (1, 1, 1)^\top$ .

### 4.3. Restricted parameter spaces

Theorem 4.5's decomposition also leads to implications when there exists parametric restrictions on  $\zeta_{(2)} = H\theta$ . Statistical models where parametric restrictions are present appear naturally in a great variety of contexts, and there is a large literature on related inferential problems, namely for a decision-theoretic approach (e.g., [23, 28]). Questions of predictive analysis under parametric restrictions are also of interest with findings obtained in [22, 17, 10]. Namely, for normal models, specifically model (2.8) with  $a = 0$ ,  $\Sigma_X = \sigma_X^2 I_d$ ,  $\Sigma_Y = \sigma_Y^2 I_d$  with  $\theta$  constrained to a convex set  $C_0$  with non-empty interior, [10] showed that the Bayes predictive density associated with the uniform prior for  $\theta$  on  $C_0$  dominates the MRE predictive density under Kullback-Leibler loss. The next results extends this finding to MMN models.

**Theorem 4.6.** *Consider  $X, Y$  distributed as in model (2.8) with  $\Sigma_X = \sigma_X^2 I_d$ ,  $\Sigma_Y = \sigma_Y^2 I_d$ , and  $d \geq 2$ . Let  $C \subset \mathbb{R}^{d-1}$  be a convex set with non-empty interior, and let  $\pi_C(\theta) = \pi_{0,U}(\zeta_{(2)}) = I_C(\zeta_{(2)})$ . Then  $\hat{q}_{\pi_C}(\cdot; X)$  dominates  $\hat{q}_U(\cdot; X)$  under KL risk and the restriction  $\theta \in \{\theta \in \mathbb{R}^d : H_2\theta \in C\}$ .*

**Proof.** As in Theorem 4.5 and the given proof, we infer that  $\hat{q}_\pi$  given in (4.28) with prior density  $\pi(\theta) = \pi_0(\zeta_{(2)})$  for  $\zeta_{(2)} = H_2\theta$  dominates  $\hat{q}_U$  if and only if  $\hat{q}'_{\pi_0}(\cdot; X'_2)$  dominates the MRE density for  $Y'_{(2)} \sim N_{d-1}(\zeta_2, \sigma_Y^2 I_{d-1})$ .<sup>2</sup> But, since this latter dominance holds precisely for density  $\pi = \pi_C$  for the uniform density choice  $\pi_0 = \pi_{0,U}$  as shown in [10], the result follows.  $\square$

The setting of  $C$  above is quite general and interesting examples includes balls and cones. As earlier, the finding is unified and general to the MMN models. Here are two applications of Theorem 4.6.

<sup>2</sup>Said otherwise, part (b) of Theorem 4.5 could have been stated for  $d \geq 2$ , but this would lead to knowingly vacuous conditions in the absence of a parametric restriction.

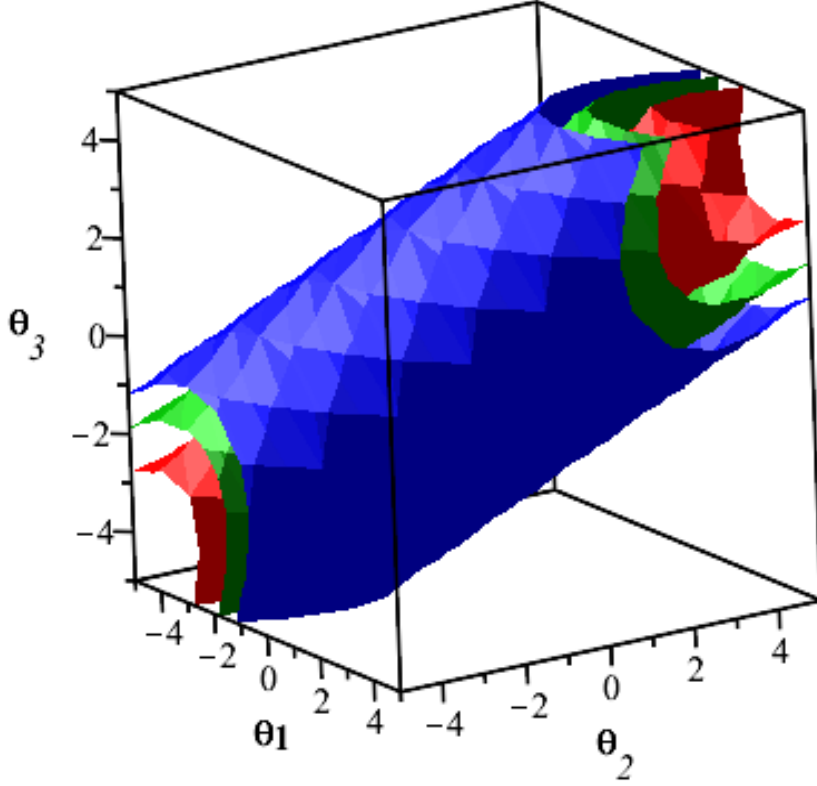


Figure 1: Contours of  $\pi(\theta)$  for  $d = 3$  and  $a = (1, 1, 1)^\top$ .

