

# Generalized maximum likelihood estimates for exponential families

Imre Csiszár · František Matúš

Received: 21 July 2006 / Revised: 5 April 2007 / Published online: 11 August 2007  
© Springer-Verlag 2007

**Abstract** For a standard full exponential family on  $\mathbb{R}^d$ , or its canonically convex subfamily, the generalized maximum likelihood estimator is an extension of the mapping that assigns to the mean  $a \in \mathbb{R}^d$  of a sample for which a maximizer  $\vartheta^*$  of a corresponding likelihood function exists, the member of the family parameterized by  $\vartheta^*$ . This extension assigns to each  $a \in \mathbb{R}^d$  with the likelihood function bounded above, a member of the closure of the family in variation distance. Its detailed description, complete characterization of domain and range, and additional results are presented, not imposing any regularity assumptions. In addition to basic convex analysis tools, the authors' prior results on convex cores of measures and closures of exponential families are used.

**Mathematics Subject Classification (2000)** Primary: 60A10; Secondary: 62H12 · 62B10

---

This paper is dedicated to the memory of Albert Perez (1920–2003).

---

This work was supported by the Hungarian National Foundation for Scientific Research under Grant T046376 and by Grant Agency of Academy of Sciences of the Czech Republic under Grant IAA 100750603.

---

I. Csiszár  
A. Rényi Institute of Mathematics, Hungarian Academy of Sciences,  
1364 Budapest, P.O. Box 127, Hungary  
e-mail: csiszar@renyi.hu

F. Matúš (✉)  
Institute of Information Theory and Automation, Academy of Sciences of the Czech Republic,  
Pod vodárenskou věží 4, 182 08 Prague, Czech Republic  
e-mail: matus@utia.cas.cz

**Keywords** Accessible face · Convex core · Exponential family · Kullback–Leibler information divergence · Log-convexity · Maximum likelihood · Partial mean · Variation distance closure

### 1 Introduction

1.1 Let  $\mu$  be a nonzero Borel measure on  $\mathbb{R}^d$  whose log-Laplace transform

$$\Lambda(\vartheta) = \Lambda_\mu(\vartheta) \triangleq \ln \int_{\mathbb{R}^d} e^{\langle \vartheta, x \rangle} \mu(dx), \quad \vartheta \in \mathbb{R}^d,$$

has the effective domain  $dom(\Lambda) = \{\vartheta : \Lambda(\vartheta) < +\infty\}$  nonempty. The full exponential family  $\mathcal{E} = \mathcal{E}_\mu$ , determined by  $\mu$  and the identity mapping as canonical statistic, consists of the probability measures (pm’s)  $Q_\vartheta = Q_{\mu, \vartheta}$ ,  $\vartheta \in dom(\Lambda)$ , with  $\mu$ -densities  $x \mapsto e^{\langle \vartheta, x \rangle - \Lambda(\vartheta)}$ . Here,  $\vartheta$  is called canonical parameter. The mapping  $\vartheta \mapsto Q_\vartheta$  may be many-to-one.

Throughout this paper the symbol  $\mathcal{E}$  denotes a nonempty convex subset of  $dom(\Lambda)$  and  $\mathcal{E}_\mathcal{E} = \mathcal{E}_{\mu, \mathcal{E}}$  denotes the *canonically convex exponential family*  $\{Q_\vartheta : \vartheta \in \mathcal{E}\}$ . The full family  $\mathcal{E}$  is regarded as the family  $\mathcal{E}_\mathcal{E}$  with  $\mathcal{E} = dom(\Lambda)$ , even if  $\mathcal{E} = \mathcal{E}_\mathcal{E}$  may hold also for other choices of  $\mathcal{E}$ . All results in this paper are stated for general  $\mathcal{E}$ , full exponential families are covered as the instance  $\mathcal{E} = dom(\Lambda)$ . Exponential families with a canonical statistic different from the identity mapping will not be addressed, but the results easily extend to those, via reduction by sufficiency to “standard” families as above.

For  $a \in \mathbb{R}^d$  a maximizer  $\vartheta^*$  of the function  $\vartheta \mapsto \langle \vartheta, a \rangle - \Lambda(\vartheta)$  subject to  $\vartheta \in \mathcal{E}$  has a statistical interpretation when  $a$  equals the mean of an i.i.d. sample from a probability measure  $Q_\vartheta$  with  $\vartheta \in \mathcal{E}$  unknown. Then  $\vartheta^*$  is a maximum likelihood estimate (MLE) of the unknown parameter, from this sample. As well known [1], a sufficient condition for the existence of MLE is the equality of  $a$  to the mean of  $Q_\theta$  for some  $\theta \in \mathcal{E}$ . Then  $\vartheta^* = \theta$  maximizes the above function even when  $\vartheta$  ranges over the whole  $dom(\Lambda)$ , and

$$[\langle \vartheta^*, a \rangle - \Lambda(\vartheta^*)] - [\langle \vartheta, a \rangle - \Lambda(\vartheta)] = D(Q_{\vartheta^*} \| Q_\vartheta), \quad \vartheta \in \mathcal{E}.$$

Here,  $D$  denotes Kullback–Leibler information divergence ( $I$ -divergence or relative entropy), defined for any pm’s  $P$  and  $Q$  by

$$D(P \| Q) \triangleq \begin{cases} \int \ln \frac{dP}{dQ} dP, & \text{if } P \ll Q, \\ +\infty, & \text{otherwise.} \end{cases}$$

The authors have recently shown [6,5] that for any  $a \in \mathbb{R}^d$  such that

$$\Psi^*(a) = \Psi_{\mu, \mathcal{E}}^*(a) \triangleq \sup_{\vartheta \in \mathcal{E}} [\langle \vartheta, a \rangle - \Lambda(\vartheta)] \tag{1}$$

is finite, thus  $a \in \text{dom}(\Psi^*)$ , there exists a unique pm  $R^*(a)$  with the property

$$\Psi^*(a) - [\langle \vartheta, a \rangle - \Lambda(\vartheta)] \geq D(R^*(a) \| Q_\vartheta), \quad \vartheta \in \mathcal{E}. \tag{2}$$

This  $R^*(a) = R_{\mu, \mathcal{E}}^*(a)$  has been called *generalized maximum likelihood estimate* (GMLE). In the case of existence of an MLE  $\vartheta^*$ , a direct substitution gives that  $R^*(a) = Q_{\vartheta^*}$ , thus  $\vartheta^*$  parameterizes the GMLE. The inequality (2) implies that for each  $\vartheta \in \mathcal{E}$  satisfying  $\langle \vartheta, a \rangle - \Lambda(\vartheta) \geq \Psi^*(a) - \varepsilon$  (such  $\vartheta$  may be called an  $\varepsilon$ -MLE), the pm  $Q_\vartheta$  belongs to an information divergence ball of radius  $\varepsilon$  centered at  $R^*(a)$ . The concept of GMLE should be of interest in those cases when no MLE exists (not even in the slightly modified sense discussed in Sect. 3, corresponding to the closure of a projection of the set  $\mathcal{E}$ ), but the GMLE exists, thus it can serve as a substitute of MLE. In these cases, the GMLE does not belong to the exponential family  $\mathcal{E}$ , see the passage after Theorem 3.2, although it always belongs to  $cl_v(\mathcal{E}_\mathcal{E})$ , the closure of  $\mathcal{E}_\mathcal{E}$  in variation distance, see Remark 1.2 below. This means that for GMLE to go beyond MLE it is necessary that  $\mathcal{E}$  be not closed in variation distance, which excludes statistical models of continuous type. For models of discrete or mixed type, however, it is not uncommon that the log-likelihood function is bounded but has no maximizer; then GMLE provides a remedy to the nonexistence of MLE.

*Remark 1.1* The notation points to the fact that  $\Psi^*$  is the convex conjugate of the function  $\Psi$  equal to  $\Lambda$  on  $\mathcal{E}$  and to  $+\infty$ , otherwise. Thus, (2) sharpens the Fenchel inequality for  $\Psi$ , stating that the left-hand side of (2) is nonnegative. Note that if a full exponential family is considered then  $\Psi = \Lambda$ ,  $\Psi^* = \Lambda^*$ .

*Remark 1.2* For any sequence  $\vartheta_n$  in  $\mathcal{E}$  with  $\langle \vartheta_n, a \rangle - \Lambda(\vartheta_n) \rightarrow \Psi^*(a) < \infty$  it follows from (2) that the sequence  $Q_{\vartheta_n}$  converges to  $R^*(a)$  in reversed information (*rI*-) divergence,  $D(R^*(a) \| Q_{\vartheta_n}) \rightarrow 0$ , and hence by the Pinsker inequality also in variation distance. This shows that the GMLE  $R^*(a)$  is uniquely determined by (2), and belongs to  $cl_v(\mathcal{E}_\mathcal{E})$ . This closure has been recently described in detail [8].

*Remark 1.3* The supremum in (1) can be finite also if  $a$  is outside the convex support of  $\mu$ , denoted by  $cs(\mu)$ , though not for a full family since  $\text{dom}(\Lambda^*) \subseteq cs(\mu)$  [1, Theorem 9.1]. Then the maximization has no direct statistical interpretation since the mean of a sample from a distribution in the family  $\mathcal{E}_\mathcal{E}$  belongs to  $cs(\mu)$  with probability 1. As explained below, it is still useful to consider MLE and GMLE also when  $a \notin cs(\mu)$ , and even when  $a$  is outside the affine hull of  $cs(\mu)$ , denoted by  $aff(\mu)$ . In that case, different parameter sets  $\mathcal{E}$  can give rise to the same family of pm's  $\mathcal{E}_\mathcal{E}$  but to different GMLE's, see Example 3.5 in Sect. 3. For this reason, the choice of  $\mathcal{E}$  is left free in this paper, unlike in [8] where from the possible parameter sets for the family  $\mathcal{E}_\mathcal{E}$  a natural one is selected.

1.2 The goal of this paper is to study properties of the GMLE in full generality. This generality is not for its own sake. Indeed, suppose a characterization of  $R^*(a)$  were required only for full families  $\mathcal{E}_\mu$  in minimal representation, that is, such that both  $cs(\mu)$  and  $\text{dom}(\Lambda_\mu)$  have nonempty interiors. As illustrated in Example 1.4

below, even in that case there may exist a set  $A$  that contains the sample mean with positive probability (for all sample sizes) such that the characterization of  $R^*(a)$  for  $a \in A$  necessitates the maximization of  $\langle \vartheta, a \rangle - \Lambda_\nu(\vartheta)$  subject to  $\vartheta \in \Xi$  where  $\nu$  is a restriction of  $\mu$ ,  $\Xi = \text{dom}(\Lambda_\mu)$  is a proper subset of  $\text{dom}(\Lambda_\nu)$  and, more remarkably, no  $a \in A$  is contained in  $\text{aff}(\nu)$ . Thus, in order to determine GMLE's for full exponential families in minimal representation in all cases of statistical interest, it is necessary to consider also non-full families with underlying measure concentrated on a proper affine subspace of  $\mathbb{R}^d$  and  $a$  not contained in that subspace. We admit that no such intrinsic need is apparent for going beyond the case when  $\text{dom}(\Lambda)$  has nonempty interior that contains some  $\vartheta \in \Xi$ , which will be called the *nondegenerate case*. We do not make that restriction mainly for mathematical completeness. This does lead to technical problems not present in the nondegenerate case, but these are not difficult to overcome via the concept of partial mean, see Sect. 2.4. Full generality is useful also for at least one application, viz. an explicit description of generalized  $rI$ -projections to canonically convex exponential families, even of pm's that do not have a mean, as indicated in Sect. 5.

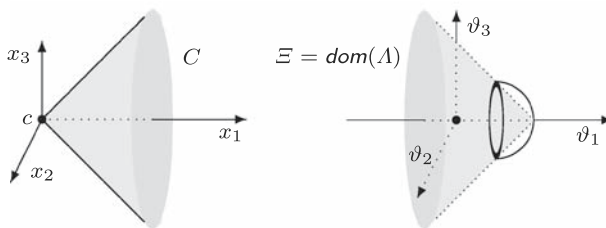
*Example 1.4* Let  $C$  be the closed convex cone in  $\mathbb{R}^3$  given by  $x_1 \geq 0$  and  $x_1^2 \geq x_2^2 + x_3^2$ . Let  $\mu$  be sum of the unit point mass  $\nu$  at the origin  $c = (0, 0, 0)$  and of the pm concentrated on the boundary  $\partial C$  of  $C$  with the density  $\frac{1}{2\pi} e^{-r} dr d\phi$ ,  $r \geq 0$ ,  $0 \leq \phi < 2\pi$ , in cylindrical coordinates. The convex support of  $\mu$  equals  $C$ , and for  $\vartheta = (\vartheta_1, \vartheta_2, \vartheta_3)$  in the interior of  $(1, 0, 0) - C$

$$\begin{aligned} \Lambda(\vartheta) &= \ln \left[ 1 + \frac{1}{2\pi} \int_0^{2\pi} d\phi \int_0^{+\infty} e^{r(\vartheta_1 - 1 + \vartheta_2 \cos \phi + \vartheta_3 \sin \phi)} dr \right] \\ &= \ln \left[ 1 + [(\vartheta_1 - 1)^2 - \vartheta_2^2 - \vartheta_3^2]^{-\frac{1}{2}} \right] \end{aligned}$$

and  $\Lambda(\vartheta) = +\infty$  otherwise. Thus  $\text{dom}(\Lambda)$  equals this shifted open cone, see Fig. 1.

Consider the full exponential family  $\mathcal{E}_\mu$ , thus let  $\Xi = \text{dom}(\Lambda)$ . If  $a$  is in the interior of  $C$ , elementary calculations show that the function  $\vartheta \mapsto \langle \vartheta, a \rangle - \Lambda(\vartheta)$  has a unique maximizer  $\vartheta^* \in \text{dom}(\Lambda)$ , the pm  $Q_{\vartheta^*}$  has the mean  $a$  and equals  $R^*(a)$ .

For  $a = (a_1, a_2, a_3) \in \partial C$  the function has no maximizer but is bounded above. In fact,  $\langle \vartheta, a \rangle - \Lambda(\vartheta) < \langle \vartheta, a \rangle \leq a_1$  for  $\vartheta \in \text{dom}(\Lambda)$  whence  $\Lambda^*(a) \leq a_1$  similarly



**Fig. 1** Illustrations of Examples 1.4 and 5.5

to [6, Example 2]. Actually,  $\Lambda^*(a) = a_1$  by limiting for  $\vartheta_n = (1 - \frac{1}{n}, 0, 0) + n^2b \in \text{dom}(\Lambda)$ ,  $n \rightarrow \infty$ , where  $b \in -\partial C$  is nonzero and orthogonal to  $a$ . Since  $\Lambda^*(a) = \Psi^*(a)$  is finite, the concept of GMLE applies, and (2) with  $R^*(a)$  replaced by  $\nu$  rewrites to

$$a_1 - [\langle \vartheta, a \rangle - \Lambda(\vartheta)] \geq D(\nu \| Q_\vartheta) = \Lambda(\vartheta), \quad \vartheta \in \mathcal{E},$$

which obviously holds. Thus,  $R_{\mu, \mathcal{E}}^*(a) = \nu$  by uniqueness of GMLE discussed in Remark 1.2.

On the other hand,  $\Lambda_\nu = 0$  and the family  $\mathcal{E}_{\nu, \mathcal{E}} = \{\nu\}$  is not in minimal representation. Obviously,  $\Psi_{\nu, \mathcal{E}}^*(a)$  equals  $a_1$  if  $a \in C$  and  $+\infty$  otherwise. In the cases  $a \in \partial C$ , the equalities  $\Psi_{\mu, \mathcal{E}}^*(a) = \Psi_{\nu, \mathcal{E}}^*(a)$  and  $R_{\nu, \mathcal{E}}^*(a) = R_{\mu, \mathcal{E}}^*(a)$  hold, and thus a characterization of  $R^*(a)$  is indeed related to the maximization of  $\langle \vartheta, a \rangle - \Lambda_\nu(\vartheta)$  subject to  $\vartheta \in \mathcal{E}$ , where  $a \notin \text{aff}(\nu)$  unless  $a = c$ . These observations are not particular to this simple example but rather illustrate general phenomena studied in this paper.

For an i.i.d. sample of any size from any pm in the family  $\mathcal{E}_\mu$ , the probability is positive that exactly one element differs from  $c$ , thus the sample mean  $a$  belongs to  $A = \partial C \setminus \{c\}$  with a positive probability. Example 5.5 of Sect. 5 exposes further illustrations of results of this paper for the full family  $\mathcal{E}_\mu$ .

1.3 The main concepts and notations are introduced in Sect. 2 including tools from convex geometry. Several auxiliary results are also collected there. Section 3 deals with the case when  $a$  belongs to a certain convex subset of  $\text{dom}(\Psi^*)$ . As will be clear later, this set contains the relative interior of  $\text{dom}(\Psi^*)$  and is exactly the set of those  $a \in \text{dom}(\Psi^*)$  for which  $R^*(a)$  belongs to  $\mathcal{E}$ , see Corollary 4.2. For  $a$  in this subset, the supremum in (1) is actually a maximum, perhaps not subject to  $\vartheta \in \mathcal{E}$  but to  $\vartheta \in \mathcal{E}_{\mu, a}$  where the latter set is the closure of a projection of  $\mathcal{E}$  such that the values  $\langle \vartheta, a \rangle - \Lambda(\vartheta)$ ,  $\vartheta \in \text{dom}(\Lambda)$ , are preserved by this projection. While this section is still of preparatory character, it has also independent interest because results well known under some regularity conditions [1, Sect. 9.4], are proved here in full generality. Much of the effort is needed to cover also the case  $a \notin \text{aff}(\mu)$ , disregarded in previous literature because of no immediate statistical meaning.

In Sect. 4, a pm satisfying (2), thus the GMLE from  $a \in \text{dom}(\Psi^*)$ , is constructed explicitly in Theorem 4.1, in contrast to [6] where  $R^*(a)$  was obtained as the limit of a Cauchy sequence in a complete space. In the case not covered in Sect. 3, a proper face of the set called convex core of  $\mu$  is identified by geometric means. For the family  $\mathcal{E}_{\nu, \mathcal{E}}$  determined by the restriction  $\nu$  of  $\mu$  to the closure of this face, results of Sect. 3 apply to obtain  $R_{\nu, \mathcal{E}}^*(a)$ . This equals the desired GMLE  $R_{\mu, \mathcal{E}}^*(a)$ . Relying on this construction, also descriptions of  $\text{dom}(\Psi^*)$  are presented in Theorem 4.9.

In Sect. 5, several properties of the GMLE mapping  $R^* : a \mapsto R^*(a)$ , defined on  $\text{dom}(\Psi^*)$ , are proved. Its range is characterized as a subset of  $cl_\nu(\mathcal{E}_\mathcal{E})$ , Theorem 5.1. In the nondegenerate case, this subset consists exactly of those pm's that have a mean. The inverse of the GMLE mapping is also addressed in Theorem 5.1. For each pm  $P$  in the range, the set  $\{a \in \text{dom}(\Psi^*) : R^*(a) = P\}$  is shown to be the singleton  $\{m(P)\}$  consisting of the partial mean  $m(P)$  of  $P$  or a (not necessarily convex) cone shifted

by  $m(P)$ . A general result about the continuity of the GMLE mapping is presented, Theorem 5.6, and special situations are discussed where the GMLE mapping is bijective or even a homeomorphism, Remark 5.9. Then, a relationship of the concepts of GMLE and generalized  $rI$ -projection is discussed. Finally, the log-convexity of  $cl_v(\mathcal{E}_{\mathcal{E}})$  and of the range of  $R^*$  is established, Theorem 5.10, in the sense of [6].

In Sect. 6, an MLE in  $cl_v(\mathcal{E}_{\mathcal{E}})$  is considered. This MLE concept extends that in Sect. 1.1 when  $a \in \text{aff}(\mu)$ , but is undefined otherwise. It is shown that if this MLE exists then it corresponds to the GMLE, but not conversely. In the special case  $\text{dom}(\Lambda) = \mathbb{R}^d$ , the GMLE and MLE coincide, and the results of this paper on GMLE's cover previous ones on nonexistence of MLE's. Finally, in Appendix the invariance of the concepts of this paper w.r.t. different parameterizations of a family is treated. The GMLE  $R^*(a)$  is shown not to depend on the parametrization when  $a \in \text{aff}(\mu)$ , though this invariance breaks down in other cases without statistical interpretation.

1.4 Exponential families are covered in depth in the monographs Chentsov [3], Barndorff-Nielsen [1], Brown [2] and Letac [11]. The standard reference for convexity is Rockafellar [15].