**Example 4.4.** Suppose  $d = 2$ ,  $a = (1, 1)^\top$ , and the parametric restriction  $\underline{c} \leq \theta_1 - \theta_2 \leq \bar{c}$ , with  $C = (\underline{c}, \bar{c})$  a strict subset of  $\mathbb{R}$ . The MRE density  $\hat{q}_U(\cdot; X)$  is that of  $MMN_2(X, a, (\sigma_X^2 + \sigma_Y^2)I_2, \mathcal{L}_3)$  distribution. In the context of Theorem 4.6, we have  $\zeta_{(2)} = \frac{\theta_1 - \theta_2}{\sqrt{2}}$  and the prior density  $\pi_C(\theta) = I_C(\theta_1 - \theta_2)$ . Theorem 4.5 tells us that the Bayes predictive density  $\hat{q}_{\pi_C}$  dominates the MRE  $\hat{q}_U$  with respect to KL loss and under the given parametric restriction.<sup>3</sup>

An explicit expression for  $\hat{q}_{\pi_C}$  is available from Lemma 4.8 with  $\pi_0$  the uniform  $U(\frac{\underline{c}}{\sqrt{2}}, \frac{\bar{c}}{\sqrt{2}})$  density for  $\zeta_{(2)}$ . As evaluated in [17], we obtain

$$\begin{aligned} \left( \frac{\sqrt{2}}{\bar{c} - \underline{c}} \right) m_{\pi_0}(z, \sigma^2) &= \int_{\underline{c}/\sqrt{2}}^{\bar{c}/\sqrt{2}} \phi(z - \zeta_{(2)}, \sigma^2) d\zeta_{(2)} \\ &= \Phi\left(\frac{z + \bar{c}/\sqrt{2}}{\sigma}\right) - \Phi\left(\frac{z + \underline{c}/\sqrt{2}}{\sigma}\right), \end{aligned}$$

<sup>3</sup>In Example 4.4, for the compact interval case say without loss of generality  $\underline{c} = -m$  and  $\bar{c} = m$ , there exists a much larger class of dominating predictive densities obtained by replacing the uniform density for  $\zeta_{(2)}$  by an even density  $\pi_0$  supported on  $(-m, m)$  that is increasing and logconcave on  $(0, m)$ . This is established as in Theorem 4.6 and making use of Theorem 3.2 in [10], which exploits a related point estimation finding in [19].

and (4.31) then yields

$$\hat{q}_{\pi_C}(y; x) = \hat{q}_U(y; x) \frac{\Phi\left(\frac{w+\bar{c}/\sqrt{2}}{\sigma_W}\right) - \Phi\left(\frac{w+\underline{c}/\sqrt{2}}{\sigma_W}\right)}{\Phi\left(\frac{x+\bar{c}/\sqrt{2}}{\sigma_X}\right) - \Phi\left(\frac{x+\underline{c}/\sqrt{2}}{\sigma_X}\right)}, y \in \mathbb{R},$$

with  $w = \frac{\sigma_X^2 y + \sigma_Y^2 x}{\sigma_X^2 + \sigma_Y^2}$ ,  $\sigma_W^2 = \frac{\sigma_X^2 \sigma_Y^2}{\sigma_X^2 + \sigma_Y^2}$ , and  $\hat{q}_U$  the MRE density which is that of a  $MMN_1(x, a, (\sigma_X^2 + \sigma_Y^2), \mathcal{L}_3)$  distribution.

**Example 4.5.** Theorem 4.6 applies for  $\theta$  restricted to a cylinder of radius, say  $m$ , with the axis along the direction  $a$ , i.e.,

$$C_m = \left\{ \theta \in \mathbb{R}^d : \left\| \theta - \frac{a^\top \theta}{a^\top a} a \right\| \leq m \right\};$$

examples of which are drawn in Figure 1. The dominating predictive density  $\hat{q}_{\pi_{C_m}}$  is Bayes with respect to the uniform prior density on  $C_m$ , which corresponds to (4.33) with  $g(t) = I_{(0,m)}(t)$ . An explicit expression for  $\hat{q}_{\pi_{C_m}}$  can be derived from Lemma 4.8 with  $\pi_0$  the uniform density on the ball  $B_m = \{t \in \mathbb{R}^{d-1} : \|t\| \leq m\}$  and marginal density

$$\begin{aligned} m_{\pi_0}(z, \sigma^2) &= \int_{B_m} \phi_{d-1}(z - \zeta_{(2)}, \sigma^2 I_{d-1}) d\zeta_{(2)} \\ &= F_{d-1, \frac{\|z\|^2}{\sigma^2}}\left(\frac{m^2}{\sigma^2}\right), \end{aligned}$$

with  $F_{\nu, \lambda}$  the cdf of a  $\chi_\nu^2(\lambda)$  distribution. From (4.31), we thus obtain

$$\hat{q}_{\pi_{C_m}}(y; x) = \hat{q}_U(y; x) \left( \frac{F_{d-1, \frac{\|H_2 w\|^2}{\sigma_W^2}}\left(\frac{m^2}{\sigma_W^2}\right)}{F_{d-1, \frac{\|H_2 x\|^2}{\sigma_X^2}}\left(\frac{m^2}{\sigma_X^2}\right)} \right), y \in \mathbb{R}^d,$$

with  $\|H_2 t\|^2 = t^\top \left( I - \frac{aa^\top}{a^\top a} \right) t$ , for  $t \in \mathbb{R}^d$ ,  $w = \frac{\sigma_X^2 y + \sigma_Y^2 x}{\sigma_X^2 + \sigma_Y^2}$ ,  $\sigma_W^2 = \frac{\sigma_X^2 \sigma_Y^2}{\sigma_X^2 + \sigma_Y^2}$ , and  $\hat{q}_U$  the MRE density which is that of a  $MMN_d(x, a, (\sigma_X^2 + \sigma_Y^2)I_d, \mathcal{L}_3)$  distribution.

## 5. Illustrations

We provide here illustrations of Theorems 4.4 and 4.5 accompanied by numerical comparisons and various observations.

**Example 5.6.** (A Bayesian minimax predictive density) In the context of Theorem 4.5, consider  $H_2$  as in Remark 4.6 combined with the harmonic prior density for  $\zeta_{(2)} \in \mathbb{R}^{d-1}$  given by  $\pi_0(\zeta_{(2)}) = \|\zeta_{(2)}\|^{-(d-3)}$  and which generates via (4.33) an “adjusted” harmonic prior density on  $\theta$  given by

$$\pi_H(\theta) = \left\| \theta - \frac{a^\top \theta}{a^\top a} a \right\|^{-(d-3)}. \quad (5.34)$$

Thus, the prior density is the product measure on  $\mathbb{R}^d$  with uniform prior on the linear subspace

spanned by  $a$  and the above harmonic measure on the  $(d-1)$ -dimensional chosen subspace orthogonal to  $a$ . Since  $\pi_0$  is superharmonic on  $\mathbb{R}^{d-1}$  for  $d \geq 4$ , it follows from Corollary 4.1 that the Bayes predictive density  $\hat{q}_{\pi_H}(\cdot; X)$  given in (4.28), as well as in (5.36) below, dominates the MRE density  $\hat{q}_U$  and is consequently minimax.

An explicit expression for  $\hat{q}_{\pi_H}$  is available from Lemma 4.31 with marginal density

$$m_{\pi_0}(z, \sigma^2) = \int_{\mathbb{R}^{d-1}} \phi_{d-1}(z - \zeta_{(2)}, \sigma^2 I_{d-1}) \frac{1}{\|\zeta_{(2)}\|^{(d-3)}} d\zeta_{(2)} = \sigma^{3-d} \mathbb{E} T^{\frac{(3-d)}{2}},$$

where  $T \sim \chi_{d-1}^2 \left( \frac{\|z\|^2}{\sigma^2} \right)$ . In particular for odd  $d \geq 5$ , as shown in the Appendix, one may obtain

$$m_{\pi_0}(z, \sigma^2) = (\|z\|^2)^{\frac{3-d}{2}} \left( 1 - e^{-\frac{\|z\|^2}{2\sigma^2}} \sum_{k=0}^{\frac{d-5}{2}} \left( \frac{\|z\|^2}{2\sigma^2} \right)^k \frac{1}{k!} \right) = r(\|z\|^2, \sigma^2) \quad (\text{say}), \quad (5.35)$$