The ML estimation in canonically convex exponential families is treated in detail in [1, Sect. 9.4]. Generalizations of MLE's for the case when the supremum in (1) is not attained, have been considered for exponential families of pm's concentrated on finite or countable sets. The former case, in a theoretical sense, is completely settled in [1, Theorem 9.16], replacing the family by its "completion"; for recent related works addressing also computational issues in the framework of contingency tables see [10, 13]. The latter case is treated in [2, Chap. 6] via "aggregate families", under rather restrictive additional conditions. In general, one natural "completion" is the closure in variation distance. Some results on MLE in that direction, slightly improved here in Sect. 6, appeared in Csiszár and Matúš [6, Sect. 6]. The closure in variation distance of any canonically convex exponential family is characterized in Csiszár and Matúš [8], using the concepts of convex core of a measure introduced in [4], and of its accessible faces introduced in [8]. The log-convexity of this closure has not been yet addressed, but its subset consisting of all pm's to which some sequence in  $\mathcal{E}_{\mathcal{E}}$  converges in  $rI$ -divergence (containing all GMLE's) has been shown to be not necessarily log-convex, Csiszár and Matúš [7].

A description of  $\text{dom}(\Psi^*)$ , even of  $\text{dom}(\Lambda^*)$ , has been elusive in general, though a relationship between them has been known under some conditions [1, p. 159, Eq. (3)]. The strongest previous results on  $\text{dom}(\Lambda^*)$  are apparently those in [6, Proposition 1].

The concept of GMLE, the main subject of this paper, is introduced in [5, 6], motivated by the study of generalized  $rI$ -projections. The nonconstructive existence proof of GMLE given in [6] is extended to infinite dimensional exponential families in [9]. A constructive proof of the existence of GMLE is first given in this paper.

One application area of GMLE's seems to be within graphical models of multivariate statistics, see [12, Chap. 6] where MLE's have been considered for various general exponential families which reduce by sufficiency to standard ones considered in this paper. In such general mixed models, nonexistence of MLE's is a rather common phenomenon, and GMLE's may provide a remedy. In the particular case of log-linear

models for multidimensional contingency tables computational strategies under the nonexistence of MLE are discussed in [14].

## 2 Preliminaries

2.1 Denote for a subset  $B$  of  $\mathbb{R}^d$  its closure by  $cl(B)$ , affine hull by  $aff(B)$ , relative interior by  $ri(B)$ , which is the interior of  $B$  in the topology of  $aff(B)$ , and denote the shift of  $aff(B)$  containing the origin by  $lin(B)$ , which is the linear space spanned by the differences  $b - c, b, c \in B$ . The orthogonal projector to  $lin(B)$  is denoted by  $\pi_B$ . The orthogonal complement of a linear subspace  $E$  of  $\mathbb{R}^d$  is denoted by  $E^\perp$ .

A face of a convex set  $C \subseteq \mathbb{R}^d$  is a nonempty convex subset  $F$  that contains  $ta + (1 - t)b$  for some  $a, b$  in  $C$  and  $0 < t < 1$  only if  $a$  and  $b$  are in  $F$ ; note the slight deviation from the terminology of [15] where also the empty set is a face.

**Lemma 2.1** *For convex subsets  $C$  and  $K$  of  $\mathbb{R}^d$  and  $a \in C + K$ , in the family of faces  $F$  of  $C$  with  $a \in ri(F) + K$  the inclusion-largest element exists.*

*Proof* If  $a$  is in  $ri(F_i) + K$  for faces  $F_i$  of  $C, i = 1, 2$ , thus  $a = f_i + k_i$  with  $f_i \in ri(F_i)$  and  $k_i \in K$ , then the element  $g = \frac{1}{2}(f_1 + f_2)$  of  $C$  belongs to  $ri(G)$  for a unique face  $G$  of  $C$  [15, Theorem 18.2]. Since  $a = g + \frac{1}{2}(k_1 + k_2)$  it follows from convexity of  $K$  that  $a \in ri(G) + K$ . The definition of face implies that the segment with endpoints  $f_1$  and  $f_2$  is contained in  $G$ . Hence, both  $ri(F_i)$  intersect  $G$ , and thus are contained in  $G$  [15, Theorem 18.1]. This implies that a face in the family with the largest dimension contains all faces of the family. □

A subset of  $\mathbb{R}^d$  is a cone if it contains  $tx$  with each of its elements  $x$  and  $t > 0$ . For a convex set  $C \subseteq \mathbb{R}^d$ , its recession cone  $rec(C)$  is the set of all  $y \in \mathbb{R}^d$  such that  $x + ty \in C$  for all  $x \in C$  and  $t \geq 0$ . In [15], the notation  $0^+C$  is used instead of  $rec(C)$ .

**Lemma 2.2** *If  $T$  is a linear mapping and  $C \subseteq \mathbb{R}^d$  is convex then  $T rec(C) \subseteq rec(TC)$ .*

*Proof* Given  $a \in T rec(C)$  and  $b \in TC$ , write  $a = Tc_a, b = Tc_b$  where  $c_a \in rec(C), c_b \in C$ . Then  $c_b + tc_a \in C$  for  $t \geq 0$ . Applying  $T, b + ta \in TC$ , and thus  $a \in rec(TC)$ . □

The barrier cone  $bar(\Gamma)$  of a set  $\Gamma \subseteq \mathbb{R}^d$  is the set of all  $x \in \mathbb{R}^d$  such that the mapping  $\vartheta \mapsto \langle \vartheta, x \rangle$  is bounded above on  $\Gamma$ . A vector  $x \in \mathbb{R}^d$  is normal to the set  $\Gamma$  at  $\theta \in \mathbb{R}^d$  if  $\langle \vartheta - \theta, x \rangle \leq 0$  for all  $\vartheta \in \Gamma$ . Such a vector obviously belongs to  $bar(\Gamma)$ . The set of those normal vectors is the normal cone of  $\Gamma$  at  $\theta$ , denoted by  $N_\theta(\Gamma)$ . In this paper, these terms are used in an extended sense, as the usual requirement  $\theta \in \Gamma$  is not imposed.

**Lemma 2.3** *For any nonempty convex  $\Gamma \subseteq \mathbb{R}^d$ ,*

$$rec(\Gamma) \subseteq rec(cl(\Gamma)) = rec(ri(\Gamma)) = N_0(bar(\Gamma)).$$

*Proof* The recession cone of  $cl(\Gamma)$  contains that of  $\Gamma$  by [15, Theorem 8.3] and equals that of  $ri(\Gamma)$  by [15, Corollary 8.3.1]. By [15, Corollary 14.2.1], where the normal cone at 0 of a cone  $K$  is called the polar of  $K$ , it holds that  $rec(cl(\Gamma)) = N_0(bar(cl(\Gamma)))$ . This proves the last equality as  $bar(cl(\Gamma)) = bar(\Gamma)$  is obvious by definition.  $\square$

2.2 A face  $F$  of a convex set  $C \subseteq \mathbb{R}^d$  is *exposed* if either  $F = C$  or, for some unit vector  $\tau$ , the maximum of  $\langle \tau, x \rangle$  subject to  $x \in C$  is attained if and only if  $x \in F$ . Such  $\tau$  is said to *expose*  $F$  in  $C$ .

The following concepts were introduced in [8] in order to characterize the variation closure of  $\mathcal{E}_{\mathcal{E}}$ . An *access sequence* to a proper face  $F$  of a convex set  $C \subseteq \mathbb{R}^d$  is an orthonormal sequence  $\tau_1, \dots, \tau_m$  such that  $\tau_i \in lin(F_{i-1})$  exposes a face  $F_i$  of  $F_{i-1}$  for  $i = 1, \dots, m$ , where  $F_0 = C$  and  $F_m = F$ . Such a sequence always exists. In particular, an access sequence of length  $m = 1$  to a proper face  $F$  exists if and only if  $F$  is an exposed face. An access sequence to  $F$  is *adapted* to a convex set  $\mathcal{E} \subseteq \mathbb{R}^d$  if  $\tau_i \in rec(\pi_{F_{i-1}}(ri(\mathcal{E})))$  for  $1 \leq i \leq m$ . The access sequence to  $F = C$  is empty, by definition, and it is adapted to any convex  $\mathcal{E}$ . Note that the notion of adaptedness depends on  $\mathcal{E}$  only through  $ri(\mathcal{E})$  or  $\pi_C(\mathcal{E})$  or  $\pi_C(ri(\mathcal{E})) = ri(\pi_C(\mathcal{E}))$ . A face  $F$  of  $C$  is  $\mathcal{E}$ -*accessible* if there exists an access sequence to  $F$  which is adapted to  $\mathcal{E}$ . More details and examples illustrating these notions can be found in [8].

**Lemma 2.4** *For intersecting faces  $F$  and  $G$  of a convex set  $C$ , if  $F$  is a  $\mathcal{E}$ -accessible face of  $C$  then  $F \cap G$  is a  $\mathcal{E}$ -accessible face of  $G$ .*

*Proof* The case  $F \supseteq G$  is trivial. Otherwise,  $F \neq C$  and there exists an access sequence  $\tau_1, \dots, \tau_m$  to the face  $F$  of  $C$  which is adapted to  $\mathcal{E}$ . Let  $C = F_0 \supset \dots \supset F_m = F$  be the corresponding chain of faces, in particular,  $\tau_i \in rec(\pi_{F_{i-1}}(ri(\mathcal{E})))$ ,  $1 \leq i \leq m$ .

The sequence  $G_i = F_i \cap G$  of faces of  $G$  decreases from  $G$  to  $F \cap G$  and  $\tau_i$  exposes  $G_i$  in  $G_{i-1}$ . If  $G_{i-1}$  contains  $G_i$  strictly then  $\tau_i$  is not orthogonal to  $lin(G_{i-1})$  and the unit vector  $\vartheta_i$  in the direction  $\pi_{G_{i-1}}(\tau_i)$  exposes  $G_i$  in  $G_{i-1}$  as well. This  $\vartheta_i$  belongs to  $\pi_{G_{i-1}}(rec(\pi_{F_{i-1}}(ri(\mathcal{E}))))$ , thus to  $rec(\pi_{G_{i-1}}(ri(\mathcal{E})))$  by Lemma 2.2 and  $\pi_{G_{i-1}}\pi_{F_{i-1}} = \pi_{G_{i-1}}$ .

Let  $i_1 < \dots < i_k$  be those indices  $1 \leq i \leq m$  for which  $G_i$  is a proper face of  $G_{i-1}$ . For  $1 \leq j \leq k$  the face  $G_{i_{j-1}}$  of  $G$  equals  $G_{i_{j-1}}$ , and hence  $\vartheta_{i_j}$  exposes  $G_{i_j}$  in  $G_{i_{j-1}}$  and belongs to  $rec(\pi_{G'}(ri(\mathcal{E})))$  where  $G' = G_{i_{j-1}}$ . Thus,  $\vartheta_1, \dots, \vartheta_k$  is an access sequence to  $F \cap G$  from  $G$ , adapted to  $\mathcal{E}$ .  $\square$

**Corollary 2.5** *If  $F$  and  $G$  are intersecting  $\mathcal{E}$ -accessible faces of a convex set  $C$  then  $F \cap G$  is also a  $\mathcal{E}$ -accessible face of  $C$ .*

*Proof* Since  $G$  is  $\mathcal{E}$ -accessible from  $C$  and  $F \cap G$  is  $\mathcal{E}$ -accessible from  $G$  by Lemma 2.4 it suffices to concatenate their access sequences adapted to  $\mathcal{E}$ . The concatenation is an access sequence to the face  $F \cap G$  of  $C$ , adapted to  $\mathcal{E}$ .  $\square$

2.3 For a  $\sigma$ -finite Borel measure  $\mu$  on  $\mathbb{R}^d$ , its *convex support*  $cs(\mu)$  and *convex core*  $cc(\mu)$  are defined as intersections of those convex sets  $C \subseteq \mathbb{R}^d$  that are  $\mu$ -full,



$\mu(\mathbb{R}^d \setminus C) = 0$ , and closed, respectively Borel. The cores have been considered previously only for finite measures. Since  $c\mathcal{S}(\mu)$  and  $cc(\mu)$  do not change when  $\mu$  is replaced by any equivalent measure, results on finite measures immediately extend to  $\sigma$ -finite ones (provided that when images of  $\sigma$ -finite measures are considered, these are postulated to be  $\sigma$ -finite). References are made below to the following assertions.

**Fact 2.6** *The closure of  $cc(\mu)$  equals  $c\mathcal{S}(\mu)$  [4, Lemma 1], thus  $cc(\mu)$  is nonempty if and only if  $\mu$  is nonzero.*

**Fact 2.7** *Every face  $F$  of  $cc(\mu)$  equals the convex core of the restriction of  $\mu$  to  $cl(F)$  [4, Lemma 3], hence this restriction is nonzero.*

**Fact 2.8** *The image of  $cc(\mu)$  under any affine transformation  $T$  equals the convex core of the image of  $\mu$  under  $T$  [4, Lemma 8], providing that image is  $\sigma$ -finite.*

**Fact 2.9** *A supporting hyperplane  $H$  of  $c\mathcal{S}(\mu)$  has positive  $\mu$ -measure if and only if  $F = H \cap cc(\mu)$  is nonempty, in which case  $\mu(H \setminus cl(F)) = 0$  [8, Lemma 1].*

**Fact 2.10** *If  $F$  and  $G$  are faces of  $cc(\mu)$  then  $cl(F) \cap cl(G) \setminus cl(F \cap G)$  has  $\mu$ -measure zero [4, Corollary 4].*

On account of Fact 2.6,  $cc(\mu)$  and  $c\mathcal{S}(\mu)$  have the same relative interior, denoted by  $ri(\mu)$ . The notations  $aff(\mu)$ , already used in Sect. 1, and  $lin(\mu)$  are analogous and selfexplaining.

For  $a \in \mathbb{R}^d$  let  $E_{\mu,a}$  denote the subspace of  $\mathbb{R}^d$  spanned by the set  $\{b - a : b \in aff(\mu)\}$ , or equivalently by  $lin(\mu)$  and  $b - a$  for any fixed  $b \in aff(\mu)$ . Thus,  $E_{\mu,a} \supseteq lin(\mu)$ , with equality if and only if  $a \in aff(\mu)$ . Let  $\pi_\mu$  and  $\pi_{\mu,a}$  denote the orthogonal projectors to  $lin(\mu)$  and  $E_{\mu,a}$ , respectively.

**Lemma 2.11** *Each  $\vartheta \in \mathbb{R}^d$  satisfies  $\langle \vartheta, a \rangle - \Lambda(\vartheta) = \langle \pi_{\mu,a}(\vartheta), a \rangle - \Lambda(\pi_{\mu,a}(\vartheta))$ .*

*Proof* If  $x \in aff(\mu)$  then  $x - a \in E_{\mu,a}$ , thus  $\langle \vartheta, x - a \rangle = \langle \pi_{\mu,a}(\vartheta), x - a \rangle$ . Hence,

$$\langle \vartheta, a \rangle - \Lambda(\vartheta) = - \ln \int_{aff(\mu)} e^{\langle \vartheta, x - a \rangle} \mu(dx)$$

does not change when  $\vartheta$  is replaced by  $\pi_{\mu,a}(\vartheta)$ . □

Recall the standing convention that  $\mathcal{E}$  denotes a nonempty convex subset of  $dom(\Lambda)$  and  $\Psi^* = \Psi_{\mu,\mathcal{E}}^*$  is defined by (1).

**Corollary 2.12**  $\Psi_{\mu,\mathcal{E}}^*(a) = \Psi_{\mu,\pi_{\mu,a}(\mathcal{E})}^*(a)$ . *If  $\Gamma \subseteq dom(\Lambda)$  and  $\pi_{\mu,a}(\Gamma) = \pi_{\mu,a}(\mathcal{E})$  then  $\Psi_{\mu,\Gamma}^*(a) = \Psi_{\mu,\mathcal{E}}^*(a)$ .*

The special instance of Lemma 2.11

$$\langle \vartheta, b \rangle - \Lambda(\vartheta) = \langle \pi_\mu(\vartheta), b \rangle - \Lambda(\pi_\mu(\vartheta)), \quad \vartheta \in \mathbb{R}^d, b \in aff(\mu), \tag{3}$$

is well-known. Its consequences include the facts that  $Q_\vartheta$  with  $\vartheta \in dom(\Lambda)$  does not change when  $\vartheta$  is replaced by  $\pi_\mu(\vartheta)$ , and  $dom(\Lambda) = \pi_\mu(dom(\Lambda)) + lin(\mu)^\perp$ .

**Lemma 2.13**  $\Psi_{\mu,ri(\Xi)}^*(a) = \Psi_{\mu,\Xi}^*(a) = \Psi_{\mu,cl(\Xi)\cap dom(\Lambda)}^*(a)$ .

*Proof* If  $\theta \in cl(\Xi) \cap dom(\Lambda)$  then some sequence  $\vartheta_n$  in  $ri(\Xi)$  converges to  $\theta$  along a segment. By continuity of  $\Lambda$  along the segment,  $\Psi_{\mu,ri(\Xi)}^*(a) \geq \langle \theta, a \rangle - \Lambda(\theta)$ . Hence,  $\Psi_{\mu,ri(\Xi)}^*(a) \geq \Psi_{\mu,cl(\Xi)\cap dom(\Lambda)}^*(a)$ . The remaining inequalities are trivial.  $\square$

2.4 For any pm  $P$  on  $\mathbb{R}^d$ , write

$$M(P) = \{ \tau \in \mathbb{R}^d : x \mapsto \langle \tau, x \rangle \text{ is } P\text{-integrable} \}$$

and define the *partial mean*  $m(P)$  of  $P$  as the unique element of the linear space  $M(P)$  that satisfies

$$\int_{\mathbb{R}^d} \langle \tau, x \rangle P(dx) = \langle \tau, m(P) \rangle, \quad \tau \in M(P).$$

Note that  $M(P) = \mathbb{R}^d$  if and only if  $P$  has a mean, in which case  $m(P)$  equals the mean of  $P$ .

**Lemma 2.14** For any pm  $P$  on  $\mathbb{R}^d$ ,  $lin(P)^\perp$  is a subspace of  $M(P)$ .

*Proof* If  $\tau \in lin(P)^\perp$  then  $x \mapsto \langle \tau, x \rangle$  is constant on  $aff(P)$ , and thus  $\tau \in M(P)$ .  $\square$

**Lemma 2.15** The partial mean  $m(P)$  of a pm  $P$  belongs to  $ri(P) + M(P)^\perp$ , contained in  $aff(P)$ .

*Proof* Let  $\pi$  denote, in this proof only, the orthogonal projector to  $M(P)$ , and  $\pi P$  the image of  $P$  under  $\pi$ . Then for  $\vartheta \in M(P)$

$$\int_{\mathbb{R}^d} \langle \vartheta, x \rangle (\pi P)(dx) = \int_{\mathbb{R}^d} \langle \vartheta, \pi(x) \rangle P(dx) = \int_{\mathbb{R}^d} \langle \vartheta, x \rangle P(dx) = \langle \vartheta, m(P) \rangle$$

while for  $\vartheta \in M(P)^\perp$  clearly

$$\int_{\mathbb{R}^d} \langle \vartheta, x \rangle (\pi P)(dx) = 0 = \langle \vartheta, m(P) \rangle.$$

Hence,  $M(\pi P) = \mathbb{R}^d$  and the mean of  $\pi P$  exists and equals  $m(P)$ . It follows that  $m(P)$  belongs to the relative interior of  $cs(\pi P)$ , and then by Fact 2.6 to

$$ri(cc(\pi P)) = ri(\pi cc(P)) = \pi(ri(cc(P))).$$

Here, the interchangeability of  $cc$  and the projector holds by Fact 2.8 and that of relative interior and projectors is obvious. This implies that  $m(P)$  is in  $\pi(ri(P))$ , and thus in  $ri(P) + M(P)^\perp$ . This is a subset of  $aff(P)$  because  $ri(P) \subseteq aff(P)$  and  $M(P)^\perp \subseteq lin(P)$ , the latter by Lemma 2.14.  $\square$

The relevance of partial means for exponential families is indicated by the following simple lemma, particularly by its special instance  $P = Q_\theta$ .