which relates to known results on the inverse moments of a chi-square variable with even degrees of freedom (e.g., [8]), as well a closed form for an incomplete gamma function which intervenes in Komaki's [14] representation of  $m_{\pi_0}$ . From (4.31) and the above, we thus have

$$\hat{q}_{\pi_H}(y; x) = \hat{q}_U(y; x) \frac{r \left( \left\| w - \frac{a^\top w}{a^\top a} a \right\|^2, \sigma_W^2 \right)}{r \left( \left\| x - \frac{a^\top x}{a^\top a} a \right\|^2, \sigma_X^2 \right)}, \quad y \in \mathbb{R}^d, \quad (5.36)$$

where  $w$  and  $\sigma_W^2$  are as given in Lemma 4.8.

Risk differences between  $\hat{q}_U$  and  $\hat{q}_{\pi_H}$  are plotted in **Figure 2a** and **Figure 2b** as a function of  $\|\zeta_{(2)}\|^2$ , or equivalently as a function of

$$t = \frac{\|\zeta_{(2)}\|^2}{d-1} = \frac{1}{d-1} \left\| \theta - \frac{a^\top \theta}{a^\top a} a \right\|^2,$$

i.e., in terms of the average squared component of  $\zeta_{(2)}$ . The actual risks depend on the underlying mixing distributions  $\mathcal{L}_1$  and  $\mathcal{L}_2$ , but not the risk differences as previously observed in Remark 4.5. Observe as well that  $t$  is independent of  $\|a\|$  and only depends on the direction  $a/\|a\|$ . Figure 2a has  $\sigma_X^2 = 1, \sigma_Y^2 = 2$  and varying  $d$ , while Figure 2b has fixed  $d = 5, \sigma_X^2 = 1$  with  $\sigma_Y^2 = c\sigma_X^2$  and varying  $c$ . As seen with **Figure 2a**, the improvement in KL risk vanishes at  $t \rightarrow \infty$ , but gains in prominence with increasing  $d$ , and with the proximity of  $\theta$  to the linear subspace spanned by  $a$ . As seen with **Figure 2b**, the KL risk difference loses in prominence with larger  $c$  which is consistent with the fact that MRE density gains in reliability when the variance  $\sigma_X^2$  of the observable decreases.

Frequentist risk ratios between  $\hat{q}_U$  and  $\hat{q}_{\pi_H}$  are plotted in **Figure 2c** for  $\sigma_X^2 = 1, \sigma_Y^2 = 2$  and varying  $d$ . These ratios depend additionally on the mixing distributions  $\mathcal{L}_1$  and  $\mathcal{L}_2$  and they are set here with  $\sqrt{\chi_1^2}$  mixing (Example 2.1 (B)), i.e.,  $X|\theta$  and  $Y|\theta$  have skew-normal distributions with densities given in (2.4) and MRE density expanded upon in part (D) of Section 2.4. We further set  $a = \mathbf{1}_d = (1, \dots, 1)^\top$ , in which case the harmonic prior density on  $\theta$  in (5.34) reduces to  $\pi_0(\theta) = \|\theta - \bar{\theta}\mathbf{1}_d\|^{-(d-3)}$  with  $\bar{\theta} = \frac{1}{d} \sum_{i=1}^d \theta_i$ . With the above settings, the constant (and minimax)

risk of the MRE density can be computed from (2.15). For instance, we obtain  $R(\theta, \hat{q}_U) \approx 1.0954$  for  $d = 5$ ,  $\approx 1.5187$  for  $d = 7$  and  $\approx 1.9403$  for  $d = 9$ . These are close to linear with the term  $\frac{d}{2} \log \frac{\sigma_S^2}{\sigma_Y^2} = \frac{d}{2} \log \frac{3}{2}$  ( $\approx 1.0137$  for  $d = 5$ ,  $\approx 1.4191$  for  $d = 7$  and  $\approx 1.8246$  for  $d = 9$ ), representing the MRE risk for the normal case with  $a = 0$ , being dominant in (2.15). As seen in **Figure 2c**, where the risk ratios are plotted with respect to  $t = \frac{1}{d-1} \|\theta - \bar{\theta} \mathbf{1}_d\|^2$ , the gains increase in  $d$  and with the closeness of the  $\theta_i$ 's to  $\bar{\theta}$ .

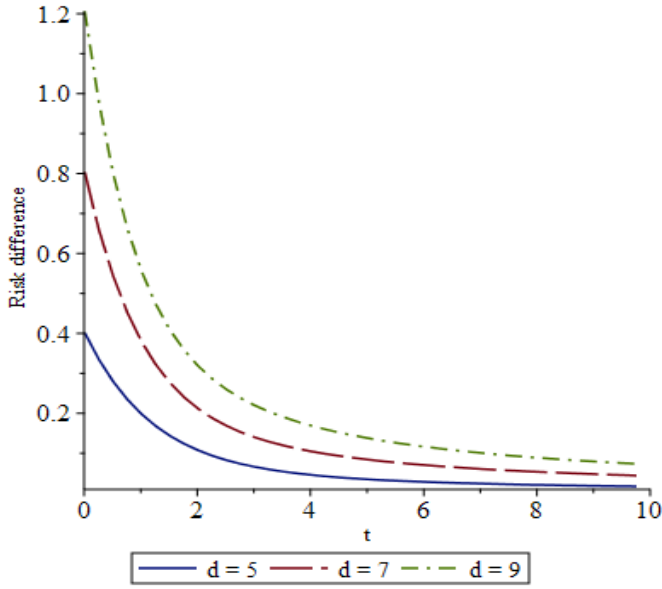
**Example 5.7.** (Plug-in type improved predictive density) In the context of Theorem 4.4, consider plug-in type predictive densities  $q_{H, \hat{\zeta}_2}(y; X) = q_1(h_1^\top y; X) \times q_{2, \hat{\zeta}_2}(H_2 y; X)$  with the choice of the James-Stein estimator  $\hat{\zeta}_2(Z_2) = \left(1 - \frac{(d-3)\sigma_X^2}{\|Z_2\|^2}\right) Z_2$  leading to the dominance of  $q_{H, \hat{\zeta}_2}$  over  $\hat{q}_U$  for  $d \geq 4$ . Both the dominating predictive density  $q_{H, \hat{\zeta}_2}$  and the actual difference in risks do depend on the choice of  $H_2$ , but the KL risk difference, as given in (4.26) and mentioned at the end of Section 4.1, is independent of the underlying mixing distributions and will thus coincide with the corresponding difference stemming for  $d - 1$  dimensional normal models and which have appeared many times in the literature. The difference in risks will be a function of  $\zeta_{(2)} = H_2 \theta$  in general, and more precisely as a function of  $\|\zeta_{(2)}\|^2$  in this case given that the James-Stein estimator is equivariant with respect to orthogonal transformations.

It is thus more interesting to look at the ratio of Kullback-Leibler risks and such ratios are presented in **Figure 2d** with the same settings as in Example 5.6, i.e., multivariate skew-normal models with  $\sqrt{\chi_1^2}$  mixing,  $\sigma_X^2 = 1$ ,  $\sigma_Y^2 = 2$ , and  $a = (1, \dots, 1)^T$ . Again here, the risk ratios are plotted with respect to  $t = \frac{1}{d-1} \|\theta - \bar{\theta} \mathbf{1}_d\|^2$ , the gains increase in  $d$  and with the closeness of the  $\theta_i$ 's to  $\bar{\theta}$ .

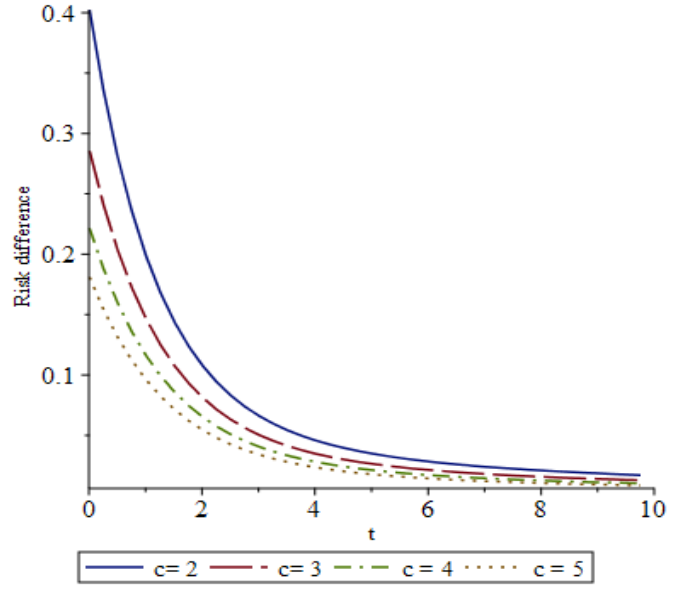
## Concluding remarks

In this work, we have addressed the problem of determining efficient predictive densities under Kullback-Leibler frequentist risk for multivariate skew-normal distributions and, more generally, for mean mixtures of multivariate normal (MMN) distributions, and provided Bayesian and plug-in type predictive densities which dominate the MRE density, and are minimax in four dimensions or more. In doing so, we have made use of a canonical transformation which leads to the decomposition of the Kullback-Leibler risk for the predictive densities being considered into two additive parts, one of which matching that of the MRE and minimax density, the other relating to a normal model and permitting improvement in view of shrinkage predictive density estimation results for such models. Further implications are provided for certain type of parametric restrictions. In addition, motivated by the relative paucity of analytical representations for Bayesian posterior and predictive densities, we have contributed such explicit representations.