**Lemma 2.16** *For  $\vartheta$  and  $\theta$  in  $\text{dom}(\Lambda)$  and a pm  $P$  with  $D(P\|Q_\theta) < \infty$ , the necessary and sufficient condition for  $D(P\|Q_\vartheta) < \infty$  is  $\vartheta - \theta \in M(P)$ . Under that condition*

$$D(P\|Q_\vartheta) - D(P\|Q_\theta) = \langle \theta - \vartheta, m(P) \rangle - \Lambda(\theta) + \Lambda(\vartheta). \tag{4}$$

*Proof* The left-hand side of (4) is, by definition, the difference of two integrals. If one of them is finite, the difference can be written as one integral which is equal to  $\int \langle \theta - \vartheta, x \rangle P(dx) - \Lambda(\theta) + \Lambda(\vartheta)$ . This is finite if and only if  $\vartheta - \theta \in M(P)$ , in which case it equals the right-hand side of (4).  $\square$

**Corollary 2.17** *The I-divergence  $D(P\|Q_\vartheta)$  is finite for each  $\vartheta \in \mathcal{E}$  if and only if it is finite for one  $\vartheta \in \mathcal{E}$  and  $\text{lin}(\mathcal{E}) \subseteq M(P)$ .*

The instance  $P = Q_\theta$  of Lemma 2.16 is that the I-divergence  $D(Q_\theta\|Q_\vartheta)$  is finite if and only if  $\vartheta - \theta$  belongs to  $M(Q_\theta)$ , in which case

$$D(Q_\theta\|Q_\vartheta) = \langle \theta - \vartheta, m(Q_\theta) \rangle - \Lambda(\theta) + \Lambda(\vartheta). \tag{5}$$

The corresponding instance of Corollary 2.17 is that if  $\theta \in \mathcal{E}$  then  $\text{lin}(\mathcal{E}) \subseteq M(Q_\theta)$  is equivalent to the finiteness of  $D(Q_\theta\|Q_\vartheta)$  for each  $\vartheta \in \mathcal{E}$ . To extend this equivalence to  $\theta$  beyond  $\mathcal{E}$ , let

$$\mathcal{E}_\mu = \text{cl}(\pi_\mu(\mathcal{E})) \cap \text{dom}(\Lambda_\mu) \quad \text{and} \quad \tilde{\mathcal{E}}_\mu = \mathcal{E}_\mu + \text{lin}(\mu)^\perp,$$

and observe that  $\mathcal{E} \subseteq \tilde{\mathcal{E}}_\mu$ .

**Lemma 2.18** *For  $\theta \in \tilde{\mathcal{E}}_\mu$ , the inclusion  $\text{lin}(\mathcal{E}) \subseteq M(Q_\theta)$  is equivalent to  $\vartheta - \theta \in M(Q_\theta)$  for each  $\vartheta \in \mathcal{E}$ , or to the finiteness of  $D(Q_\theta\|Q_\vartheta)$  for each  $\vartheta \in \mathcal{E}$ .*

*Proof* Since clearly  $\text{lin}(\mathcal{E}_\mu) = \text{lin}(\pi_\mu(\mathcal{E}))$ , and  $M(Q_\theta)$  contains  $\text{lin}(Q_\theta)^\perp = \text{lin}(\mu)^\perp$  by Lemma 2.14, the condition  $\text{lin}(\mathcal{E}) \subseteq M(Q_\theta)$  is equivalent to  $\text{lin}(\mathcal{E}_\mu) \subseteq M(Q_\theta)$ . On account of  $\pi_\mu(\theta) \in \mathcal{E}_\mu$  the latter is, in turn, equivalent to  $\pi_\mu(\vartheta) - \pi_\mu(\theta) \in M(Q_\theta)$ ,  $\vartheta \in \mathcal{E}$ , and using Lemma 2.14 again, to  $\vartheta - \theta \in M(Q_\theta)$ ,  $\vartheta \in \mathcal{E}$ . The second equivalence follows from Lemma 2.16 with  $P = Q_\theta$ .  $\square$

**Lemma 2.19** *If  $\mathcal{E}$  intersects  $\text{ri}(\text{dom}(\Lambda))$  and  $\theta \in \tilde{\mathcal{E}}_\mu$  then  $M(Q_\theta)$  contains  $\text{lin}(\mathcal{E})$  if and only if it contains  $\text{lin}(\text{dom}(\Lambda))$ .*

*Proof* By Lemma 2.18, the second assumption and  $\text{lin}(\mathcal{E}) \subseteq M(Q_\theta)$  imply that  $\mathcal{E}$  is contained in  $[\theta + M(Q_\theta)] \cap \text{dom}(\Lambda)$ . This set is a face of  $\text{dom}(\Lambda)$  due to Lemma 2.21 below. The first assumption implies that this face cannot be proper. Hence,  $\theta + M(Q_\theta)$  contains  $\text{dom}(\Lambda)$ , and thus the assertion follows.  $\square$

**Corollary 2.20** *In the nondegenerate case, when  $\mathcal{E}$  intersects the interior of  $\text{dom}(\Lambda)$ , a pm  $Q_\theta$  with  $\theta \in \tilde{\mathcal{E}}_\mu$  has a mean if and only if  $\text{lin}(\mathcal{E}) \subseteq M(Q_\theta)$ .*

The instance  $\mathcal{E} = \{\theta\}$  of this corollary covers the well-known fact that the pm  $Q_\theta$  indexed by an interior point  $\theta$  of  $dom(\Lambda)$  has a mean.

**Lemma 2.21** *If  $\theta \in dom(\Lambda)$  then  $[\theta + M(Q_\theta)] \cap dom(\Lambda)$  is a face of  $dom(\Lambda)$ .*

*Proof* For  $\vartheta, \tau \in dom(\Lambda)$  and  $0 < t < 1$

$$\begin{aligned} & t \int_{\mathbb{R}^d} \langle \vartheta - \theta, x \rangle Q_\theta(dx) + (1 - t) \int_{\mathbb{R}^d} \langle \tau - \theta, x \rangle Q_\theta(dx) \\ &= \int_{\mathbb{R}^d} \langle t\vartheta + (1 - t)\tau - \theta, x \rangle Q_\theta(dx) \end{aligned}$$

where each integral is finite or  $-\infty$  by [8, Lemma 4]. If  $\Gamma = [\theta + M(Q_\theta)] \cap dom(\Lambda)$  contains  $t\vartheta + (1 - t)\tau$  then the integral on the right is finite, and hence all integrals are finite. Thus,  $\vartheta, \tau \in \Gamma$ , and  $\Gamma$  is a face. □

2.5 The *subdifferential* at  $\theta \in \mathbb{R}^d$  of a convex function  $f$  on  $\mathbb{R}^d$  is the set  $\partial f(\theta)$  consisting of those  $x \in \mathbb{R}^d$  that satisfy  $f(\vartheta) \geq f(\theta) + \langle \vartheta - \theta, x \rangle$  for all  $\vartheta \in \mathbb{R}^d$  or, equivalently,  $f^*(x) = \langle \theta, x \rangle - f(\theta)$ , where  $f^*$  is the convex conjugate of  $f$ . Subdifferentials appear in this paper only when comparing results with previous ones. In the following lemma,  $\Psi$  denotes the function defined in Remark 1.1 that equals  $\Lambda$  on the convex nonempty subset  $\mathcal{E}$  of  $dom(\Lambda)$  and  $+\infty$  elsewhere.

**Lemma 2.22** *If  $\mathcal{E}$  intersects  $ri(dom(\Lambda))$  then  $dom(\Psi^*) = dom(\Lambda^*) + bar(\mathcal{E})$  and*

$$\partial\Psi(\theta) = \begin{cases} \partial\Lambda(\theta) + N_\theta(\mathcal{E}), & \theta \in \mathcal{E}, \\ \emptyset, & \text{otherwise.} \end{cases}$$

*Remark 2.23* Analogous results appear in [1, Sect. 9.4], see (3) and (4) there, under the assumption that  $\mathcal{E}$  and  $dom(\Lambda)$  have nonempty interiors. Later, the first equality will be shown to hold for any convex nonempty  $\mathcal{E} \subset dom(\Lambda)$ , see Remark 4.10 in Sect. 4.

*Proof* The function  $\Psi$  is the sum of  $\Lambda$  and the indicator function  $\delta(\cdot|\mathcal{E})$  of  $\mathcal{E}$ , equal to 0 on  $\mathcal{E}$  and  $+\infty$  elsewhere. As  $\mathcal{E}$  is contained in  $dom(\Lambda)$  but not in its relative boundary  $dom(\Lambda) \setminus ri(dom(\Lambda))$ , the relative interior of  $\mathcal{E}$  is contained in that of  $dom(\Lambda)$  [15, Corollary 6.5.2]. Then, the first equality follows from [15, Theorem 16.4] using that  $dom(\delta^*(\cdot|\mathcal{E}))$  equals  $bar(\mathcal{E})$ , and the second one from [15, Theorem 23.8] since the subdifferential of  $\delta(\cdot|\mathcal{E})$  at  $\theta$  equals  $N_\theta(\mathcal{E})$  if  $\theta \in \mathcal{E}$  and  $\emptyset$  otherwise. □

### 3 The GMLE when $a \in ri(\mu) + bar(\mathcal{E})$

3.1 Given a nonempty convex set  $\mathcal{E}$  contained in  $dom(\Lambda)$  and  $a \in \mathbb{R}^d$ , let

$$\mathcal{E}_{\mu,a} = cl(\pi_{\mu,a}(\mathcal{E})) \cap dom(\Lambda) \quad \text{and} \quad \tilde{\mathcal{E}}_{\mu,a} = \mathcal{E}_{\mu,a} + E_{\mu,a}^\perp,$$

see Sect. 2.3 for notation. These are convex subsets of  $dom(\Lambda)$ , and  $\mathcal{E} \subseteq \tilde{\mathcal{E}}_{\mu,a}$ . Then

$$\Psi_{\mu,\mathcal{E}}^*(a) = \Psi_{\mu,\mathcal{E}_{\mu,a}}^*(a) = \Psi_{\mu,\tilde{\mathcal{E}}_{\mu,a}}^*(a) \tag{6}$$

where the first equality follows from Corollary 2.12 and Lemma 2.13, and the second one from Corollary 2.12 applied to  $\tilde{\mathcal{E}}_{\mu,a}$  in the role of  $\mathcal{E}$ . If  $a \in aff(\mu)$  then  $\mathcal{E}_{\mu,a}$  is equal to  $\mathcal{E}_\mu = cl(\pi_\mu(\mathcal{E})) \cap dom(\Lambda)$  and  $\tilde{\mathcal{E}}_{\mu,a}$  to  $\tilde{\mathcal{E}}_\mu$ , see Sect. 2.4, while always  $\pi_\mu(\mathcal{E}_{\mu,a}) \subseteq \mathcal{E}_\mu$ , and hence  $\tilde{\mathcal{E}}_{\mu,a} \subseteq \tilde{\mathcal{E}}_\mu$ . In the sequel, notations introduced in Sect. 2 are used without reference.

The following theorem admits to define a mapping that assigns to each element  $a$  of  $ri(\mu) + bar(\mathcal{E})$  the parameter  $\theta^*(a) = \theta_{\mu,\mathcal{E}}^*(a)$  equal to the unique maximizer of the function  $\vartheta \mapsto \langle \vartheta, a \rangle - \Lambda(\vartheta)$  subject to  $\vartheta \in \mathcal{E}_{\mu,a}$ . Necessity of the assumption  $a \in ri(\mu) + bar(\mathcal{E})$  for the existence of such a maximizer is shown in Theorem 3.2.

**Theorem 3.1** *If  $a \in ri(\mu) + bar(\mathcal{E})$  then  $\Psi^*(a)$  is finite and a unique  $\theta \in \mathcal{E}_{\mu,a}$  exists such that  $\Psi^*(a) = \langle \theta, a \rangle - \Lambda(\theta)$ .*

*Proof* If  $a$  belongs to  $ri(\mu) + bar(\mathcal{E})$  then  $a - b \in bar(\mathcal{E})$  for some  $b \in ri(\mu)$ , thus

$$\langle \vartheta, a - b \rangle \leq r, \quad \vartheta \in \mathcal{E}, \tag{7}$$

for some  $r \in \mathbb{R}$ . By (3) and [8, Lemma 9] applied to this  $b$ , for some  $s > 0$  and  $t \in \mathbb{R}$

$$\langle \vartheta, b \rangle - \Lambda(\vartheta) = \langle \pi_\mu(\vartheta), b \rangle - \Lambda(\pi_\mu(\vartheta)) \leq t - s \|\pi_\mu(\vartheta)\|, \quad \vartheta \in \mathbb{R}^d. \tag{8}$$

Note that though Lemma 9 of [8] is formulated only for finite measures, its proof holds verbatim also for  $\sigma$ -finite measures  $\mu$  with  $dom(\Lambda_\mu)$  nonempty. Combining (7) and (8),

$$\langle \vartheta, a \rangle - \Lambda(\vartheta) - r \leq \langle \vartheta, b \rangle - \Lambda(\vartheta) \leq t - s \|\pi_\mu(\vartheta)\|, \quad \vartheta \in \mathcal{E}, \tag{9}$$

whence  $\Psi^*(a) \leq t + r$ .

Consider a sequence  $\vartheta_n$  in  $\mathcal{E}$  with  $\langle \vartheta_n, a \rangle - \Lambda(\vartheta_n)$  converging to the finite number  $\Psi^*(a)$ . By (8) and (7)

$$\langle \vartheta_n, a \rangle - \Lambda(\vartheta_n) - t \leq \langle \vartheta_n, a - b \rangle \leq r.$$

Thus the sequence  $\langle \vartheta_n, a - b \rangle$  is bounded, and by (9) the sequence  $\pi_\mu(\vartheta_n)$  is bounded, too. As  $lin(\mu)$  and  $b - a$  span  $E_{\mu,a}$ , it follows that the sequence  $\pi_{\mu,a}(\vartheta_n)$  is bounded. Hence, going to a subsequence if necessary,  $\pi_{\mu,a}(\vartheta_n)$  converges to some  $\theta \in cl(\pi_{\mu,a}(\mathcal{E}))$ . By Lemma 2.11,  $\langle \pi_{\mu,a}(\vartheta_n), a \rangle - \Lambda(\pi_{\mu,a}(\vartheta_n))$  converges to  $\Psi^*(a)$ . Then,  $\Psi^*(a) \leq \langle \theta, a \rangle - \Lambda(\theta)$  by semicontinuity of  $\Lambda$ . In particular,  $\theta \in dom(\Lambda)$ , thus  $\theta \in \mathcal{E}_{\mu,a}$ . The opposite inequality follows from (6).

To prove uniqueness, it suffices to show that  $\theta, \tau \in \mathcal{E}_{\mu,a}$  and

$$\Psi^*(a) = \langle \theta, a \rangle - \Lambda(\theta) = \langle \tau, a \rangle - \Lambda(\tau)$$

imply  $\theta = \tau$ . By (6), both  $\theta$  and  $\tau$  maximize the function  $\vartheta \mapsto \langle \vartheta, a \rangle - \Lambda(\vartheta)$  over  $\tilde{\mathcal{E}}_{\mu,a}$ , and hence  $\Lambda$  is not strictly convex on the segment connecting  $\theta$  and  $\tau$ . It follows from (3) and the strict convexity of  $\Lambda$  on  $\text{lin}(\mu)$  that  $\pi_\mu(\theta)$  and  $\pi_\mu(\tau)$  coincide. This and (3) imply that if  $b \in \text{aff}(\mu)$  then  $\langle \theta, b \rangle - \Lambda(\theta)$  equals  $\langle \tau, b \rangle - \Lambda(\tau)$ , and hence  $\langle \theta, a - b \rangle = \langle \tau, a - b \rangle$ . Since  $\theta, \tau \in \mathcal{E}_{\mu,a} \subseteq E_{\mu,a}$ , it follows that  $\theta = \tau$ .  $\square$

**Theorem 3.2** *For any  $a \in \mathbb{R}^d$  and  $\theta \in \tilde{\mathcal{E}}_{\mu,a}$  the following statements are equivalent.*

- (i)  $\Psi^*(a) = \langle \theta, a \rangle - \Lambda(\theta)$ .
- (ii)  $\text{lin}(\mathcal{E}) \subseteq M(Q_\theta)$  and  $a - m(Q_\theta)$  belongs to  $N_\theta(\mathcal{E})$ .
- (iii)  $a \in \text{ri}(\mu) + \text{bar}(\mathcal{E})$  and  $\theta^*(a) = \pi_{\mu,a}(\theta)$ .
- (iv)  $\langle \theta, a \rangle - \Lambda(\theta) \geq \langle \vartheta, a \rangle - \Lambda(\vartheta) + D(Q_\theta \| Q_\vartheta)$  for all  $\vartheta \in \mathcal{E}$ .

By Theorem 3.2, in case  $a \in \text{ri}(\mu) + \text{bar}(\mathcal{E})$  the GMLE  $R^*(a)$  exists and equals  $Q_{\theta^*(a)}$ . Indeed,  $\theta = \theta^*(a)$  satisfies (iii), then (i) and (iv) combine to the GMLE inequality (2) with  $R^*(a)$  replaced by  $Q_{\theta^*(a)}$ , and then  $R^*(a) = Q_{\theta^*(a)}$  by uniqueness of GMLE, discussed in Remark 1.2. This result is complemented in Sect. 4, where Corollary 4.2 shows that if  $a \notin \text{ri}(\mu) + \text{bar}(\mathcal{E})$  and the GMLE  $R^*(a)$  exists then  $R^*(a) \notin \mathcal{E}$ .

Note that in the nondegenerate case, the inclusion in (ii) holds if and only if  $Q_\theta$  has a mean, due to Corollary 2.20 and  $\tilde{\mathcal{E}}_{\mu,a} \subseteq \tilde{\mathcal{E}}_\mu$ . If, in addition,  $\theta$  is in the interior of  $\mathcal{E}$  then the normal cone  $N_\theta(\mathcal{E})$  is the singleton  $\{0\}$ , and the condition (ii) requires the mean of  $Q_\theta$  to equal  $a$ .

*Proof* (i) $\Rightarrow$ (ii): On account of [8, Lemma 4], if  $\theta \in \text{dom}(\Lambda)$  and  $\tau$  is a unit vector such that  $\theta + t\tau \in \text{dom}(\Lambda)$  for some  $t > 0$  then  $\int \langle \tau, x \rangle Q_\theta(dx)$  exists, either finite or  $-\infty$ , and is equal to the one-sided directional derivative of  $\Lambda$  at  $\theta$  in the direction  $\tau$ . The assumption (i) and (6) imply  $\langle \theta, a \rangle - \Lambda(\theta) \geq \langle \vartheta, a \rangle - \Lambda(\vartheta)$  for all  $\vartheta \in \tilde{\mathcal{E}}_{\mu,a}$ . Hence, the derivative of the function  $\vartheta \mapsto \langle \vartheta, a \rangle - \Lambda(\vartheta)$  at  $\theta$  in any direction pointing from  $\theta$  to some  $\vartheta \in \tilde{\mathcal{E}}_{\mu,a}$  is not positive. It follows that for all  $\vartheta \in \tilde{\mathcal{E}}_{\mu,a}$ , the integral  $\int \langle \vartheta - \theta, x \rangle Q_\theta(dx)$  is finite, thus  $\vartheta - \theta \in M(Q_\theta)$ , and

$$0 \geq \langle \vartheta - \theta, a \rangle - \int_{\mathbb{R}^d} \langle \vartheta - \theta, x \rangle Q_\theta(dx) = \langle \vartheta - \theta, a - m(Q_\theta) \rangle.$$

This means that  $\text{lin}(\tilde{\mathcal{E}}_{\mu,a})$  is contained in  $M(Q_\theta)$  and  $a - m(Q_\theta)$  is normal to  $\tilde{\mathcal{E}}_{\mu,a}$  at  $\theta$ . Since  $\mathcal{E} \subseteq \tilde{\mathcal{E}}_{\mu,a}$  the validity of (ii) follows.

(i), (ii) $\Rightarrow$ (iii): If  $a - m(Q_\theta)$  belongs to  $N_\theta(\mathcal{E}) \subseteq \text{bar}(\mathcal{E})$  then  $a \in m(Q_\theta) + \text{bar}(\mathcal{E})$ . The assumption  $\text{lin}(\mathcal{E}) \subseteq M(Q_\theta)$  implies, using Lemma 2.14, that  $M(Q_\theta)^\perp \subseteq \text{lin}(\mathcal{E})^\perp \subseteq \text{bar}(\mathcal{E})$ . By this and  $\text{ri}(Q_\theta) = \text{ri}(\mu)$ , Lemma 2.15 gives that  $m(Q_\theta) \in \text{ri}(\mu) + \text{bar}(\mathcal{E})$ . As  $\text{bar}(\mathcal{E})$  is a convex cone,  $a \in \text{ri}(\mu) + \text{bar}(\mathcal{E})$  follows. Then  $\theta^*(a) = \pi_{\mu,a}(\theta)$  is a consequence of (i) and Lemma 2.11, due to the uniqueness assertion of Theorem 3.1.

(iii) $\Rightarrow$ (i): By the definition of  $\theta^*(a)$ , (iii) implies that the equality in (i) holds for  $\pi_{\mu,a}(\theta)$  in the role of  $\theta$ . Then, by Lemma 2.11, it holds also for  $\theta$ .

(ii) $\Leftrightarrow$ (iv): Due to  $\theta \in \tilde{\mathcal{E}}_{\mu,a} \subseteq \tilde{\mathcal{E}}_\mu$  and Lemma 2.18, the inclusion in (ii) is equivalent to the finiteness of  $D(Q_\theta \| Q_\vartheta)$  for each  $\vartheta \in \mathcal{E}$ . Then the second condition in (ii), thus  $0 \geq \langle \vartheta - \theta, a - m(Q_\theta) \rangle, \vartheta \in \mathcal{E}$ , is equivalent to (iv) on account of (5).

(iv) $\Rightarrow$ (i): This obtains by taking supremum over  $\vartheta \in \mathcal{E}$  in (iv), where the  $I$ -divergence is nonnegative, and by (6).  $\square$

3.2 Theorem 3.2 makes it possible to describe the range of the mapping  $\theta^*$ , based on the observation that the necessary and sufficient conditions for  $a \in ri(\mu) + bar(\mathcal{E})$  and  $\theta^*(a) = \theta$  are

$$\theta \in \mathcal{E}_{\mu,a}, \quad lin(\mu) \subseteq M(Q_\theta) \quad \text{and} \quad a - m(Q_\theta) \in N_\theta(\mathcal{E}). \tag{10}$$

Let

$$\mathcal{E}_\mu^M = \{\theta \in \mathcal{E}_\mu : lin(\mathcal{E}) \subseteq M(Q_\theta)\},$$

let  $\mathbf{E}$  denote the family of linear subspaces  $E$  of  $\mathbb{R}^d$  that contain  $lin(\mu)$  as a subspace of codimension 1, and for  $E \in \mathbf{E}$  denote

$$\mathcal{E}_E^{M,\cap} = \{\theta \in cl(\pi_E(\mathcal{E})) \cap dom(\Lambda) : lin(\mathcal{E}) \subseteq M(Q_\theta) \text{ and } N_\theta(\mathcal{E}) \cap [E \setminus lin(\mu)] \neq \emptyset\}.$$

**Lemma 3.3** *The range  $\theta^*(ri(\mu) + bar(\mathcal{E}))$  of  $\theta^*$  is equal to  $\mathcal{E}_\mu^M \cup \bigcup_{E \in \mathbf{E}} [\mathcal{E}_E^{M,\cap} \setminus lin(\mu)]$ . For  $\theta$  in this range,  $\{a \in ri(\mu) + bar(\mathcal{E}) : \theta^*(a) = \theta\}$  equals  $m(Q_\theta) + K_\mu(\theta)$  where*

$$K_\mu(\theta) = N_\theta(\mathcal{E}) \cap \begin{cases} lin(\mu) \cup \bigcup \{E \in \mathbf{E} : \theta \in \mathcal{E}_E^{M,\cap}\}, & \theta \in \mathcal{E}_\mu^M, \\ E \setminus lin(\mu), & \theta \in \mathcal{E}_E^{M,\cap} \setminus lin(\mu), E \in \mathbf{E}. \end{cases}$$

*Proof* By Theorem 3.2,  $\theta \in \mathbb{R}^d$  is in the range if and only if (10) holds for this  $\theta$  and some  $a \in \mathbb{R}^d$ . Since  $m(Q_\theta) \in aff(\mu)$  by Lemma 2.15, the conditions (10) and  $a \in aff(\mu)$  are equivalent to the conditions

$$\theta \in \mathcal{E}_\mu^M \quad \text{and} \quad a - m(Q_\theta) \in N_\theta(\mathcal{E}) \cap lin(\mu), \tag{11}$$

using that  $a \in aff(\mu)$  implies  $\mathcal{E}_{\mu,a} = \mathcal{E}_\mu$ , while (10) and  $a \notin aff(\mu)$  are equivalent to

$$\theta \in \mathcal{E}_E^{M,\cap} \quad \text{and} \quad a - m(Q_\theta) \in N_\theta(\mathcal{E}) \cap [E \setminus lin(\mu)], \tag{12}$$

where  $E = E_{\mu,a}$ . Hence, the range is contained in  $\Gamma = \mathcal{E}_\mu^M \cup \bigcup_{E \in \mathbf{E}} \mathcal{E}_E^{M,\cap}$ .

If  $\theta \in \mathcal{E}_\mu^M$  then the conditions (11) are trivially met by  $a = m(Q_\theta)$ , and if  $\theta \in \mathcal{E}_E^{M,\cap}$  for some  $E \in \mathbf{E}$  then there exists  $a \in \mathbb{R}^d$  satisfying the conditions (12), by the definition of  $\mathcal{E}_E^{M,\cap}$ . Hence, the range contains  $\Gamma$ . The proof of the first assertion of the lemma is completed by noting that  $\mathcal{E}_E^{M,\cap} \setminus \mathcal{E}_\mu^M$  is equal to  $\mathcal{E}_E^{M,\cap} \setminus lin(\mu)$ , due to the obvious inclusions

$$\mathcal{E}_E^{M,\cap} \cap lin(\mu) \subseteq \pi_\mu(\mathcal{E}_E^{M,\cap}) \subseteq \mathcal{E}_\mu^M.$$

It remains to show that for  $\theta$  in the range of  $\theta^*$  just determined, the set of those  $a \in \mathbb{R}^d$  that satisfy either (11), or (12) for some  $E \in \mathbf{E}$ , is equal to  $m(Q_\theta) + K_\mu(\theta)$ .

If  $\theta \in \mathcal{E}_\mu^M$  then the set of all  $a \in \mathbb{R}^d$  satisfying (11) is equal to  $m(Q_\theta) + [N_\theta(\mathcal{E}) \cap \text{lin}(\mu)]$ , and the set of all  $a \in \mathbb{R}^d$  satisfying (12) for some  $E \in \mathcal{E}$  has a similar form, with  $\text{lin}(\mu)$  replaced by the union of  $E \setminus \text{lin}(\mu)$  over those linear spaces  $E \in \mathcal{E}$  for which  $\theta \in \mathcal{E}_E^{M,\cap}$ . This establishes the second assertion for the case  $\theta \in \mathcal{E}_\mu^M$ . If  $\theta$  is in the range but not in  $\mathcal{E}_\mu^M$  then it belongs to one of the mutually disjoint sets  $\mathcal{E}_E^{M,\cap} \setminus \text{lin}(\mu)$ ,  $E \in \mathcal{E}$ . Then no  $a \in \mathbb{R}^d$  satisfies (11), while (12) is satisfied by exactly those  $a$  that belong to  $m(Q_\theta) + K_\mu(\theta)$  with  $K_\mu(\theta)$  equal to  $N_\theta(\mathcal{E}) \cap [E \setminus \text{lin}(\mu)]$ .  $\square$

**Corollary 3.4** *The range of the mapping  $a \mapsto Q_{\theta^*(a)}$  defined on  $\text{ri}(\mu) + \text{bar}(\mathcal{E})$  coincides with the family  $\{Q_\tau : \tau \in \mathcal{E}_\mu^M\}$ . If  $\tau \in \mathcal{E}_\mu^M$  then the set  $\{a \in \text{ri}(\mu) + \text{bar}(\mathcal{E}) : Q_{\theta^*(a)} = Q_\tau\}$  equals  $m(Q_\tau) + K_\mu^*(\tau)$  where*

$$K_\mu^*(\tau) = \bigcup \{K_\mu(\theta) : \theta \in \theta^*(\text{ri}(\mu) + \text{bar}(\mathcal{E})) \text{ and } \pi_\mu(\theta) = \tau\}.$$

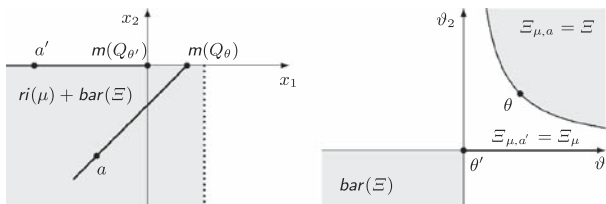
*Proof* Since  $Q_{\theta^*(a)} = Q_{\pi_\mu \theta^*(a)}$ , the first assertion follows by projecting the range of  $\theta^*$ , using that  $\pi_\mu(\mathcal{E}_E^{M,\cap}) \subseteq \mathcal{E}_\mu^M$ . The second assertion follows directly from Lemma 3.3.  $\square$

It is easy to see that the sets  $K_\mu(\theta)$  and  $K_\mu^*(\tau)$  are cones, both contained in  $\text{bar}(\mathcal{E})$ .

**Example 3.5** Let  $\mu$  be the measure on  $\mathbb{R}^2$  equal to the sum of the point masses at  $(1, 0)$  and  $(-1, 0)$ . Then  $\text{aff}(\mu) = \text{lin}(\mu)$  is the horizontal axis,  $\Lambda(\vartheta) = \ln(e^{\vartheta_1} + e^{-\vartheta_1})$  is finite for all  $\vartheta = (\vartheta_1, \vartheta_2) \in \mathbb{R}^2$ , and  $Q_\vartheta$  has the mean  $m(Q_\vartheta) = (\tanh(\vartheta_1), 0)$ .

(i) Let  $\mathcal{E} = \{\vartheta : \vartheta_1 > 0, \vartheta_1 \vartheta_2 \geq 1\}$ , see Fig. 2. If  $a = (a_1, a_2) \in \mathbb{R}^2$  then  $\mathcal{E}_{\mu,a}$  equals  $\mathcal{E}$  if  $a_2 \neq 0$  and the halfaxis  $\mathcal{E}_\mu = \{(\vartheta_1, 0) : \vartheta_1 \geq 0\}$ , otherwise. Further,  $\text{ri}(\mu) + \text{bar}(\mathcal{E})$  consists of all  $a$  with  $a_1 < 1, a_2 \leq 0$ . It is strictly contained in its closure equal to  $\text{dom}(\Psi^*)$ , and strictly contains its interior. Each  $a \in \text{ri}(\mu) + \text{bar}(\mathcal{E})$  with  $a_2 < 0$  can be represented as  $a = (\tanh(t) - r, -rt^2)$  with  $t, r > 0$ , and then for  $\theta = (t, t^{-1}) \in \mathcal{E}_{\mu,a}$  the vector  $a - m(Q_\theta) = -r(1, t^2)$  is normal to  $\mathcal{E}$  at  $\theta$ , thus  $\theta^*(a) = \theta$  by Theorem 3.2. If  $a = (a_1, 0), 0 < a_1 < 1$ , then  $a = m(Q_\theta)$  for  $\theta = (t, 0) \in \mathcal{E}_\mu$  with suitable  $t > 0$ , and (ii) of Theorem 3.2 trivially holds, thus  $\theta^*(a) = \theta$ . Finally, if  $a' = (a_1, 0), a_1 \leq 0$ , then  $\theta' = (0, 0)$  with  $m(Q_{\theta'}) = (0, 0)$  satisfies (ii) of Theorem 3.2, hence  $\theta^*(a') = \theta'$ . Note that while  $\theta'$  is in  $\mathcal{E}_{\mu,a'} = \mathcal{E}_\mu$ , it is not equal to the projection of any  $\vartheta \in \mathcal{E}$ , even though  $\mathcal{E}$  is a closed set.

(ii) Let  $\mathcal{E} = \{\vartheta : |\vartheta_2| < \vartheta_1\}$ , see Fig. 3. Then  $\mathcal{E}_{\mu,a}$  equals  $\text{cl}(\mathcal{E})$  if  $a_2 \neq 0$  and  $\mathcal{E}_\mu$ , the halfaxis as in (i), otherwise. Further,  $\text{ri}(\mu) + \text{bar}(\mathcal{E})$  consists of all  $a \in \mathbb{R}^2$



**Fig. 2** Illustration of Example 3.5 (i)



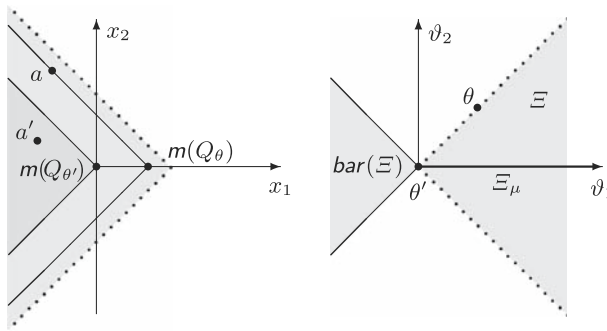


Fig. 3 Illustration of Example 3.5 (ii)

satisfying  $a_1 + |a_2| < 1$ . This is an open set, and  $dom(\Psi^*)$  equals its closure. If  $a_1 + |a_2| \leq 0$  then for  $\theta' = (0, 0)$  the vector  $a - m(Q_{\theta'}) = a$  is normal to  $\mathcal{E}$  at  $\theta'$ , hence  $\theta^*(a) = \theta'$ . If  $0 < a_1 + |a_2| < 1$  then  $a$  can be represented by some  $t, r > 0$  as  $(\tanh(t), 0)$  or  $(\tanh(t) - r, r)$  or  $(\tanh(t) - r, -r)$  according to  $a_2$  is 0, positive or negative, respectively. Then, for  $\theta$  equal to  $(t, 0)$ ,  $(t, t)$  or  $(t, -t)$  the vector  $a - m(Q_\theta)$  is equal to  $(0, 0)$ ,  $(-r, r)$  or  $(-r, -r)$ , respectively. In all cases the vector is normal to  $\mathcal{E}$  at  $\theta$ , and hence  $\theta^*(a) = \theta$ . It follows from Theorem 3.2 that the range of  $\theta^*$  is union of the three halflines  $\mathcal{E}_\mu$ ,  $\{(r, r) : r > 0\}$  and  $\{(r, -r) : r > 0\}$ , see also Lemma 3.3. It follows also that the set  $\{a : Q_{\theta^*(a)} = Q_\tau\}$  is the sum of  $m(Q_\tau)$  and the cone  $K_\mu^*(\tau)$  equal to the union of  $\{(0, 0)\}$  with the two halflines  $\{(-r, 0) : r > 0\}$  and  $\{(-r, -r) : r > 0\}$ , see also Corollary 3.4. The range of  $\theta^*$  and  $K_\mu^*(\tau)$  are not convex.

Note that while the canonically convex family  $\mathcal{E}_\mathcal{E}$  is the same in (i) and (ii), the sets  $ri(\mu) + bar(\mathcal{E})$  where  $\theta^* = \theta_{\mu, \mathcal{E}}^*$  is defined are different. Even for  $a$  in their intersection, not only  $\theta^*(a)$  but also the pm  $Q_{\theta^*(a)}$  does depend on the choice of the parameter set  $\mathcal{E}$  in general.

3.3 In the final part of this section, the above results are discussed and related to previous ones in the special cases  $\mathcal{E} = dom(\Lambda)$  and  $\mathcal{E} \subseteq lin(\mu)$ , sharing the feature that the set  $\mathcal{E}_{\mu, a}$  in the definition of  $\theta^*(a)$  does not depend on  $a$ . These cases cover those treated in the literature, where  $lin(\mu) = \mathbb{R}^d$  is assumed except for a few occasions addressing full families.

Consider first  $\mathcal{E} = dom(\Lambda)$ , the case of a full family. Then the set  $ri(\mu) + bar(\mathcal{E})$  on which the mapping  $\theta^*$  is defined is equal to  $ri(\mu)$ . This follows from

$$bar(dom(\Lambda)) \subseteq rec(dom(\Lambda^*)) \subseteq rec(ri(dom(\Lambda^*))) \quad \text{and} \quad ri(dom(\Lambda^*)) = ri(\mu) \tag{13}$$

where the first inclusion holds by [1, Theorem 5.19], the second one by Lemma 2.3, and for the last equality see for example [6, Proposition 1]. Further,  $\mathcal{E}_{\mu, a} = dom(\Lambda) \cap lin(\mu)$  for each  $a \in ri(\mu)$ , thus the mapping  $\theta^*$  ranges in this intersection. Using that  $\tilde{\mathcal{E}}_{\mu, a} = \mathcal{E} = dom(\Lambda)$  for each  $a \in \mathbb{R}^d$ , Theorem 3.2 yields the result of [1, Theorem 9.30] that an MLE  $\vartheta^*$ , attaining the maximum of  $\langle \vartheta, a \rangle - \Lambda(\vartheta)$  subject

to  $\vartheta \in \mathbb{R}^d$ , exists if and only if  $a$  belongs to  $ri(\mu)$ . By the equivalence (i) $\Leftrightarrow$ (iii),  $\theta^*(a) + lin(\mu)^\perp$  is the set of all such MLE's. It also follows from Theorem 3.2 that the conditions

$$lin(dom(\Lambda)) \subseteq M(Q_\theta) \quad \text{and} \quad a - m(Q_\theta) \in N_\theta(dom(\Lambda)) \tag{14}$$

are necessary and sufficient for  $a \in \mathbb{R}^d$  and  $\theta \in dom(\Lambda)$  to satisfy  $\Lambda^*(a) = \langle \theta, a \rangle - \Lambda(\theta)$ . Hence, each  $\theta$  in the range of  $\theta^*$  satisfies the inclusion in (14). Considering  $a = m(Q_\theta)$ , it follows that the range of  $\theta^*$  consists of all  $\theta \in dom(\Lambda) \cap lin(\mu)$  that satisfy the inclusion in (14), in which case  $\theta^*(m(Q_\theta)) = \theta$ .

In terms of subdifferentials, recalled in Sect. 2.5,

$$\partial\Lambda(\theta) = \begin{cases} m(Q_\theta) + N_\theta(dom(\Lambda)), & \text{if } \theta \in dom(\Lambda) \text{ and } lin(dom(\Lambda)) \subseteq M(Q_\theta), \\ \emptyset, & \text{otherwise,} \end{cases} \tag{15}$$

and the inverse of  $\theta^*$  is the restriction of the mapping  $\theta \mapsto \partial\Lambda(\theta)$  to  $dom(\Lambda) \cap lin(\mu)$  in the sense

$$\{a \in ri(\mu) : \theta^*(a) = \theta\} = \partial\Lambda(\theta), \quad \theta \in dom(\Lambda) \cap lin(\mu).$$

When  $dom(\Lambda)$  has nonempty interior, or equivalently  $lin(dom(\Lambda)) = \mathbb{R}^d$ , other well-known results follow. The range of  $\theta^*$  consists of those  $\theta \in dom(\Lambda) \cap lin(\mu)$  for which  $Q_\theta$  has a mean. The subdifferential  $\partial\Lambda(\theta)$  is nonempty if and only if  $\theta \in dom(\Lambda)$  and  $Q_\theta$  has a mean. If  $\theta$  is an interior point of  $dom(\Lambda)$  then  $\partial\Lambda(\theta) = \{m(Q_\theta)\}$ , due to  $N_\theta(dom(\Lambda)) = \{0\}$ , and thus  $\Lambda$  is differentiable at  $\theta$  with its gradient equal to the mean of  $Q_\theta$ , existing by Corollary 2.20. Since the sets  $\partial\Lambda(\theta)$  with  $\theta \in dom(\Lambda) \cap lin(\mu)$  are disjoint and cover  $ri(\mu)$ , each  $a \in ri(\mu)$  is the mean of some pm in  $\mathcal{E}_\mu$  if and only if the family  $\mathcal{E}_\mu$  is steep, in the sense that no  $Q_\vartheta \in \mathcal{E}$  with  $\vartheta$  on the boundary of  $dom(\Lambda)$  has a mean.

Consider next the case  $\mathcal{E} \subseteq lin(\mu)$ . Then,  $\mathcal{E}_{\mu,a} = \mathcal{E}_\mu = cl(\mathcal{E}) \cap dom(\Lambda)$  for each  $a \in \mathbb{R}^d$ , and

$$\mathcal{E}_E^{M,\cap} = \{\theta \in \mathcal{E}_\mu^M : N_\theta(\mathcal{E}) \cap [E \setminus lin(\mu)] \neq \emptyset\}, \quad E \in \mathcal{E}.$$

It follows, as a special case of Lemma 3.3, that the range of  $\theta^*$  is equal to  $\mathcal{E}_\mu^M$ , and for  $\theta \in \mathcal{E}_\mu^M$  the set of all  $a \in ri(\mu) + bar(\mathcal{E})$  with  $\theta^*(a) = \theta$  is equal to  $m(Q_\theta) + N_\theta(\mathcal{E})$ .

In accordance to previous notations, let  $\Psi$  respectively  $\Psi_{\mathcal{E}_\mu}$  denote the function equal to  $\Lambda$  on the set  $\mathcal{E}$  respectively  $\mathcal{E}_\mu$ , and to  $+\infty$  elsewhere. Then, for  $a \in ri(\mu) + bar(\mathcal{E})$ , a maximizer of  $\vartheta \mapsto \langle \vartheta, a \rangle - \Lambda(\vartheta)$  over  $\mathcal{E}$  is a parameter  $\theta \in \mathbb{R}^d$  satisfying  $\Psi^*(a) = \langle \theta, a \rangle - \Psi(\theta)$ , while  $\theta^*(a) = \theta$  is equivalent to  $\Psi^*(a) = \langle \theta, a \rangle - \Psi_{\mathcal{E}_\mu}(\theta)$ . These conditions are, in turn, equivalent to  $a \in \partial\Psi(\theta)$  respectively  $a \in \partial\Psi_{\mathcal{E}_\mu}(\theta)$ , see Sect. 2.5. In particular, the subdifferentials coincide if  $\theta \in \mathcal{E}$  while  $\partial\Psi(\theta)$  is empty

otherwise. By the special case of Lemma 3.3 mentioned above,

$$\partial\Psi_{\mathcal{E}_\mu}(\theta) = \begin{cases} m(Q_\theta) + N_\theta(\mathcal{E}), & \theta \in \mathcal{E}_\mu^M, \\ \emptyset, & \text{otherwise,} \end{cases} \tag{16}$$

which implies a similar representation of  $\partial\Psi(\theta)$ , with  $\mathcal{E}_\mu^M$  replaced by the set  $\mathcal{E}^M$  of those  $\theta \in \mathcal{E}$  that satisfy  $\text{lin}(\mathcal{E}) \subseteq M(Q_\theta)$ . In the nondegenerate case, when  $\mathcal{E}$  intersects the interior of  $\text{dom}(\Lambda)$ , the last inclusion is equivalent to existence of the mean of the pm  $Q_\theta$  by Corollary 2.20.