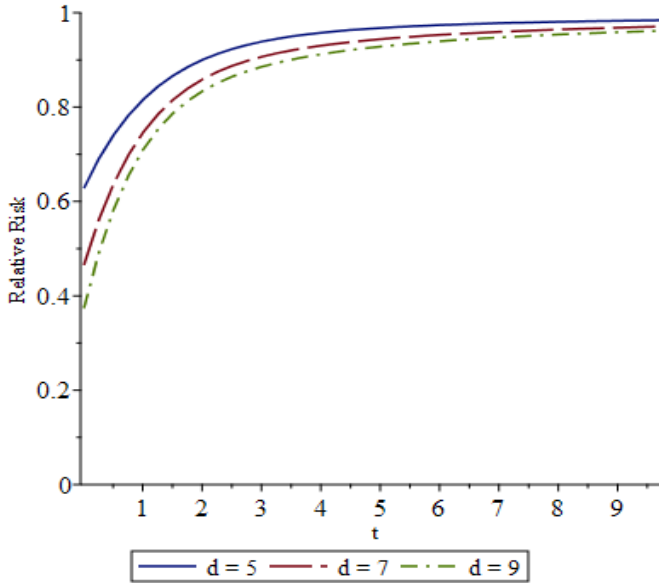
This work represents, to the best of our knowledge, a first foray of the study of predictive density estimation for MMN distributions. The findings are thus novel and they are also unified. The canonical transformation technique may well find further applications in predictive analysis, such as for mean-variance mixture of normal distributions. Extensions to other choices of loss (e.g.,  $\alpha$ -divergence) and to unknown covariance structures would be most interesting to investigate as well.



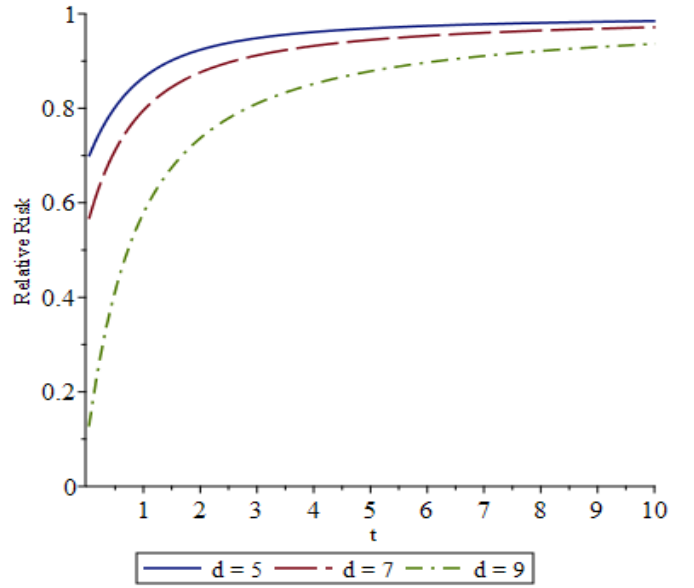
(a) The KL Risk difference between  $\hat{q}_U$  and  $\hat{q}_{\pi_0}$  as a function of  $t = \frac{\|\zeta_{(2)}\|^2}{d-1} = \frac{1}{d-1} \left\| \theta - \frac{a^\top \theta}{a^\top a} a \right\|^2$ , for  $\sigma_Y^2 = 2, \sigma_X^2 = 1$ .



(b) The KL Risk difference between  $\hat{q}_U$  and  $\hat{q}_{\pi_0}$  as a function of  $t = \frac{\|\zeta_{(2)}\|^2}{d-1} = \frac{1}{d-1} \left\| \theta - \frac{a^\top \theta}{a^\top a} a \right\|^2$ , for  $d = 5, \sigma_X^2 = 1$ , and  $c = \frac{\sigma_Y^2}{\sigma_X^2} = 2, 3, 4, 5$ .