The above considerations have also the following consequences. The set  $\text{ri}(\mu) + \text{bar}(\mathcal{E})$  partitions into the nonempty subdifferentials  $\partial\Psi_{\mathcal{E}_\mu}(\theta) = m(Q_\theta) + N_\theta(\mathcal{E})$ ,  $\theta \in \mathcal{E}_\mu^M$ , where in the nondegenerate case the last condition means that the pm  $Q_\theta$  parameterized by  $\theta \in \mathcal{E}_\mu$  has a mean. Further, an MLE exists if and only if  $a$  belongs to the subdifferential  $\partial\Psi(\theta)$  at some  $\theta \in \mathcal{E}$  which extends and strengthens the first assertion of [1, Theorem 9.18]. Finally, the condition  $a \in \partial\Psi_{\mathcal{E}_\mu}(\theta)$  is equivalent to  $\theta \in \partial\Psi_{\mathcal{E}_\mu}^*(a)$  by [15, Theorem 23.5], and  $\Psi^*$  is equal to  $\Psi_{\mathcal{E}_\mu}^*$  by Lemma 2.13. As  $\theta = \theta^*(a)$  is uniquely determined by  $a \in \text{ri}(\mu) + \text{bar}(\mathcal{E})$ , this implies that  $\Psi^*$  is differentiable on  $\text{ri}(\mu) + \text{bar}(\mathcal{E})$ , with the gradient  $\nabla\Psi^*(a) = \theta^*(a)$ . This extends the third assertion of [1, Theorem 9.18] where  $\text{lin}(\mu) = \mathbb{R}^d$  and  $\mathcal{E}$  with nonempty interior are assumed.

In the nondegenerate case, and more generally when  $\mathcal{E}$  intersects  $\text{ri}(\text{dom}(\Lambda))$ , the above representations of  $\partial\Psi_{\mathcal{E}_\mu}(\theta)$  and  $\partial\Psi(\theta)$  follow also directly from (15). To see this, recall that in that case for  $\theta \in \mathcal{E}$  Lemma 2.22 gives that  $\partial\Psi(\theta) = \partial\Lambda(\theta) + N_\theta(\mathcal{E})$ , and Lemma 2.19 gives that the inclusion in (15) holds if and only if  $\text{lin}(\mathcal{E}) \subseteq M(Q_\theta)$  thus  $\theta \in \mathcal{E}^M$ .

### 4 The GMLE in general

By Fact 2.7, the restriction of  $\mu$  to the closure of a face  $F$  of  $\text{cc}(\mu)$  is a nonzero measure whose convex core is equal to  $F$ . For this restriction  $\nu = \mu^{\text{cl}(F)}$ , with  $\text{ri}(\nu)$ ,  $\text{aff}(\nu)$ ,  $\text{lin}(\nu)$  equal to  $\text{ri}(F)$ ,  $\text{aff}(F)$ ,  $\text{lin}(F)$ , in the sequel the convenient notations  $\Lambda_F$ ,  $Q_{F,\vartheta}$ ,  $\Psi_{F,\mathcal{E}}^*$ ,  $\theta_{F,\mathcal{E}}^*$ , etc. are used instead of  $\Lambda_\nu$ ,  $Q_{\nu,\vartheta}$ ,  $\Psi_{\nu,\mathcal{E}}^*$ ,  $\theta_{\nu,\mathcal{E}}^*$ , etc.

When  $a \in \text{cc}(\mu) + \text{bar}(\mathcal{E})$ , among the faces  $G$  of  $\text{cc}(\mu)$  satisfying  $a \in \text{ri}(G) + \text{bar}(\mathcal{E})$  there exists the inclusion-largest one, by Lemma 2.1. This face is denoted by  $G^*(a) = G_{\mu,\mathcal{E}}^*(a)$ .

**Theorem 4.1** *If  $a \in \text{dom}(\Psi^*)$  then  $a \in \text{cc}(\mu) + \text{bar}(\mathcal{E})$ , the face  $G = G^*(a)$  of  $\text{cc}(\mu)$  is  $\mathcal{E}$ -accessible,  $\Psi_{G,\mathcal{E}}^*(a)$  is equal to  $\Psi^*(a) = \Psi_{\mu,\mathcal{E}}^*(a)$ , and with  $\theta = \theta_{G,\mathcal{E}}^*(a)$*

$$\Psi^*(a) - [\langle \vartheta, a \rangle - \Lambda(\vartheta)] \geq D(Q_{G,\theta} \| Q_\vartheta), \quad \vartheta \in \mathcal{E}. \tag{17}$$

This proves the existence of GMLE by explicitly identifying it. Indeed, comparison of (2) and (17) gives, referring to Remark 1.2 for uniqueness, that  $Q_{G,\theta} = R^*(a)$ .

**Corollary 4.2** *The GMLE  $R_{\mu,\mathcal{E}}^*(a)$  belongs to  $\mathcal{E} = \mathcal{E}_\mu$  if and only if  $a \in \text{ri}(\mu) + \text{bar}(\mathcal{E})$ .*

In particular, for a full exponential family  $\mathcal{E}$  it follows that  $R^*(a) \in \mathcal{E}$  if and only if  $a \in \text{ri}(\mu)$ , or equivalently, if and only if an MLE exists, see the passage containing Eq. (13) in Sect. 3.3.

*Remark 4.3* The inequality (17) holds even for  $\vartheta \in \tilde{\mathcal{E}}_{\mu,a}$ . In fact, by Lemma 2.11, the bracket and the pm  $Q_\vartheta$  do not change when  $\vartheta \in \mathcal{E}$  is replaced by  $\pi_{\mu,a}(\vartheta)$ , thus (17) holds for  $\vartheta$  in  $\pi_{\mu,a}(\mathcal{E})$ . By limiting along segments, continuity of  $\Lambda$  along segments and lower semicontinuity of  $I$ -divergence, (17) holds for  $\vartheta$  in  $\mathcal{E}_{\mu,a}$  and, in turn, for  $\vartheta \in \tilde{\mathcal{E}}_{\mu,a}$  by Lemma 2.11. A further immediate extension of (17) is

$$\Psi^*(a) - [\langle \vartheta, a \rangle - \Lambda_F(\vartheta)] \geq D(Q_{G,\theta} \| Q_{F,\vartheta}), \quad \vartheta \in \tilde{\mathcal{E}}_{F,a}, \tag{18}$$

where  $F$  is any face of  $\text{cc}(\mu)$  containing  $G = G^*(a)$ . Indeed, the obvious identity

$$D(P \| Q_{F,\vartheta}) - \Lambda_F(\vartheta) = D(P \| Q_\vartheta) - \Lambda(\vartheta),$$

valid for any face  $F$  of  $\text{cc}(\mu)$ , pm  $P \ll \mu^{cl(F)}$ , and  $\vartheta \in \text{dom}(\Lambda)$ , implies with  $P = Q_{G,\theta}$  that for faces  $F \supseteq G$  of  $\text{cc}(\mu)$ , (18) with  $\vartheta \in \mathcal{E}$  is equivalent to (17). Then, the extension to  $\vartheta \in \tilde{\mathcal{E}}_{F,a}$  follows by limiting as above.

Theorem 4.1 will be proved via induction on the dimension of  $\text{aff}(\mu)$ . Five lemmas are sent forth, the first three of elementary geometric content.

**Lemma 4.4** *If  $B \subseteq \text{aff}(\mu)$  then  $a \in B + \text{bar}(\mathcal{E})$  is equivalent to  $a \in B + \text{bar}(\pi_{\mu,a}(\mathcal{E}))$ .*

*Proof* If  $x \in E_{\mu,a}$  then  $\langle \vartheta, x \rangle = \langle \pi_{\mu,a}(\vartheta), x \rangle$  for all  $\vartheta \in \mathbb{R}^d$ , hence  $x \in \text{bar}(\mathcal{E})$  takes place if and only if  $x \in \text{bar}(\pi_{\mu,a}(\mathcal{E}))$ . Applying this to  $x = a - b$  with  $b \in B$ , when  $x \in E_{\mu,a}$  due to  $B \subseteq \text{aff}(\mu)$ , the assertion follows.  $\square$

**Lemma 4.5** *For a nonempty convex set  $C$  and cone  $K$  containing 0, a halfspace  $H_{\leq} = \{x : \langle \theta, x - a \rangle \leq 0\}$  with  $\theta \neq 0$  contains  $C + K$  if and only if  $C \subseteq H_{\leq}$  and  $\theta \in N_0(K)$ .*

*Proof* If  $C + K \subseteq H_{\leq}$  then  $C \subseteq H_{\leq}$  as  $0 \in K$ . In addition, for any  $c$  from the nonempty set  $C$ , the inequality  $\langle \theta, c + y - a \rangle \leq 0$  holds for all  $y \in K$ . Since  $K$  is a cone, it follows that  $\langle \theta, y \rangle \leq 0, y \in K$ , and thus  $\theta \in N_0(K)$ . The converse implication is obvious.  $\square$

**Lemma 4.6** *For  $C, K$  and  $H_{\leq}$  as in Lemma 4.5, if  $C + K \subseteq H_{\leq}$  and  $\text{ri}(C) + K$  intersects the boundary hyperplane  $H$  of  $H_{\leq}$  then  $C \subseteq H$ .*

*Proof* Suppose some  $x \in \text{ri}(C) + K$ , say  $x = c + y$  where  $c \in \text{ri}(C)$  and  $y \in K$ , belongs to  $H$ , thus  $\langle \theta, c + y - a \rangle = 0$ . By Lemma 4.5,  $C \subseteq H_{\leq}$  and  $\theta \in N_0(K)$ , thus  $\langle \theta, c - a \rangle \leq 0$  and  $\langle \theta, y \rangle \leq 0$ . Hence  $\langle \theta, c - a \rangle = 0$ , that is,  $H$  contains  $c \in \text{ri}(C)$ . This and  $C \subseteq H_{\leq}$  imply the assertion.  $\square$

**Lemma 4.7** *If  $a \in \text{dom}(\Psi_{\mu,\Gamma}^*)$  for a convex subset  $\Gamma$  of  $\text{dom}(\Lambda)$ , and  $cs(\mu)$  is contained in a halfspace  $H_{\leq} = \{x : \langle \theta, x - a \rangle \leq 0\}$  with a nonzero  $\theta \in \text{rec}(ri(\Gamma))$ , then the boundary hyperplane  $H$  of  $H_{\leq}$  intersects  $cc(\mu)$  in a face  $F$  of  $cc(\mu)$  such that  $\Psi_{\mu,\Gamma}^*(a) = \Psi_{F,\Gamma}^*(a)$ .*

*Proof* If  $\vartheta \in ri(\Gamma)$  and  $t \geq 0$  then  $\vartheta + t\theta \in ri(\Gamma)$  due to  $\theta \in \text{rec}(ri(\Gamma))$ , thus

$$\Psi_{\mu,\Gamma}^*(a) \geq \langle \vartheta + t\theta, a \rangle - \Lambda(\vartheta + t\theta) = -\ln \int_{\mathbb{R}^d} e^{\langle \vartheta, x - a \rangle + t\langle \theta, x - a \rangle} \mu(dx).$$

Since  $H_{\leq}$  contains  $cs(\mu)$ , it has full  $\mu$ -measure. The integrand is bounded on  $H_{\leq}$  by  $e^{\langle \vartheta, x - a \rangle}$  which is  $\mu$ -integrable due to  $\vartheta \in \Gamma \subseteq \text{dom}(\Lambda)$ . By dominated convergence, the integral converges to  $\int_H e^{\langle \vartheta, x - a \rangle} \mu(dx)$  when  $t \rightarrow \infty$ . It follows that

$$\Psi_{\mu,\Gamma}^*(a) \geq -\ln \int_H e^{\langle \vartheta, x - a \rangle} \mu(dx), \quad \vartheta \in ri(\Gamma),$$

and as  $\Psi_{\mu,\Gamma}^*(a)$  is finite by assumption, this implies  $\mu(H) > 0$ . Hence,  $H$  is a supporting hyperplane of  $cs(\mu)$ . By Fact 2.9,  $F = H \cap cc(\mu)$  is a face of  $cc(\mu)$  and  $\mu(H \setminus cl(F)) = 0$ . Therefore, the above inequality rewrites to

$$\Psi_{\mu,\Gamma}^*(a) \geq -\ln \int_{cl(F)} e^{\langle \vartheta, x - a \rangle} \mu(dx) = \langle \vartheta, a \rangle - \Lambda_F(\vartheta), \quad \vartheta \in ri(\Gamma).$$

This and Lemma 2.13 imply that  $\Psi_{\mu,\Gamma}^*(a) \geq \Psi_{F,ri(\Gamma)}^*(a) = \Psi_{F,\Gamma}^*(a)$ , while  $\Lambda \geq \Lambda_F$  implies that  $\Psi_{\mu,\Gamma}^*(a) \leq \Psi_{F,\Gamma}^*(a)$ . □

**Lemma 4.8** *If  $a$  belongs to  $\text{dom}(\Psi_{\mu,\Xi}^*)$  but not to  $ri(\mu) + \text{bar}(\Xi)$  then some unit vector  $\tau \in \text{rec}(\pi_\mu(ri(\Xi)))$  exposes a proper face  $F$  of  $cc(\mu)$  such that  $\Psi_{\mu,\Xi}^*(a) = \Psi_{F,\Xi}^*(a)$  and  $F$  contains each face  $G$  of  $cc(\mu)$  with  $a \in ri(G) + \text{bar}(\Xi)$ .*

*Proof* Let  $\Gamma$  denote the convex subset  $\pi_{\mu,a}(\Xi)$  of  $\text{dom}(\Lambda)$ . By Corollary 2.12,  $\Psi_{\mu,\Xi}^*(a)$  equals  $\Psi_{\mu,\Gamma}^*(a)$ . This and Lemma 4.4 with  $B = ri(\mu)$  imply that the assumptions equivalently mean that  $a$  is in  $\text{dom}(\Psi_{\mu,\Gamma}^*)$  but not in  $ri(\mu) + \text{bar}(\Gamma)$ . On account of the latter, there exists a halfspace  $H_{\leq} = \{x \in \mathbb{R}^d : \langle \theta, x - a \rangle \leq 0\}$  with nonzero  $\theta$  that contains  $ri(\mu) + \text{bar}(\Gamma)$ . By Lemma 4.5 with  $C = ri(\mu)$  and  $K = \text{bar}(\Gamma)$ , the halfspace  $H_{\leq}$  contains  $ri(\mu)$  and hence also  $cs(\mu)$ , and  $\theta$  belongs to  $N_0(\text{bar}(\Gamma))$  that by Lemma 2.3 is equal to  $\text{rec}(ri(\Gamma))$ . Then, Lemma 4.7 implies that the intersection of  $cc(\mu)$  with the boundary hyperplane  $H$  of  $H_{\leq}$  is a face  $F$  of  $cc(\mu)$ , and  $\Psi_{\mu,\Gamma}^*(a) = \Psi_{F,\Gamma}^*(a)$ .

As  $\Gamma$  is a subset of  $E_{\mu,a}$ , the linear span of  $\{x - a : x \in ri(\mu)\}$ , so is also the recession cone  $\text{rec}(ri(\Gamma))$ , thus the nonzero vector  $\theta$  contained in that cone satisfies  $\langle \theta, x - a \rangle \neq 0$  for some  $x \in ri(\mu)$ . The latter implies  $\langle \theta, x - y \rangle \neq 0$  for all  $y \in H$ , in particular, for  $y$  in  $F = H \cap cc(\mu) \neq \emptyset$ . Thus  $\theta$  is not orthogonal to  $\text{lin}(\mu)$ . Using that  $cc(\mu) \subseteq H_{\leq}$ , it follows that  $F$  is a proper face of  $cc(\mu)$ , exposed by a unit vector

proportional to  $\theta$ , and hence also by the unit vector  $\tau$  proportional to  $\pi_\mu(\theta)$ . This vector satisfies

$$\tau \in \pi_\mu(\text{rec}(ri(\Gamma))) \subseteq \text{rec}(\pi_\mu(ri(\Gamma))) = \text{rec}(\pi_\mu(ri(\mathcal{E}))),$$

where the inclusion holds by Lemma 2.2 and the equality by obvious interchange of projections and relative interiors.

The equality  $\Psi_{\mu, \mathcal{E}}^*(a) = \Psi_{F, \mathcal{E}}^*(a)$  is a consequence of the two equalities in the first passage of the proof and of  $\Psi_{F, \Gamma}^*(a) = \Psi_{F, \mathcal{E}}^*(a)$ , obtained from Corollary 2.12 with  $\mu^{cl(F)}$  in the role of  $\mu$ .

For the last assertion, consider a face  $G$  of  $cc(\mu)$  with  $a \in ri(G) + bar(\mathcal{E})$ , equivalent to  $a \in ri(G) + bar(\Gamma)$  by Lemma 4.4 with  $B = ri(G)$ . Since  $G + bar(\Gamma)$  is contained in  $cc(\mu) + bar(\Gamma) \subseteq H_{\leq}$  and  $ri(G) + bar(\Gamma)$  contains  $a \in H$ , Lemma 4.6 with  $C = G$  and  $K = bar(\Gamma)$  implies  $G \subseteq H$ . This proves that  $G$  is a subset of  $F$ . □

*Proof of Theorem 4.1* If  $a \in ri(\mu) + bar(\mathcal{E})$  then obviously  $a \in cc(\mu) + bar(\mathcal{E})$  and the face  $G_{\mu, \mathcal{E}}^*(a)$  coincides with  $cc(\mu)$  and is  $\mathcal{E}$ -accessible. Since  $a$  and  $\theta_{\mu, \mathcal{E}}^*(a) = \theta$  satisfy (iii) of Theorem 3.2 they also satisfy (i) and (iv) which immediately rewrites to

$$\Psi^*(a) - [\langle \vartheta, a \rangle - \Lambda(\vartheta)] \geq D(Q_\theta \| Q_\vartheta), \quad \vartheta \in \mathcal{E},$$

proving (17). In particular, the assertion of Theorem 4.1 holds when the dimension of  $aff(\mu)$  is zero, since then  $dom(\Psi_{\mu, \mathcal{E}}^*)$  trivially equals  $ri(\mu) + bar(\mathcal{E})$ . Induction argument on the dimension is applied, assuming validity of the assertion for any measure  $\nu$  with  $aff(\nu)$  of smaller dimension than  $aff(\mu)$ .

In the induction step, it suffices to consider the case when  $a$  is in  $dom(\Psi^*)$  but not in  $ri(\mu) + bar(\mathcal{E})$ . Then Lemma 4.8 implies existence of a proper face  $F$  of  $cc(\mu)$  exposed by some unit vector  $\tau \in \text{rec}(\pi_\mu(ri(\mathcal{E})))$  such that  $\Psi^*(a) = \Psi_{\mu, \mathcal{E}}^*(a)$  is equal to  $\Psi_{F, \mathcal{E}}^*(a)$ . In particular,  $a$  belongs to  $dom(\Psi_{F, \mathcal{E}}^*)$ , thus the induction hypothesis applies to  $\nu = \mu^{cl(F)}$ . Hence,  $a \in F + bar(\mathcal{E})$ , the largest face  $G$  of  $F$  satisfying  $a \in ri(G) + bar(\mathcal{E})$  is  $\mathcal{E}$ -accessible,  $\Psi_{F, \mathcal{E}}^*(a) = \Psi_{G, \mathcal{E}}^*(a)$ , and with  $\theta = \theta_{G, \mathcal{E}}^*(a)$

$$\Psi_{F, \mathcal{E}}^*(a) - [\langle \vartheta, a \rangle - \Lambda_F(\vartheta)] \geq D(Q_{G, \theta} \| Q_{F, \vartheta}), \quad \vartheta \in \mathcal{E}.$$

As  $a$  belongs to  $F + bar(\mathcal{E})$  it belongs  $cc(\mu) + bar(\mathcal{E})$ , proving the first assertion. Then, the face  $G^*(a)$  of  $cc(\mu)$  is well defined. By the last assertion of Lemma 4.8, this face is contained in  $F$  whence its maximality implies  $G^*(a) = G$ . The  $\mathcal{E}$ -accessibility of  $G$  from  $cc(\mu)$  is a consequence of that of  $G$  from  $F$  and  $\tau \in \text{rec}(\pi_\mu(ri(\mathcal{E})))$ , proving the second assertion. The equality  $\Psi_{G, \mathcal{E}}^*(a) = \Psi^*(a)$  follows by combining two above equalities of the same type. Finally, since  $\Psi_{F, \mathcal{E}}^*(a) = \Psi^*(a)$ , the last display means that (18) holds for  $\vartheta \in \mathcal{E}$ . As shown in Remark 4.3, this is equivalent to (17). □

**Theorem 4.9**  $dom(\Psi^*) = cc(\mu) + bar(\mathcal{E}) = \bigcup ri(G) + bar(\mathcal{E})$  where the union runs over the  $\mathcal{E}$ -accessible faces  $G$  of  $cc(\mu)$ .