(c) Kullback-Leibler risk ratio between  $\hat{q}_U$  and  $\hat{q}_{\pi_0}$  as a function of  $t = \frac{1}{d-1} \|\theta - \bar{\theta} \mathbf{1}_d\|^2$ , for  $\sigma_Y^2 = 2, \sigma_X^2 = 1$ .



(d) Kullback-Leibler risk ratio between  $\hat{q}_U$  and  $q_1 \times q_{2, \hat{\zeta}_{JS}}$  as a function of  $t = \frac{1}{d-1} \|\theta - \bar{\theta} \mathbf{1}_d\|^2$ , for  $\sigma_X^2 = 1, \sigma_Y^2 = 2$ , where  $\hat{\zeta}_{JS}$  is James-Stein estimator.

Figure 2: KL risk performance of the different predictive density estimators with the MRE

# Acknowledgements

Éric Marchand's research is supported in part by the Natural Sciences and Engineering Research Council of Canada. Pankaj Bhagwat is grateful to the ISM (Institut des sciences mathématiques) for financial support. Thanks to Jean-Philippe Burelle for useful discussions on geometric representations related to prior density (4.33).

# Appendix

**Lemma 5.9.** *For all  $B, c \in \mathbb{R}$ ,  $A \in \mathbb{R}_+$ , we have*

$$\int_0^\infty \Phi(ct) e^{-\frac{t^2}{2A} + Bt} dt = e^{\frac{AB^2}{2}} \sqrt{2\pi A} \Phi_2\left(\frac{cAB}{\sqrt{1+c^2A}}, B\sqrt{A}; \frac{c\sqrt{A}}{\sqrt{1+c^2A}}\right). \quad (5.37)$$

**Proof.** We have

$$\begin{aligned} e^{\frac{-AB^2}{2}} (2\pi A)^{-1/2} \int_0^\infty \Phi(ct) e^{-\frac{t^2}{2A} + Bt} dt &= \int_0^\infty \Phi(ct) \frac{1}{\sqrt{A}} \phi\left(\frac{t-AB}{\sqrt{A}}\right) dt \\ &= \mathbb{P}(U - cT \leq 0, -T \leq 0), \end{aligned}$$

with  $U, T$  independently distributed as  $N(0, 1)$  and  $N(\theta_T = AB, \sigma_T^2 = A)$ , respectively. The result follows since

$$(U - cT, -T)^\top \sim N_2\left(\begin{pmatrix} -cAB \\ AB \end{pmatrix}, \begin{bmatrix} 1 + c^2A & cA \\ cA & A \end{bmatrix}\right). \quad \square$$

**Proof of (5.35).** With the standard representation  $T|K \sim \chi_{d-1+2K}^2$  with  $K \sim \text{Poisson}\left(\frac{\|z\|^2}{2\sigma^2}\right)$ , we have

$$\begin{aligned} \mathbb{E} T^{(3-d)/2} &= \sum_{k=0}^\infty e^{-\frac{\|z\|^2}{2\sigma^2}} \frac{1}{k!} \left(\frac{\|z\|^2}{2\sigma^2}\right)^k \mathbf{E} (\chi_{d-1+2k}^2)^{\frac{(3-d)}{2}} \\ &= \frac{1}{2^{\frac{d-3}{2}}} e^{-\frac{\|z\|^2}{2\sigma^2}} \sum_{k=0}^\infty \left(\frac{\|z\|^2}{2\sigma^2}\right)^k \frac{1}{\Gamma\left(\frac{d-1}{2} + k\right)} \\ &= e^{-\frac{\|z\|^2}{2\sigma^2}} \left(\frac{\|z\|^2}{2\sigma^2}\right)^{-\frac{d-3}{2}} \sum_{k=\frac{d-3}{2}}^\infty \left(\frac{\|z\|^2}{2\sigma^2}\right)^k \frac{1}{k!}, \end{aligned}$$

which yields (5.35). □

# References

- [1] Abdi, M., Madadi, M., Balakrishnan, N. & Jamalizadeh, A. (2021). Family of mean-mixtures of multivariate normal distributions: Properties, inference and assessment of multivariate skewness. *Journal of Multivariate Analysis*, **181**, 104679.