*Proof* The union is obviously contained in  $cc(\mu) + bar(\mathcal{E})$  and contains  $dom(\Psi^*)$  by Theorem 4.1. Thus, it suffices to prove that  $cc(\mu) + bar(\mathcal{E})$  is a subset of  $dom(\Psi^*)$ . This follows from the inclusions  $cc(\mu) \subseteq dom(\Psi^*)$  and  $bar(\mathcal{E}) \subseteq rec(dom(\Psi^*))$ . The former holds since  $cc(\mu) \subseteq dom(\Lambda^*)$  [6, Proposition 1(i)] and the latter is a consequence of  $\mathcal{E} = dom(\Psi)$ , see Remark 1.1, and  $bar(dom(\Psi)) \subseteq rec(dom(\Psi^*))$  [1, Theorem 5.19].  $\square$

*Remark 4.10* In the special case of a full exponential family, Theorem 4.9 implies that  $dom(\Lambda^*)$  is equal to  $cc(\mu) + bar(dom(\Lambda))$ . Once this special case is known, the first equality of Theorem 4.9 is equivalent to  $dom(\Psi^*) = dom(\Lambda^*) + bar(\mathcal{E})$ , due to the obvious inclusion  $bar(dom(\Lambda)) \subseteq bar(\mathcal{E})$ . In particular, the latter equality does not require the condition under which it has been proven in Lemma 2.22, let alone the condition under which it has been known previously, see Remark 2.23.

**Proposition 4.11** *A set  $A \subseteq \mathbb{R}^d$  is equal to  $dom(\Lambda_\mu^*)$  for some nonzero Borel measure  $\mu$  with  $dom(\Lambda_\mu) \neq \emptyset$  if and only if  $A = C + K$  where  $C \subseteq \mathbb{R}^d$  is a nonempty convex set with at most countably many faces and  $K$  is a nonempty closed convex cone contained in  $rec(cl(C))$ .*

For example, closed balls cannot be written in this form if  $d \geq 2$ .

*Proof* For a measure  $\mu$  with the above property the cone  $K = bar(dom(\Lambda_\mu))$  is a subset of  $rec(dom(\Lambda_\mu^*))$  by (13) which is contained in  $rec(cl(dom(\Lambda_\mu^*)))$  by Lemma 2.3. Here,  $cl(dom(\Lambda_\mu^*))$  equals  $cs(\mu)$  by [6, Proposition 1] which is the closure of  $C = cc(\mu)$  by Fact 2.6. Hence,  $K \subseteq rec(cl(C))$ . By [4, Theorem 1],  $C$  has at most a countable number of faces. Theorem 4.9 gives the equality  $dom(\Lambda_\mu^*) = C + K$  and one implication follows.

In the opposite direction, assuming nonempty  $C$  and  $K$  have the above properties, an inspection of the proof of [4, Theorem 1] shows that there exists a Borel pm  $\nu$  such that  $cc(\nu) = C$  and  $dom(\Lambda_\nu) = \mathbb{R}^d$ . If  $K = \{0\}$  then, by Theorem 4.9,  $dom(\Lambda_\mu^*) = C + K$  where  $\mu = \nu$ . Otherwise, a sequence of unit vectors  $b_n \in rec(cl(C))$  exists such that  $K$  is the inclusion-smallest closed convex cone containing them. Let  $a \in ri(C)$  and  $\nu_n$  be the pm sitting on the halfline  $\{a + tb_n : t \geq 0\}$  with a density  $\frac{dt}{(t+1)^2}$ . Then  $dom(\Lambda_{\nu_n})$  is the halfspace given by  $\{\vartheta \in \mathbb{R}^d : \langle \vartheta, b_n \rangle \leq 0\}$ . Let  $\mu$  equal to  $\nu + \sum 2^{-n} \nu_n$ . Then,  $dom(\Lambda_\mu)$  equals the intersection of the halfspaces and, in turn,  $bar(dom(\Lambda_\mu))$  is the cone  $K$  by [15, Sect. 14]. Since the halflines are contained in  $ri(\nu)$ , using that  $rec(cl(C)) = rec(ri(C))$  by Lemma 2.3, the convex core of  $\mu$  is  $C$ . Theorem 4.9 implies that again  $dom(\Lambda_\mu^*) = C + K$ .  $\square$

### 5 Properties of GMLE

5.1 Recall that notations involving the underlying measure in an index are conveniently shortened when that measure is the restriction of  $\mu$  to the closure of a face of  $cc(\mu)$ , replacing the restriction by the face. For example,  $\mathcal{E}_G$  is a shorthand for  $\mathcal{E}_{\mu^{cl(G)}}$ , thus denotes  $cl(\pi_G(\mathcal{E})) \cap dom(\Lambda_G)$ , the set  $\mathcal{E}_G^M$  consists of those  $\theta \in \mathcal{E}_G$  that satisfy the inclusion  $lin(\mathcal{E}) \subseteq M(Q_{G,\theta})$ , and  $K_G^*(\tau)$  is given as in Corollary 3.4 with  $\mu$  replaced by  $\mu^{cl(G)}$ .

By [8, Theorem 2], the variation closure  $cl_v(\mathcal{E}_\Xi)$  of any canonically convex exponential family  $\mathcal{E}_\Xi = \{Q_\vartheta : \vartheta \in \Xi\}$  is equal to the union of the families  $\mathcal{E}_{G, \Xi_G}$  over all  $\Xi$ -accessible faces  $G$  of  $cc(\mu)$ . While there  $\Xi \subseteq \text{lin}(\mu)$  is assumed, the result holds with general  $\Xi \subseteq \text{dom}(\Lambda)$  since the set  $\Xi_G$  depends on  $\Xi$  only through  $\pi_\mu(\Xi)$ .

**Theorem 5.1** *The range of the GMLE mapping  $R^*$  is equal to the set of pm's  $P \in cl_v(\mathcal{E}_\Xi)$  with  $\text{lin}(\Xi) \subseteq M(P)$ , thus the union of the families  $\mathcal{E}_{G, \Xi_G^M}$  over all  $\Xi$ -accessible faces  $G$  of  $cc(\mu)$ . For  $P = Q_{G, \tau}$  in this range, where  $G$  is a  $\Xi$ -accessible face of  $cc(\mu)$  and  $\tau \in \Xi_G^M$ ,*

$$\{a : R^*(a) = P\} = [m(P) + K_G^*(\tau)] \setminus \bigcup [ri(F) + \text{bar}(\Xi)],$$

where the union is over those  $\Xi$ -accessible faces  $F$  of  $cc(\mu)$  that properly contain  $G$ . This set is a shifted cone with the apex  $m(P)$ .

As a consequence,  $\text{dom}(\Psi^*)$  partitions into shifted cones with apices  $m(P)$ , perhaps reducing to singletons  $\{m(P)\}$ , extending the result of Corollary 3.4 about such partitioning of the subset  $ri(\mu) + \text{bar}(\Xi)$  of  $\text{dom}(\Psi^*)$ .

**Corollary 5.2** *In the nondegenerate case, the range of the GMLE mapping consists of those pm's from  $cl_v(\mathcal{E}_\Xi)$  that have means.*

*Proof* If  $\Xi$  intersects the interior of  $\text{dom}(\Lambda)$  then it intersects also the interior of  $\text{dom}(\Lambda_G)$  for each  $\Xi$ -accessible face  $G$  of  $cc(\mu)$ . Corollary 2.20 applied to restrictions of  $\mu$  implies the equivalence of the inclusion  $\text{lin}(\Xi) \subseteq M(P)$  to the existence of the mean of  $P$ , and thus the assertion follows from Theorem 5.1. □

The proof of Theorem 5.1 is preceded by two lemmas.

**Lemma 5.3** *If  $G$  is a proper  $\Xi$ -accessible face of  $cc(\mu)$  and a pm  $P$  is concentrated on  $\text{aff}(G)$  then  $m(P) \notin ri(\mu) + \text{bar}(\Xi)$ .*

*Proof* By the first assumption, some unit vector  $\theta \in \text{rec}(\pi_\mu(ri(\Xi)))$  exposes a proper face  $F$  of  $cc(\mu)$  such that  $G \subseteq F$ . Let  $H_{\leq}$  denote the halfspace  $\{x : \langle \theta, x - a \rangle \leq 0\}$  that contains  $cc(\mu)$  and whose boundary hyperplane  $H$  intersects  $cc(\mu)$  in  $F$ . By Lemma 2.15,  $m(P) \in \text{aff}(P)$ , where  $\text{aff}(P) \subseteq \text{aff}(G)$  due to the second assumption. Hence,  $m(P)$  belongs to  $\text{aff}(F)$ , and thus to  $H$ . On the other hand, let  $C = cc(\mu)$  and  $K = \text{bar}(\pi_\mu(\Xi))$ . By Lemma 2.3,  $\text{rec}(ri(\pi_\mu(\Xi)))$  equals  $N_0(K)$ , whence  $\theta$  belongs to  $N_0(K)$ . This and  $C \subseteq H_{\leq}$  imply  $C + K \subseteq H_{\leq}$  on account of Lemma 4.5. Since  $C \not\subseteq H$ , it follows from Lemma 4.6 that  $ri(C) + K$  does not intersect  $H$ . Therefore,  $m(P) \notin ri(\mu) + \text{bar}(\pi_\mu(\Xi))$ . The proof is completed applying Lemma 4.4 to  $B = ri(\mu)$  and  $a = m(P)$ , as  $\pi_\mu$  and  $\pi_{\mu, a}$  are identical projections due to  $a \in \text{aff}(\mu)$ . □

**Lemma 5.4** *If  $P \in cl_v(\mathcal{E}_\Xi)$  satisfies  $\text{lin}(\Xi) \subseteq M(P)$  then  $m(P)$  belongs to  $ri(P) + \text{bar}(\Xi)$ , this set is contained in  $\text{dom}(\Psi^*)$ , and  $R^*(m(P)) = P$ .*

*Proof* If  $P \in cl_v(\mathcal{E}_\Xi)$  then  $P = Q_{G, \theta}$  for a  $\Xi$ -accessible face  $G$  of  $cc(\mu)$  and some  $\theta \in \Xi_G$ . The latter and the assumption  $\text{lin}(\Xi) \subseteq M(Q_{G, \theta})$  mean that  $\theta$  belongs to  $\Xi_G^M$ .



Lemma 3.3 implies that the set  $\{a \in ri(G) + bar(\mathcal{E}) : \theta_{G, \mathcal{E}}^*(a) = \theta\}$  contains  $m(P)$ . Since  $ri(P) = ri(G)$ , this proves that  $m(P)$  belongs to  $ri(P) + bar(\mathcal{E})$ , a subset of  $dom(\Psi^*)$  due to Theorem 4.9. In addition,  $R^*(m(P)) = P$  follows from Theorem 4.1 if it is shown that  $a = m(P)$  satisfies  $G^*(a) = G$ , or equivalently that if a face  $F$  of  $cc(\mu)$  properly contains  $G$  then  $a \notin ri(F) + bar(\mathcal{E})$ . Lemma 2.4 implies that the  $\mathcal{E}$ -accessible face  $G$  of  $cc(\mu)$  is a  $\mathcal{E}$ -accessible face also of  $F$ , and the assertion follows from Lemma 5.3 applied to  $\nu = \mu^{cl(F)}$  instead of  $\mu$ , since  $P$  is concentrated on  $cl(F) \subseteq aff(F)$  and  $ri(\nu) = ri(F)$ .  $\square$

*Proof of Theorem 5.1* If  $P = R^*(a)$  for some  $a \in dom(\Psi^*)$  then  $P \in cl_\nu(\mathcal{E}_\mathcal{E})$  by Remark 1.2. On account of (2), the  $I$ -divergence  $D(P \| Q_\vartheta)$  is finite for  $\vartheta \in \mathcal{E}$ , thus Corollary 2.17 implies  $lin(\mathcal{E}) \subseteq M(P)$ . This and Lemma 5.4 prove the first assertion, equivalent to the second one due to the representation of  $cl_\nu(\mathcal{E}_\mathcal{E})$  cited above.

For  $P = Q_{G, \tau}$  with  $\mathcal{E}$ -accessible  $G$  and  $\tau \in \mathcal{E}_G^M$ , an element  $a$  of  $dom(\Psi^*)$  satisfies  $R^*(a) = P$  if and only if  $G^*(a) = G$  and  $Q_{\theta_{G, \mathcal{E}}^*(a)}^* = P$ , by Theorem 4.1. The first condition means, by the definition of  $G^*(a)$ , that  $a$  belongs to  $ri(G) + bar(\mathcal{E})$  but not to the union of  $ri(F) + bar(\mathcal{E})$  for the  $\mathcal{E}$ -accessible faces  $F$  of  $cc(\mu)$  containing  $G$  properly. The second condition is equivalent to  $a \in m(P) + K_G^*(\tau)$ , by Corollary 3.4 applied to  $\mu^{cl(G)}$  in the role of  $\mu$ . Since  $m(P)$  belongs to  $ri(P) + bar(\mathcal{E})$  by Lemma 5.4, and  $K_G^*(\tau)$  is contained in  $bar(\mathcal{E})$ ,  $ri(G) + bar(\mathcal{E})$  contains  $m(P) + K_G^*(\tau)$ , and the claimed representation of the set  $\{a : R^*(a) = P\}$  follows.

This set contains  $m(P)$  by Lemma 5.4, hence for the last assertion it suffices to prove that if  $a$  is in the intersection of  $m(P) + K_G^*(\tau)$  and  $ri(F) + bar(\mathcal{E})$  for a  $\mathcal{E}$ -accessible face  $F \supset G$  of  $cc(\mu)$  then  $b_t = m(P) + t[a - m(P)]$  belongs to that intersection for all  $t > 0$ . Obviously,  $b_t \in m(P) + K_G^*(\tau)$  since  $K_G^*(\tau)$  is a cone. Writing  $a = c + d$  with some  $c \in ri(F)$  and  $d \in bar(\mathcal{E})$ , for  $0 < r < \min\{t, 1\}$

$$b_t = [rc + (1 - r)m(P)] + rd + (t - r)(a - m(P)).$$

Since  $m(P) \in ri(P) + bar(\mathcal{E}) \subseteq F + bar(\mathcal{E})$ , the above bracket belongs to  $ri(F) + bar(\mathcal{E})$ . This,  $d \in bar(\mathcal{E})$  and  $a - m(P) \in K_G^*(\tau) \subseteq bar(\mathcal{E})$  imply  $b_t \in ri(F) + bar(\mathcal{E})$ .  $\square$

*Example 5.5* Let us consider the situation of Example 1.4 with its notations. There, it was shown by elementary means that for  $P = \nu$  the set  $\{a : R^*(a) = P\}$  equals  $\partial C$ , which is a nonconvex cone. This equality can be simply derived also from Theorem 4.1. In fact, the convex core of  $\mu$  consists of the origin  $c$  and the interior of  $C$  [4, Example 1]. Its only proper face is  $G = \{c\}$ . Any  $a \in \partial C$  belongs to  $dom(\Lambda^*) = C$  and since it is not in  $ri(\mu) + bar(\mathcal{E}) = ri(\mu)$  it follows that  $G^*(a) = G$ . Theorem 4.1 implies  $R^*(a) = P$ . But  $\{a : R^*(a) = P\}$  is contained in  $\partial C$  by Corollary 4.2.

For an illustration, the equality  $\{a : R^*(a) = P\} = \partial C$  is now derived directly from Theorem 5.1. The point mass  $P = \nu$  belongs to the range of  $R^*$  because  $P = Q_{G, \tau}$  for  $\tau = (0, 0, 0)$  where the face  $G$  of  $cc(\mu)$  is  $\mathcal{E}$ -accessible and  $\tau \in \mathcal{E}_G^M$ , which is a singleton. By Theorem 5.1,

$$\{a : R^*(a) = P\} = [c + K_G^*(\tau)] \setminus [ri(C) + bar(\mathcal{E})] = K_G^*(\tau) \setminus ri(C)$$

and it suffices to show that  $K_G^*(\tau) = C$ . To this end, observe that the mapping  $\theta_{G,\mathcal{E}}^*$  is defined on  $ri(G) + bar(\mathcal{E}) = C$ . Obviously,  $\theta_{G,\mathcal{E}}^*(c) = \tau$ . For  $b \in C \setminus G$ ,  $\mathcal{E}_{G,b} = cl(\pi_{G,b}(\mathcal{E}))$  is the orthogonal projection of  $(1, 0, 0) - C$  to  $E_{G,b}$ , the line that passes through  $b$  and  $c$ , and it is easy to see directly from the definition of  $\theta_{G,\mathcal{E}}^*$  that  $\theta_{G,\mathcal{E}}^*(b)$  is the projection of  $(1, 0, 0)$  to this line. Therefore, the range of  $\theta_{G,\mathcal{E}}^*$  is the union of  $\{\tau\}$  with a closed halfsphere centered at  $(\frac{1}{2}, 0, 0)$ , see Fig. 1. This follows also from the first assertion of Lemma 3.3, noting that  $\mathcal{E}_E^{M,\cap}$  is nonempty only for those lines  $E$  that intersect  $C$  in a halfline and equals the intersection of  $E$  with the halfsphere. By Corollary 3.4,  $K_G^*(\tau)$  is the union of  $K_G(\theta)$  over  $\theta$  in the range  $\theta_{G,\mathcal{E}}^*(C)$ . By Lemma 3.3,  $K_G(\theta) = G$  for  $\theta = \tau$ , and for  $\theta$  in the halfsphere,  $K_G(\theta) = N_\theta(\mathcal{E}) \cap [E \setminus G] = E \cap [C \setminus G]$  where  $E$  is the line through  $\tau$  and  $\theta$ . Hence,  $K_G^*(\tau) = C$  indeed.

5.2 The following result establishes continuity of the GMLE mapping when its domain  $dom(\Psi^*)$  is endowed with the topology that corresponds to the Euclidean topology of the graph of the function  $\Psi^*$ , and its range is endowed with the topology of variation distance.

**Theorem 5.6** *If a convergent sequence  $a_n$  and its limit  $a$  are in  $dom(\Psi^*)$ , the convergence of  $\Psi^*(a_n)$  to  $\Psi^*(a)$  implies that  $R^*(a_n)$  converges to  $R^*(a)$  in variation distance, and also in reversed I-divergence.*

*Proof* Let  $\epsilon > 0$ . Since  $\Psi^*(a)$  is finite there exist  $\vartheta \in \mathcal{E}$  such that

$$\frac{\epsilon^2}{18} > \Psi^*(a) - [\langle \vartheta, a \rangle - \Lambda(\vartheta)].$$

This and the hypotheses on the convergence imply

$$\frac{2\epsilon^2}{9} > \frac{\epsilon^2}{6} + \Psi^*(a) - [\langle \vartheta, a \rangle - \Lambda(\vartheta)] > \Psi^*(a_n) - [\langle \vartheta, a_n \rangle - \Lambda(\vartheta)],$$

the right inequality holding for sufficiently large  $n$ . From these inequalities and (2) it follows by the Pinsker inequality  $2D(P\|Q) \geq |P - Q|^2$ , where  $|P - Q|$  denotes the variation distance of pm's  $P, Q$ , that  $\frac{\epsilon}{3} > |R^*(a) - Q_\vartheta|$  and  $\frac{2\epsilon}{3} > |R^*(a_n) - Q_\vartheta|$ . Thus, the variation distance between  $R^*(a)$  and  $R^*(a_n)$  is eventually below  $\epsilon$ , proving that  $R^*(a_n)$  converges to  $R^*(a)$  in variation distance. The  $rl$ -convergence,  $D(R^*(a)\|R^*(a_n)) \rightarrow 0$ , is a consequence of the following lemma.  $\square$

**Lemma 5.7** *If a variation convergent sequence of pm's  $P_n$  and its limit  $P$  belong to the range of  $R^*$  then  $D(P\|P_n) \rightarrow 0$ .*

*Proof* By Theorem 5.1,  $P_n = Q_{F_n, \vartheta_n}$  and  $P = Q_{F, \vartheta}$  for faces  $F_n$  and  $F$  of  $cc(\mu)$ , with  $\vartheta_n \in cl(\pi_{F_n}(\mathcal{E}))$  and  $\vartheta \in cl(\pi_F(\mathcal{E}))$ . The convergence in variation and [8, Theorem 4] imply that  $F_n \supseteq F$  eventually,  $Q_{F, \vartheta_n}$  converges to  $P$  in variation, and  $P_n(cl(F)) \rightarrow 1$ . If  $F_n \supseteq F$  holds then  $\pi_F(\vartheta_n)$  belongs to  $\pi_F(cl(\pi_{F_n}(\mathcal{E}))) \subseteq cl(\pi_F(\mathcal{E}))$ . Therefore,  $\pi_F(\vartheta_n) - \vartheta$  is in  $lin(\pi_F(\mathcal{E}))$  which is a subset of

$lin(\mathcal{E}) + lin(F)^\perp$ . Here,  $lin(\mathcal{E}) \subseteq M(P)$  by Theorem 5.1, and  $lin(F)^\perp \subseteq M(P)$  by Lemma 2.14. Hence,  $\pi_F(\vartheta_n) - \vartheta \in M(P)$ . It follows that  $D(P\|Q_{F,\vartheta_n})$  is eventually finite. By [8, Theorem 3(i)], applied to the sequence  $Q_{F,\vartheta_n}$  converging in variation,  $D(P\|Q_{F,\vartheta_n}) \rightarrow 0$ . This  $I$ -divergence equals  $D(P\|P_n) + \ln P_n(cl(F))$ , and the assertion follows since  $P_n(cl(F)) \rightarrow 1$ .  $\square$

*Example 5.8* Let  $\mu$  be sum of the point mass  $\nu$  at  $a = (0, 0)$  and the Lebesgue measure on the segment with endpoints  $(0, 1)$  and  $(1, 1)$ . Then

$$\Lambda(\vartheta) = \ln \left[ 1 + e^{\vartheta_2} \cdot \frac{e^{\vartheta_1} - 1}{\vartheta_1} \right], \quad \vartheta = (\vartheta_1, \vartheta_2) \in \mathbb{R}^2,$$

and  $Q_\vartheta$  has the mean

$$e^{\vartheta_2} \cdot \left[ 1 + e^{\vartheta_2} \cdot \frac{e^{\vartheta_1} - 1}{\vartheta_1} \right]^{-1} \left( \frac{\vartheta_1 e^{\vartheta_1} - e^{\vartheta_1} + 1}{\vartheta_1^2}, \frac{e^{\vartheta_1} - 1}{\vartheta_1} \right).$$

The pm's  $Q_{\theta_n}, \theta_n = (-n e^n, n)$ , in the full exponential family  $\mathcal{E}_\mu$ , are the GMLE's  $R^*(a_n)$  for  $a_n = m(Q_{\theta_n})$ , and  $R^*(a) = \nu$ . Then  $a_n \rightarrow a$  and  $D(R^*(a)\|R^*(a_n)) = \Lambda(\theta_n) \rightarrow 0$ , while  $\Lambda^*(\theta_n)$  converges to 1, different from  $\Lambda^*(a) = 0$ . This shows that the converse of the implication in Theorem 5.6 fails even for a full exponential family with bounded support.

*Remark 5.9* If  $\mathcal{E} = \mathcal{E}_\mu$  is a full exponential family such that  $dom(\Lambda) = \mathbb{R}^d$  then  $cl_\nu(\mathcal{E})$  is equal to the union of the full families  $\mathcal{E}_G$  over all faces  $G$  of  $cc(\mu)$ , called the *extension*  $ext(\mathcal{E})$  of  $\mathcal{E}$  [6]. As all pm's in this extension have a mean, by Theorem 5.1 the range of the GMLE mapping is equal to  $ext(\mathcal{E})$  in this case. The barrier cone of  $\mathcal{E} = dom(\Lambda) = \mathbb{R}^d$  is the singleton  $\{0\}$ , thus  $G^*(a)$  in Theorem 4.1 equals the unique face  $G$  of  $cc(\mu)$  with  $a \in ri(G)$ . By Theorems 4.9 and 5.1,  $dom(\Lambda^*) = cc(\mu)$  and the mapping  $R^*: cc(\mu) \rightarrow ext(\mathcal{E})$  is bijective. It assigns to each  $a \in cc(\mu)$  the unique pm in  $ext(\mathcal{E})$  whose mean is equal to  $a$ . Note that in the considered case the GMLE can always be identified with a true MLE when  $ext(\mathcal{E})$  rather than  $\mathcal{E}$  is taken as the model family, see [6, Corollary 12], or the next section. Under the additional assumptions that the set  $cc(\mu)$  is bounded and locally simplicial, the stronger result holds that the GMLE mapping is a homeomorphism between  $cc(\mu)$  and  $ext(\mathcal{E})$ , when the latter is endowed with the topology of variation distance. Indeed, boundedness of  $cc(\mu)$  implies continuity of the inverse mapping that assigns to  $P \in ext(\mathcal{E})$  its mean, the locally simplicial property implies continuity of the function  $\Lambda^*$  on  $dom(\Lambda^*) = cc(\mu)$ , by [15, Theorem 10.2], and then the mapping  $a \mapsto R^*(a)$  is continuous by Theorem 5.6. In the special instance when  $\mu$  is concentrated on a finite set, the convex hull of this set and  $ext(\mathcal{E})$  have been known to be homeomorphic under an MLE with the model family  $ext(\mathcal{E})$  [1, Theorem 9.15] and the uniqueness of MLE in  $ext(\mathcal{E})$  in terms of means appeared in [1, Theorem 9.16].

5.3 The concept of GMLE is related to that of generalized  $rI$ -projection [6]. A set of pm's is *log-convex* [6] if it contains the log-convex combinations of pairs of

nonsingular pm's in the set. Here, the log-convex combinations  $\overline{P^t Q^{1-t}}$ ,  $0 < t < 1$ , of pm's  $P$  and  $Q$  with  $\mu$ -densities  $p$  and  $q$  have the  $\mu$ -densities  $p^t q^{1-t} / \int p^t q^{1-t} d\mu$ .

According to [6, Theorem 1], if a set  $\mathcal{S}$  of pm's is log-convex and  $P$  is any pm such that  $D(P\|\mathcal{S}) \triangleq \inf_{Q \in \mathcal{S}} D(P\|Q)$  is finite, there exists a unique pm  $Q^*$  that satisfies

$$D(P\|Q) \geq D(P\|\mathcal{S}) + D(Q^*\|Q), \quad Q \in \mathcal{S}. \tag{19}$$

This pm  $Q^*$  has been called the generalized  $rI$ -projection of  $P$  to  $\mathcal{S}$ , because (19) implies variation convergence to  $Q^*$  of any sequence of pm's  $Q_n \in \mathcal{S}$  with  $D(P\|Q_n) \rightarrow D(P\|\mathcal{S})$ .

The canonically convex exponential families are obviously log-convex, thus the mentioned result applies to  $\mathcal{S} = \mathcal{E}_{\mathcal{E}}$ . In that case, as observed in [6], the generalized  $rI$ -projection  $Q^*$  of  $P$  is equal to the GMLE  $R^*(a) = R_{\mu, \mathcal{E}}^*(a)$  if the pm  $P$  has a mean and the mean equals  $a$ . More generally, for any pm  $P$  on  $\mathbb{R}^d$  with  $D(P\|\mathcal{E}_{\mathcal{E}})$  finite, Lemma 2.16 implies that the nonempty set  $\Gamma$  of all  $\vartheta \in \mathcal{E}$  with finite  $D(P\|Q_{\vartheta})$  is convex, and for  $\vartheta \in \Gamma$  the equality

$$D(P\|Q_{\vartheta}) - D(P\|\mathcal{E}_{\Gamma}) = \Psi_{\Gamma}^*(m(P)) - [\langle \vartheta, m(P) \rangle - \Lambda(\vartheta)]$$

follows from (4) by taking supremum over  $\theta \in \Gamma$ . In particular,  $\Psi_{\Gamma}^*(m(P))$  is finite, and (19) for  $\mathcal{S} = \mathcal{E}_{\mathcal{E}}$  is equivalent to

$$\Psi_{\Gamma}^*(m(P)) - [\langle \vartheta, m(P) \rangle - \Lambda(\vartheta)] \geq D(Q^*\|Q_{\vartheta}), \quad \vartheta \in \mathcal{E}_P.$$

Comparing this with (2),  $Q^*$  equals the GMLE  $R_{\mu, \Gamma}^*(m(P))$  by Remark 1.2. Thus, a generalized  $rI$ -projection to a canonically convex exponential family is always equal to a GMLE and can be explicitly described via Theorem 4.1.

The extension  $ext(\mathcal{E})$  of an exponential family  $\mathcal{E}$  is always log-convex [6, Theorem 2]. On the other hand, the  $rI$ -closure of  $\mathcal{E}$ , the set of all pm's to which some sequence in  $\mathcal{E}$  converges in reversed  $I$ -divergence, need not be log-convex [7].

**Theorem 5.10**  *$cl_v(\mathcal{E}_{\mathcal{E}})$  and the range of the GMLE mapping are log-convex.*

*Proof* Consider two pm's  $Q_{F, \vartheta}$  and  $Q_{G, \theta}$  in  $cl_v(\mathcal{E})$  where  $F, G$  are  $\mathcal{E}$ -accessible faces of  $cc(\mu)$ ,  $\vartheta \in \mathcal{E}_F$  and  $\theta \in \mathcal{E}_G$ . If these pm's are not mutually singular then the  $\mu$ -measure of  $cl(F) \cap cl(G)$  is positive, and then  $F \cap G \neq \emptyset$  by Fact 2.10. By Corollary 2.5,  $F \cap G$  is a  $\mathcal{E}$ -accessible face of  $cc(\mu)$ .

Since the  $\mu$ -densities of the two pm's are proportional to  $e^{\langle \vartheta, x \rangle}$  on  $cl(F)$ , respectively to  $e^{\langle \theta, x \rangle}$  on  $cl(G)$ , and 0 elsewhere, the  $\mu$ -density of their log-convex combination is proportional to  $e^{\langle t\vartheta + (1-t)\theta, x \rangle}$  on  $cl(F) \cap cl(G)$  and 0 elsewhere. By Fact 2.10, the log-convex combination is equal to  $Q_{F \cap G, t\vartheta + (1-t)\theta}$ . This pm does not change if  $t\vartheta + (1-t)\theta$  is replaced by its projection to  $lin(F \cap G)$ . The projector  $\pi_{F \cap G}$  maps  $\vartheta$  into

$$\pi_{F \cap G}(cl(\pi_F(\mathcal{E})) \cap dom(\Lambda_F)) \subseteq cl(\pi_{F \cap G}(\mathcal{E})) \cap dom(\Lambda_{F \cap G}) = \mathcal{E}_{F \cap G}.$$

Analogously, also  $\pi_{F \cap G}(\theta)$  is in  $\mathcal{E}_{F \cap G}$ , and thus so is  $\pi_{F \cap G}(t\vartheta + (1-t)\theta)$  due to convexity of this set. It follows that the log-convex combinations of  $Q_{F,\vartheta}$  and  $Q_{G,\theta}$  belong to the subfamily  $\mathcal{E}_{F \cap G, \mathcal{E}_{F \cap G}}$  of  $cl_v(\mathcal{E}_{\mathcal{E}})$ .

On account of Theorem 5.1, for the log-convexity of the range it suffices to prove that if  $lin(\mathcal{E})$  is contained in  $M(P)$  and  $M(Q)$  for two nonsingular pm's  $P$  and  $Q$  then it is contained also in  $M(P^t Q^{1-t})$ ,  $0 < t < 1$ . This holds because for  $\mu$ -densities  $p, q$  and  $f \geq 0$

$$t \int f p d\mu + (1-t) \int f q d\mu \geq \int f p^t q^{1-t} d\mu$$

which implies that if  $f$  is  $P$ - and  $Q$ -integrable then it is  $\overline{P^t Q^{1-t}}$ -integrable as well. □

### 6 MLE in the variation closure

For  $a \in \mathbb{R}^d$  define  $\Phi(a) = \Phi_{\mu, \mathcal{E}}(a)$  by

$$\Phi(a) = \sup \{ \ell_a(F, \vartheta) : F \text{ is a } \mathcal{E}\text{-accessible face of } cc(\mu) \text{ and } \vartheta \in \mathcal{E}_F \} \quad (20)$$

where  $\ell_a(F, \vartheta)$  equals  $\langle \vartheta, a \rangle - \Lambda_F(\vartheta)$  if  $a \in aff(F)$  and  $-\infty$  elsewhere. Note that the pairs  $(F, \vartheta)$  considered here parameterize bijectively the pm's in the variation closure of  $\mathcal{E}_{\mathcal{E}}$ . Equivalently,

$$\Phi(a) = \sup \{ \Psi_{F, \mathcal{E}}^*(a) : F \text{ is a } \mathcal{E}\text{-accessible face of } cc(\mu) \text{ with } a \in aff(F) \}, \quad (21)$$

because if  $a \in aff(F)$ , thus  $\mathcal{E}_{F,a} = \mathcal{E}_F$ , then  $\Psi_{F, \mathcal{E}}^*(a) = \Psi_{F, \mathcal{E}_F}^*(a)$  by Lemma 2.13 and Corollary 2.12.

If  $\Phi(a)$  is finite and the supremum in (20) is actually a maximum, a maximizing pair  $(F^*, \vartheta^*)$  has a statistical interpretation whenever  $a$  equals the mean of an i.i.d. sample from an unknown pm  $Q_{F,\vartheta}$  in  $cl_v(\mathcal{E}_{\mathcal{E}})$ :  $(F^*, \vartheta^*)$  is an MLE from the sample of the unknown parameter  $(F, \vartheta)$ . The implicit understanding behind this interpretation is that the  $\mu$ -density of  $Q_{F,\vartheta}$  is equal to  $e^{(\vartheta, x) - \Lambda_F(\vartheta)}$  on  $aff(F)$  and 0 elsewhere. This understanding, though not the same as in previous sections, is legitimate as densities can differ on a set of  $\mu$ -measure 0, and Fact 2.9 implies by a simple induction that  $aff(F) \setminus cl(F)$  has  $\mu$ -measure 0. An MLE in the present sense can exist only if  $a \in aff(\mu)$ , for otherwise  $\Phi(a) = -\infty$ .

**Theorem 6.1** *For  $a \in aff(\mu)$  the equality  $\Phi(a) = \Psi^*(a)$  holds, and the supremum in (20) is a maximum if and only if  $a$  belongs to the intersection of  $dom(\Psi^*)$  and  $aff(G)$  for  $G = G^*(a)$ . In this case,  $(G, \theta_{G, \mathcal{E}}^*(a))$  is the unique maximizer in (20).*

Interpreting  $a$  as the mean of an i.i.d. sample from an unknown pm  $Q_{F,\vartheta}$  in  $cl_v(\mathcal{E}_{\mathcal{E}})$  as above, Theorem 6.1 gives a necessary and sufficient condition for the existence of an MLE of the unknown parameter  $(F, \vartheta)$  from this sample. If the condition holds then the MLE is unique, equals  $(G, \theta_{G, \mathcal{E}}^*(a))$  with  $G = G^*(a)$ , and thus parameterizes the

GMLE  $R^*(a)$ , see Theorem 4.1. In case of a full exponential family  $\mathcal{E}$  with  $dom(\Lambda) = \mathbb{R}^d$ , each  $a$  in  $dom(\Lambda^*) = cc(\mu)$  meets that necessary and sufficient condition and therefore the GMLE for the family  $\mathcal{E}$  can be identified with the MLE in  $cl_\nu(\mathcal{E}) = ext(\mathcal{E})$ , see Remark 5.9.

The proof of Theorem 6.1 is preceded by three lemmas that complement Theorem 4.1 and give more insight into  $\Phi(a)$  expressed through (21).

**Lemma 6.2** *If  $a \in ri(\mu) + bar(\mathcal{E})$  and  $F$  is a proper face of  $cc(\mu)$  then  $\Psi^*(a) < \Psi_{F,\mathcal{E}}^*(a)$ .*

*Proof* The first assumption implies by Theorem 3.1 that  $\Psi^*(a) = \langle \theta, a \rangle - \Lambda(\theta)$  for  $\theta = \theta^*(a)$  in  $\mathcal{E}_{\mu,a} = cl(\pi_{\mu,a}(\mathcal{E})) \cap dom(\Lambda)$ . This and the second assumption give that  $\Psi^*(a)$  is less than  $\langle \theta, a \rangle - \Lambda_F(\theta)$ , and thus less than  $\Psi_{F,\mathcal{E}_{\mu,a}}^*(a)$ . The proof is completed by observing that  $\Psi_{F,\mathcal{E}_{\mu,a}}^*(a)$ ,  $\Psi_{F,\pi_{\mu,a}(\mathcal{E})}^*(a)$  and  $\Psi_{F,\mathcal{E}}^*(a)$  are equal, by Lemma 2.13 and Corollary 2.12, applied to  $\mu^{cl(F)}$  in the role of  $\mu$ . □

**Lemma 6.3** *For  $a \in dom(\Psi^*)$  and a face  $F$  of  $cc(\mu)$ ,  $\Psi^*(a) = \Psi_{F,\mathcal{E}}^*(a)$  if and only if  $F$  contains  $G^*(a)$ .*

*Proof* If  $F$  contains  $G = G^*(a)$ , then  $\Psi_{F,\mathcal{E}}^*(a)$  is between  $\Psi^*(a)$  and  $\Psi_{G,\mathcal{E}}^*(a)$  which coincide by Theorem 4.1. If  $F \not\supseteq G$  then only the case  $a \in F + bar(\mathcal{E})$  is of interest, for otherwise  $\Psi_{F,\mathcal{E}}^*(a)$  is infinite by Theorem 4.9 while  $\Psi^*(a)$  is finite by the first assumption. Then, by Lemma 2.1, the largest face  $G'$  of  $F$  with  $a \in ri(G') + bar(\mathcal{E})$  is well defined, and it is contained in  $G = G^*(a)$  by the definition of the latter. Thus,  $G'$  is a subset of  $G \cap F$ , and therefore a proper face of  $G$ . By Lemma 6.2 applied to  $\mu^{cl(G)}$  rather than  $\mu$ ,  $\Psi_{G',\mathcal{E}}^*(a) < \Psi_{G,\mathcal{E}}^*(a)$ . By Theorem 4.1,  $\Psi_{G,\mathcal{E}}^*(a) = \Psi^*(a)$  and  $\Psi_{G',\mathcal{E}}^*(a) = \Psi_{F,\mathcal{E}}^*(a)$ , thus  $\Psi^*(a) < \Psi_{F,\mathcal{E}}^*(a)$ . □

**Lemma 6.4** *If  $F$  is a  $\mathcal{E}$ -accessible face of  $cc(\mu)$  and  $a \in aff(F)$  then  $\Psi_{F,\mathcal{E}}^*(a) = \Psi^*(a)$ .*

*Proof* As always  $\Psi_{F,\mathcal{E}}^*(a) \geq \Psi^*(a)$ , it may be assumed that  $a \in dom(\Psi^*)$ . Then, induction argument is applied on the dimension of  $aff(\mu)$ . The assertion is trivial in the case  $F = cc(\mu)$ , in particular, it holds if the dimension is zero.

If the  $\mathcal{E}$ -accessible face  $F$  of  $cc(\mu)$  is proper, there exists a nonempty access sequence to  $F$  adapted to  $\mathcal{E}$ . Then, its first element  $\theta$  belongs to  $rec(\pi_\mu(ri(\mathcal{E})))$  and exposes a proper face  $G$  of  $cc(\mu)$  that contains  $F$  as its  $\mathcal{E}$ -accessible face. It suffices to show that  $\Psi^*(a) = \Psi_{\mu,\mathcal{E}}^*(a)$  is equal to  $\Psi_{G,\mathcal{E}}^*(a)$ , for then the induction hypothesis applies to  $\mu^{cl(G)}$ , yielding  $\Psi_{G,\mathcal{E}}^*(a) = \Psi_{F,\mathcal{E}}^*(a)$ , and the assertion follows.