- [2] Adcock, C.J. & Shutes, K. (2012). On the multivariate extended skew-normal, normal-exponential and normal-gamma distributions. *Journal of Statistical Theory and Practice*, **6**, 636–664.
- [3] Aitchison, J. & Dunsmore, I.R. (1975). *Statistical Prediction Analysis*. Cambridge University Press.
- [4] Aitchison, J. (1975). Goodness of prediction fit. *Biometrika*, **62**, 547-554.
- [5] Arrellano-Valle, R.B. & Azzalini, A. (2021). A formulation for continuous mixtures of multivariate normal distributions. *Journal of Multivariate Analysis*, **185**, 104780.
- [6] Azzalini, A. & Dalla Valle, A. (1996). The multivariate skew-normal distribution. *Biometrika*, **83**, 715–726.
- [7] Barndorff-Nielsen, O., Kent, J. & Sørensen, O. (1982). Normal variance-mean mixtures and  $z$  distributions. *International Statistical Review*, **50**, 145-159.
- [8] Bock, M.E., Judge, G.G. & Yancey, T.A. (1984). A simple form for the inverse moments of a non-central  $\chi^2$  and  $F$  random variables and certain confluent hypergeometric functions. *Journal of Econometrics*, **25**, 217–234.
- [9] Contreras-Reyes, J.E. & Arellano-Valle, R.B. (2012). Kullback-Leibler divergence measure for multivariate skew-normal distributions. *Entropy*, **14**, 1606-1626.
- [10] Fourdrinier, D., Marchand, É., Righi, A., & Strawderman, W.E. (2011). On improved predictive density estimation with parametric constraints. *Electronic Journal of Statistics*, **5**, 172-191.
- [11] Fourdrinier, D., Strawderman, W.E. & Wells, M. T. (2018). *Shrinkage estimation*. Springer series in statistics. Springer. New York, Dordrecht, Heidelberg, London.
- [12] George, E., Marchand, É., Mukherjee, G. & Paul, D. (2019). New and evolving roles of shrinkage in large-scale prediction and inference. BIRS Workshop Report.
- [13] George, E. I., Liang, F. and Xu, X. (2006). Improved minimax predictive densities under Kullback-Leibler loss. *Annals of Statistics*, **34**, 78-91.
- [14] Komaki, F. (2001). A shrinkage predictive distribution for multivariate normal observables. *Biometrika*, **88**, 859–864.
- [15] Kiefer, J. (1957). Invariance, minimax sequential estimation, and continuous time processes. *Annals of Mathematical Statistics*, **28**, 573–601.
- [16] Kubokawa, T., Marchand, É. & Strawderman, W.E. (2015). On predictive density estimation for location families under integrated squared error loss. *Journal of Multivariate Analysis*, **142**, 57–74.
- [17] Kubokawa, T., Marchand, É., Strawderman, W.E., Turcotte, J.P. (2013). Minimavity in predictive density estimation with parametric constraints. *Journal of Multivariate Analysis*, **116**, 382-397.
- [18] Kubokawa, T., Strawderman, W.E. & Yuasa, R. (2020). Shrinkage estimation of location parameters in a multivariate skew-normal distribution. *Communications in Statistics: Theory and Methods*, **49**, 2008–2024.
- [19] Kubokawa, T. (2005) Estimation of bounded location and scale parameters. *Journal of the Japanese Statistical Society*, **35**, 221–249.
- [20] Liang, F. & Barron, A. (2004). Exact minimax strategies for predictive density estimation, data compression, and model selection. *IEEE Transactions on Information Theory*, **50**, 2708–2726



- [21] LMoudden, A. & Marchand, É. (2021). Bayesian estimation and prediction for certain type of mixtures. *Communications in Statistics: Theory and Methods*, DOI:10.1080/03610926.2021. 1913185.
- [22] Marchand, É. & Sadeghkhan, N. (2018). On predictive density estimation with additional information. *Electronic Journal of Statistics*, **12**, 4209-4238.
- [23] Marchand, É. & Strawderman, W.E. (2004). Estimation in Restricted Parameter Spaces: A Review, *Festschrift for Herman Rubin*. Institute of Mathematical Statistics Lecture Notes-Monograph Series, pp. 21-44.
- [24] Murray, G.D. (1977). A note on the estimation of probability density functions, *Biometrika*, **64**, 150–152.
- [25] Negarestanu, H., Jamalizadeh, A., Shafiei, S. & Balakrishnan, N. (2019). Mean mixtures of normal distributions: Properties, inference and applications. *Metrika*, **82**, 501–528.
- [26] Stein, C. (1956). Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, I, pp. 197-206. University of California Press.
- [27] Stein, C. (1981). Estimation of the mean of a multivariate normal distribution. *Annals of Statistics*, **9**, 1135–1151.
- [28] van Eeden, C. (2006). *Restricted parameter space problems: Admissibility and minimaxity properties*. Lecture Notes in Statistics, **188**, Springer.