To this end, observe that Lemma 4.7 can be applied to  $a$ ,  $\Gamma = \pi_\mu(\mathcal{E})$  and  $\theta$ . In fact,  $a \in aff(F)$  implies  $a \in aff(\mu)$ , thus  $\pi_{\mu,a}(\mathcal{E}) = \Gamma$ , hence  $\Psi_{\mu,\mathcal{E}}^*(a) = \Psi_{\mu,\Gamma}^*(a)$  by Corollary 2.12. Thus,  $a \in dom(\Psi^*)$  implies  $a \in dom(\Psi_{\mu,\Gamma}^*)$ . Since  $\theta$  exposes  $G$ , and  $a$  belongs to  $aff(F) \subseteq aff(G)$ , the halfspace  $H_{\leq} = \{x \in \mathbb{R}^d : \langle \theta, x - a \rangle \leq 0\}$  contains  $cs(\mu)$ . The hypothesis  $\theta \in rec(ri(\Gamma))$  is satisfied since  $\pi_\mu$  commutes with relative interiors. Hence, Lemma 4.7 gives the equality  $\Psi_{\mu,\Gamma}^*(a) = \Psi_{G,\Gamma}^*(a)$ . Here, the left-hand side equals  $\Psi_{\mu,\mathcal{E}}^*(a)$  as observed above and, similarly, the right-hand side equals  $\Psi_{G,\mathcal{E}}^*(a)$ , on account of  $a \in aff(G)$  and the second assertion of Corollary 2.12 applied to  $\mu^{cl(F)}$  in the role of  $\mu$ . □

*Proof of Theorem 6.1* The first assertion immediately follows from Lemma 6.4 and (21). If  $a \in \text{dom}(\Psi^*)$  then  $\Psi^*(a) = \Psi_{G,\mathcal{E}}^*(a)$  by Theorem 4.1, where  $G = G^*(a)$ . If additionally  $a \in \text{aff}(G)$  then  $\mathcal{E}_{G,a} = \mathcal{E}_G$  thus the maximizer  $\theta_{G,\mathcal{E}}^*(a)$  of  $\langle \theta, a \rangle - \Lambda_G(\theta)$  subject to  $\theta \in \mathcal{E}_{G,a}$  belongs to  $\mathcal{E}_G$ . Then,  $(G, \theta_{G,\mathcal{E}}^*(a))$  is a maximizer in (20).

If  $(F^*, \vartheta^*)$  attains the maximum in (20) then  $a \in \text{aff}(F^*)$ ,  $\vartheta^* \in \mathcal{E}_{F^*}$  and, using the first assertion,  $\langle \vartheta^*, a \rangle - \Lambda_{F^*}(\vartheta^*) = \Psi^*(a)$ . Thus,  $a$  belongs to  $\text{dom}(\Psi^*)$ . As  $\Psi^*(a) = \Psi_{F^*,\mathcal{E}}^*(a)$  by Lemma 6.4, and  $a \in \text{aff}(F^*)$  implies  $\mathcal{E}_{F^*} = \mathcal{E}_{F^*,a}$ , it follows by Theorem 3.2 that  $\vartheta^* = \theta_{F^*,\mathcal{E}}^*(a)$  and  $a \in \text{ri}(F^*) + \text{bar}(\mathcal{E})$ . The latter implies by the definition of  $G = G^*(a)$  that  $F^* \subseteq G$ . This and Lemma 6.3 prove that  $F^* = G$ . Finally,  $(F^*, \vartheta^*)$  is equal to  $(G, \theta_{G,\mathcal{E}}^*(a))$ , proving also  $a \in \text{aff}(G)$  and uniqueness of the maximizer. □

### 7 Appendix

Some concepts of this paper depend on the chosen parametrization of the underlying exponential families, but many are actually invariant. This is discussed in detail here. Let  $\mu$  and  $\nu$  be nonzero,  $\sigma$ -finite Borel measures on  $\mathbb{R}^d$  whose log-Laplace transforms have nonempty effective domains, and let  $\mathcal{E} \subseteq \text{dom}(\Lambda_\mu)$  and  $\Gamma \subseteq \text{dom}(\Lambda_\nu)$  be nonempty and convex sets.

**Lemma 7.1** *If  $Q_{\mu,\xi}$  coincides with  $Q_{\nu,\eta}$  for some  $\xi \in \text{dom}(\Lambda_\mu)$  and  $\eta \in \text{dom}(\Lambda_\nu)$  then for  $\tau = \eta - \xi$  and  $t = \Lambda_\nu(\eta) - \Lambda_\mu(\xi)$*

- (i)  $\Lambda_\mu(\vartheta) + t = \Lambda_\nu(\vartheta + \tau)$ ,  $\vartheta \in \mathbb{R}^d$ ,
- (ii)  $\text{dom}(\Lambda_\nu) = \text{dom}(\Lambda_\mu) + \tau$ ,
- (iii)  $Q_{\mu,\vartheta} = Q_{\nu,\vartheta+\tau}$ ,  $\vartheta \in \text{dom}(\Lambda_\mu)$ ,
- (iv)  $\Lambda_\mu^*(a) = \Lambda_\nu^*(a) - [\langle \tau, a \rangle - t]$ ,  $a \in \mathbb{R}^d$ .

*Proof* By assumption,

$$\frac{d\mu}{d\nu}(x) = \frac{dQ_{\nu,\eta}(x)}{d\nu} \bigg/ \frac{dQ_{\mu,\xi}(x)}{d\mu} = e^{\langle \eta, x \rangle - \Lambda_\nu(\eta) - [\langle \xi, x \rangle - \Lambda_\mu(\xi)]} = e^{\langle \tau, x \rangle - t},$$

and then, for  $\vartheta \in \mathbb{R}^d$

$$\Lambda_\mu(\vartheta) = \ln \int_{\mathbb{R}^d} e^{\langle \vartheta, x \rangle} \mu(dx) = \ln \int_{\mathbb{R}^d} e^{\langle \vartheta, x \rangle} e^{\langle \tau, x \rangle - t} \nu(dx) = \Lambda_\nu(\vartheta + \tau) - t,$$

thus (i) and (ii) hold. For  $\vartheta \in \text{dom}(\Lambda_\mu)$

$$\ln \frac{dQ_{\mu,\vartheta}}{d\nu}(x) = \ln \frac{dQ_{\mu,\vartheta}}{d\mu}(x) + \ln \frac{d\mu}{d\nu}(x) = \langle \vartheta, x \rangle - \Lambda_\mu(\vartheta) + \langle \tau, x \rangle - t$$

which by (i) equals  $\langle \vartheta + \tau, x \rangle - \Lambda_\nu(\vartheta + \tau)$ . Hence, (iii) follows. To see (iv), write

$$\begin{aligned} \Lambda_\mu^*(a) &= \sup_{\vartheta \in \mathbb{R}^d} [\langle \vartheta, a \rangle - (\Lambda_\nu(\vartheta + \tau) - t)] \\ &= \sup_{\vartheta \in \mathbb{R}^d} [\langle \vartheta + \tau, a \rangle - \Lambda_\nu(\vartheta + \tau)] - \langle \tau, a \rangle + t \end{aligned}$$

using (i). □

If the families  $\mathcal{E}_{\mu, \mathcal{E}}$  and  $\mathcal{E}_{\nu, \Gamma}$  intersect then  $\mathcal{E}_\mu$  and  $\mathcal{E}_\nu$  intersect, and then  $\mathcal{E}_\mu = \mathcal{E}_\nu$  by (ii) and (iii). In addition,  $Q_{\mu, \vartheta} = Q_{\nu, \vartheta + \tau}$  holds for all  $\vartheta \in \text{dom}(\Lambda_\mu)$  with a unique  $\tau$  in  $\text{lin}(\mu) = \text{lin}(\nu)$ . This  $\tau$  is denoted in the sequel by  $\tau_{\mu, \nu}^*$ .

Recall that the mapping  $\theta_{\mu, \text{dom}(\Lambda_\mu)}^*$  is defined on  $\text{ri}(\mu)$  that equals  $\text{ri}(\text{dom}(\Lambda_\mu^*))$ , see (13).

**Proposition 7.2** *Suppose  $\mathcal{E}_\mu = \mathcal{E}_\nu$ .*

- (i) *For  $a \in \mathbb{R}^d$ , there exists  $\theta \in \text{dom}(\Lambda_\mu)$  satisfying  $\Lambda_\mu^*(a) = \langle \theta, a \rangle - \Lambda_\mu(\theta)$  if and only if  $\vartheta \in \text{dom}(\Lambda_\nu)$  exists such that  $\Lambda_\nu^*(a) = \langle \vartheta, a \rangle - \Lambda_\nu(\vartheta)$ , in which case  $Q_{\mu, \theta} = Q_{\nu, \vartheta}$ .*
- (ii) *If  $a \in \text{ri}(\mu) = \text{ri}(\nu)$ ,  $\theta = \theta_{\mu, \text{dom}(\Lambda_\mu)}^*(a)$  and  $\vartheta = \theta_{\nu, \text{dom}(\Lambda_\nu)}^*(a)$  then  $\vartheta = \theta + \tau_{\mu, \nu}^*$  and  $Q_{\mu, \theta} = Q_{\nu, \vartheta}$ .*
- (iii)  *$\text{dom}(\Lambda_\mu^*) = \text{dom}(\Lambda_\nu^*)$ , and for any  $a$  in this set the GMLE's  $R_{\mu, \text{dom}(\Lambda_\mu)}^*(a)$  and  $R_{\nu, \text{dom}(\Lambda_\nu)}^*(a)$  coincide.*

*Proof* The assumption implies that  $Q_{\mu, \xi} = Q_{\nu, \eta}$  for some  $\xi \in \text{dom}(\Lambda_\mu)$  and  $\eta \in \text{dom}(\Lambda_\nu)$ . These parameters are chosen to have  $\tau = \eta - \xi$  equal to  $\tau_{\mu, \nu}^*$ .

- (i) By Lemma 7.1 (i) and (iv), the equality  $\Lambda_\mu^*(a) = \langle \theta, a \rangle - \Lambda_\mu(\theta)$  rewrites to  $\Lambda_\nu^*(a) = \langle \theta + \tau, a \rangle - \Lambda_\nu(\theta + \tau)$ , which means that  $\theta$  satisfies the first equality if and only if  $\vartheta = \theta + \tau$  satisfies the second one. Here,  $\theta \in \text{dom}(\Lambda_\mu)$  if and only if  $\vartheta \in \text{dom}(\Lambda_\nu)$  by Lemma 7.1 (ii). Then (iii) of this lemma implies the equality of the pm's.
- (ii) As argued in Sect. 3.3,  $\Lambda_\mu^*(a) = \langle \theta, a \rangle - \Lambda_\mu(\theta)$  holds for the unique element  $\theta = \theta_{\mu, \text{dom}(\Lambda_\mu)}^*(a)$  of  $\text{dom}(\Lambda_\mu) \cap \text{lin}(\mu)$ . This,  $\tau \in \text{lin}(\mu)$  and the arguments of the proof of (i) imply the assertions.
- (iii) By Theorem 4.9,  $\text{dom}(\Lambda_\mu^*)$  equals  $\text{cc}(\mu) + \text{bar}(\text{dom}(\Lambda_\mu))$ . The first assertion follows since  $\text{cc}(\mu) = \text{cc}(\nu)$  and the barrier cones coincide on account of Lemma 7.1 (ii). The second one is a special case of Proposition 7.4 (iii) with  $\mathcal{E} = \text{dom}(\Lambda_\mu)$  and  $\Gamma = \text{dom}(\Lambda_\nu)$  since  $\text{dom}(\Lambda_\mu^*) \subseteq \text{aff}(\mu)$ .

□

**Lemma 7.3** *Under the assumption and notations of Lemma 7.1,*

- (v)  $\tau + \mathcal{E} \subseteq \text{dom}(\Lambda_\nu)$ ,
- (vi)  $\mathcal{E}_{\mu, \mathcal{E}} = \mathcal{E}_{\nu, \tau + \mathcal{E}}$ ,
- (vii)  $\Psi_{\mu, \mathcal{E}}^*(a) = \Psi_{\nu, \tau + \mathcal{E}}^*(a) - [\langle \tau, a \rangle - t]$ ,  $a \in \mathbb{R}^d$ ,
- (viii)  $\text{dom}(\Psi_{\mu, \mathcal{E}}^*) = \text{dom}(\Psi_{\nu, \tau + \mathcal{E}}^*)$ ,
- (ix)  $R_{\mu, \mathcal{E}}^*(a) = R_{\nu, \tau + \mathcal{E}}^*(a)$ ,  $a \in \text{dom}(\Psi_{\mu, \mathcal{E}}^*)$ .



*Proof* All references here are to Lemma 7.1. Obviously, (v) follows from (ii), and then, (vi) from (iii). A proof of (vii) is a simple variation of that of (iv) and is omitted. Clearly, (viii) follows from (vii). Finally, combining (vii), (i) and (iii), (2) rewrites to

$$\Psi_{\nu, \tau + \mathcal{E}}^*(a) + \Lambda_\nu(\vartheta + \tau) - \langle \vartheta + \tau, a \rangle \geq D(R_{\mu, \mathcal{E}}^*(a) \| Q_{\nu, \vartheta + \tau}), \quad \vartheta \in \mathcal{E}, \quad (22)$$

where  $a \in \text{dom}(\Psi_{\mu, \mathcal{E}}^*) = \text{dom}(\Psi_{\nu, \tau + \mathcal{E}}^*)$  by (viii). The uniqueness of GMLE applied to  $\nu$  and  $\tau + \mathcal{E}$  establishes (ix), see Remark 1.2.  $\square$

The assertions of Proposition 7.2 extend to general canonically convex exponential families only under the additional assumption that  $a$  belongs to  $\text{aff}(\mu) = \text{aff}(\nu)$ , see also Example 3.5.

**Proposition 7.4** *Suppose  $\mathcal{E}_{\mu, \mathcal{E}} = \mathcal{E}_{\nu, \Gamma}$  and  $a \in \text{aff}(\mu)$ .*

- (i) *There exists  $\theta \in \mathcal{E}$  satisfying  $\Psi_{\mu, \mathcal{E}}^*(a) = \langle \theta, a \rangle - \Lambda_\mu(\theta)$  if and only if  $\vartheta \in \Gamma$  exists such that  $\Psi_{\nu, \Gamma}^*(a) = \langle \vartheta, a \rangle - \Lambda_\mu(\vartheta)$ , in which case  $Q_{\mu, \theta} = Q_{\nu, \vartheta}$ .*
- (ii) *The conditions  $a \in \text{ri}(\mu) + \text{bar}(\mathcal{E})$  and  $a \in \text{ri}(\nu) + \text{bar}(\Gamma)$  are equivalent, and imply  $\vartheta = \theta + \tau_{\mu, \nu}^*$  and  $Q_{\mu, \theta} = Q_{\nu, \vartheta}$  where  $\theta = \theta_{\mu, \mathcal{E}}^*(a)$  and  $\vartheta = \theta_{\nu, \Gamma}^*(a)$ .*
- (iii) *The conditions  $a \in \text{dom}(\Psi_{\mu, \mathcal{E}}^*)$  and  $a \in \text{dom}(\Psi_{\nu, \Gamma}^*)$  are equivalent, and imply that  $R_{\mu, \mathcal{E}}^*(a)$  and  $R_{\nu, \Gamma}^*(a)$  coincide.*

*Proof* By the first assumption,  $\text{lin}(\mu) = \text{lin}(\nu)$  and  $Q_{\mu, \xi} = Q_{\nu, \eta}$  for some  $\xi \in \text{dom}(\Lambda_\mu)$  and  $\eta \in \text{dom}(\Lambda_\nu)$  with  $\eta - \xi$  equal to  $\tau = \tau_{\mu, \nu}^*$ .

- (i) By Lemma 7.3 (vi),  $\mathcal{E}_{\mu, \mathcal{E}} = \mathcal{E}_{\nu, \tau + \mathcal{E}}$ , and thus the first assumption implies that  $\pi_\nu(\tau + \mathcal{E}) = \pi_\nu(\Gamma)$ . Then,  $\Psi_{\mu, \mathcal{E}}^*(a) = \langle \theta, a \rangle - \Lambda_\mu(\theta)$  holds with some  $\theta \in \mathcal{E}$  if and only if  $\Psi_{\nu, \tau + \mathcal{E}}^*(a) = \langle \theta + \tau, a \rangle - \Lambda_\mu(\theta + \tau)$ , on account of Lemma 7.3 (vii) and Lemma 7.1 (i). By Corollary 2.12 and  $a \in \text{aff}(\nu)$ , this is equivalent to  $\Psi_{\nu, \pi_\nu(\tau + \mathcal{E})}^*(a) = \langle \theta + \tau, a \rangle - \Lambda_\mu(\theta + \tau)$  and, in turn, to  $\Psi_{\nu, \Gamma}^*(a) = \langle \vartheta, a \rangle - \Lambda_\mu(\vartheta)$  with some  $\vartheta \in \Gamma$  satisfying  $\pi_\nu(\theta + \tau) = \pi_\nu(\vartheta)$ . Hence,  $\pi_\nu(\theta) = \pi_\nu(\vartheta) - \tau_{\mu, \nu}^*$ , and the equality of pm's follows from Lemma 7.1 (iii).
- (ii) By  $a \in \text{aff}(\mu)$  and Lemma 4.4,  $a$  belongs to  $\text{ri}(\mu) + \text{bar}(\mathcal{E})$  if and only if it belongs to  $\text{ri}(\mu) + \text{bar}(\pi_\mu(\mathcal{E}))$ . Since  $\text{ri}(\mu) = \text{ri}(\nu)$  and the barrier cones of  $\pi_\mu(\mathcal{E})$  and  $\pi_\nu(\Gamma)$  coincide on account of  $\pi_\nu(\tau + \mathcal{E}) = \pi_\nu(\Gamma)$ , the first assertion obtains by applying again Lemma 4.4 with  $\mu$  replaced by  $\nu$  and  $a \in \text{aff}(\nu)$ . The second assertion is a consequence of  $\pi_\nu(\theta) = \pi_\nu(\vartheta) - \tau_{\mu, \nu}^*$ , obtained in (i) above, and the uniqueness of  $\theta$  and  $\vartheta$  in  $\text{lin}(\mu)$ , by Theorem 3.1.
- (iii) The first assertion can be proved analogously to the first assertion of (ii) above, by Theorem 4.9 and  $\text{cc}(\mu) = \text{cc}(\nu)$ . For the equality of GMLE's, rewrite (22) to

$$\Psi_{\nu, \Gamma}^*(a) + \Lambda_\nu(\vartheta) - \langle \vartheta, a \rangle \geq D(R_{\mu, \mathcal{E}}^*(a) \| Q_{\nu, \vartheta}), \quad \vartheta \in \Gamma,$$

using Corollary 2.12,  $\pi_\nu(\tau + \mathcal{E}) = \pi_\nu(\Gamma)$  and  $a \in \text{aff}(\nu)$ , and recall the uniqueness of GMLE.  $\square$

## References

1. Barndorff-Nielsen, O.: Information and Exponential Families in Statistical Theory. Wiley, New York (1978)
2. Brown, L.D.: Fundamentals of Statistical Exponential Families. Inst. of Math. Statist. Lecture Notes–Monograph Series, Vol. 9 (1986)
3. Chentsov, N.N.: Statistical Decision Rules and Optimal Inference. Translations of Mathematical Monographs, Amer. Math. Soc., Providence–Rhode Island, 1982 (Russian original: Nauka, Moscow, 1972)
4. Csiszár, I., Matúš, F.: Convex cores of measures on  $\mathbb{R}^d$ . *Studia Sci. Math. Hungar.* **38**, 177–190 (2001)
5. Csiszár, I., Matúš, F.: Information closure of exponential families and generalized maximum likelihood estimates. In: Proc. 2002 IEEE Int. Symp. Inform. Theory, p. 434 (2002)
6. Csiszár, I., Matúš, F.: Information projections revisited. *IEEE Trans. Inform. Theory* **49**, 1474–1490 (2003)
7. Csiszár, I., Matúš, F.: On information closures of exponential families: a counterexample. *IEEE Trans. Inform. Theory* **50**, 922–924 (2004)
8. Csiszár, I., Matúš, F.: Closures of exponential families. *Ann. Probab.* **33**, 582–600 (2005)
9. Csiszár, I., Matúš, F.: Generalized maximum likelihood estimates for infinite dimensional exponential families. In: Proceedings Prague Stochastics 2006. Prague, Czech Republic, pp. 288–297 (2006)
10. Eriksson, N., Fienberg, S.E., Rinaldo, A., Sullivant, S.: Polyhedral conditions for the nonexistence of the MLE for hierarchical log-linear models. *J. Symbol. Comput.* **41**, 222–233 (2006)
11. Letac, G.: Lectures on Natural Exponential Families and their Variance Functions. *Monografias de Matemática 50*. Instituto de Matemática Pura e Aplicada, Rio de Janeiro (1992)
12. Lauritzen, S.L.: Graphical Models. Clarendon, Oxford (1996)
13. Rinaldo, A.: On maximum likelihood estimation in log-linear models. Technical Report 833, Department of Statistics, Carnegie Mellon University (2006)
14. Rinaldo, A.: Computing maximum likelihood estimates in log-linear models. Technical Report 835. Department of Statistics, Carnegie Mellon University (2006)
15. Rockafellar, R.T.: Convex Analysis. Princeton University Press, Princeton (1970)