

## CHARACTERIZING PRIORS BY POSTERIOR EXPECTATIONS IN MULTIPARAMETER EXPONENTIAL FAMILIES

THEOPHILOS CACOULLOS

(Received Jan. 20, 1986; revised Apr. 7, 1986)

### Summary

Exploiting the notion of identifiability of mixtures of exponential families with respect to a vector parameter  $\theta$ , it is shown that the posterior expectation of  $\theta$  characterizes the prior distribution of  $\theta$ . The result is applied to normal and negative multinomial distributions.

### 1. Introduction

Priors play a fundamental role in Bayesian statistics. Their importance is enhanced by the fact that the posterior mean of a parameter turns out to be its Bayes estimator (predictor) under quadratic loss. It is therefore important to be able to recover the prior underlying a Bayes predictor.

The problem of characterizing the prior distribution of a parameter  $\theta$  in a given family of distributions  $f(\cdot|\theta)$  has been the object of several papers. In particular, results are available for the one-parameter exponential family

$$(1.1) \quad f(x|\theta) = a(x)b(\theta)e^{\theta(x)}$$

where, as a rule, the posterior expectation  $m(x) = E[\theta|X=x]$  characterizes the prior  $\pi(\theta)$  (see Cacoullos and Papageorgiou [3], [4] and references therein). A rigorous treatment of the multiparameter regular exponential family, (2.1) below, is given by Diaconis and Ylvisaker [5]; however, they are mainly concerned with conjugate priors which are characterized by linear posterior expectations  $m^*(x) = E\{E(X|\theta)|X=x\} = ax + b$ .

In point of fact, under rather general conditions, it is possible to characterize a large family of priors, not necessarily conjugate, relaxing at the same time the condition of linearity of  $m(x)$ . This is made

Key words and phrases: Characterizations, priors, exponential family, posterior expectation.

possible by using the concept of identifiability of mixtures (Teicher [7] and Barndorff-Nielsen [1]) and exploiting some known results on it. That is, if it is possible to determine the absolute distribution  $f(\cdot)$ ,

$$(1.2) \quad f(x) = \int f(x|\theta)\pi(\theta)d\theta$$

for certain  $f(\cdot|\theta)$ , then, by identifiability,  $\pi(\cdot)$  is also uniquely determined. Here  $f(\cdot)$  is determined from  $f(\cdot|\theta)$  and  $m(x)$ .

This approach was shown to be an efficient tool in characterizing prior (mixing) distributions for continuous [4] (discrete [3]) mixtures when both  $x$  and  $\theta$  are one-dimensional.

This note is motivated by the preceding remarks and the fact that, Teicher [7], the mixture (1.2), in general, is not identifiable when  $x$  is one-dimensional and  $\theta$  is a vector parameter. This is exemplified by the case of mixtures of a normal distribution  $N(\mu, \sigma^2)$  with respect to both parameters  $\mu$  and  $\sigma^2$ . In this situation, the identifiability of the mixture requires a sample of at least two observations. In general, the identifiability of mixtures of  $d$ -parameter exponential families requires at least  $d$  observations. However, both the sample space and the parameter space can be taken as open sets in  $R^d$ —the  $d$ -dimensional space of the sufficient statistic (see Diaconis and Ylvisaker [5], p. 271, where the sample space  $\mathcal{X}=\{x\}$  is an open set in  $R^d$ , and the parameter set  $\Theta = \left\{ \theta \mid \int e^{\theta \cdot x} d\mu(x) \right\}$  is also an open (and convex) set in  $R^d$ , the so called *natural parameter space*). Our approach below makes use of this parametrization.

## 2. The main result

In view of the remarks of the preceding section, without loss of generality, we can restrict our attention to a regular  $d$ -parameter exponential family of distributions, with probability density

$$(2.1) \quad f(x|\theta) = a(\theta)b(x)e^{\theta \cdot x},$$

where  $\theta = (\theta_1, \theta_2, \dots, \theta_d)$  and  $x = (x_1, \dots, x_d)$  plays the role of the sufficient statistic (a dot  $(\cdot)$  indicates the scalar product of two vectors). The main result can now be stated as follows.

**THEOREM 1.** *Let  $X = (X_1, X_2, \dots, X_d)$  be a random vector having an absolutely continuous distribution with density (2.1) and  $\pi(\theta)$  a prior distribution of  $\theta$  so that the continuous mixture defined by*

$$(2.2) \quad f(x) = \int f(x|\theta)\pi(\theta)d\theta$$

is identifiable. Then the posterior expectation vector

$$(2.3) \quad m(x) = E[\theta | X=x]$$

determines the (absolute) density  $f(\cdot)$  of  $X$ ; hence  $\pi(\theta)$  is also characterized.

PROOF. Differentiating  $f(x|\theta)$  w.r.t.  $x_i, i=1, \dots, d$ , gives

$$\nabla f(x|\theta) = \frac{\nabla b(x)}{b(x)} f(x|\theta) + \theta f(x|\theta), \quad \left( \nabla f = \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_d} \right) \right).$$

Multiplying both sides by  $\pi(\theta)$  and then integrating out  $\theta$  yields, after interchanging differentiation and integration,

$$(2.4) \quad \frac{\nabla f(x)}{f(x)} = \frac{\nabla b(x)}{b(x)} + m(x).$$

The solution of (2.4) determines the (marginal) distribution of  $X$ . By the identifiability of the mixture (2.2), the prior  $\pi(\theta)$  is also characterized.

*Remark 1.* As already mentioned, here the posterior mean  $m(x)$  is not restricted to be a linear function. However, it should be noted that the posterior mean vector  $m(x)$  can be an arbitrary function only to the extent that it is consistent with a given  $f(x|\theta)$  and a chosen prior  $\pi(\theta)$ . Here our primary concern is the possibility of identifying the prior from such a suitable (consistent)  $m(x)$ , not necessarily linear. The solvability of (2.4) requires, for example, that  $m(x)$  satisfy the condition

$$(2.5) \quad \frac{\partial m_i(x)}{\partial x_j} = \frac{\partial m_j(x)}{\partial x_i} \quad \text{for all } i \neq j.$$

The importance of Theorem 1 lies in the possibility of characterizing an infinite variety of priors, namely, all the priors which define an identifiable mixture in (2.1). Thus, once one chooses any such prior, this determines an  $m(x)$ , which in turn (working backwards) characterizes the prior.

*A characterization of  $N(\mu, \sigma^2)$ .* Let  $Y_1, Y_2, \dots, Y_n$  be a random sample of size  $n > 1$  from a normal  $N(\mu, \sigma^2)$  and consider the sufficient statistic  $X = (X_1, X_2)$  for  $\theta = (\theta_1, \theta_2)$  where we set (see [1], p. 121)

$$(2.6) \quad X_1 = \sum_{i=1}^n Y_i, \quad X_2 = -\frac{1}{2} \sum_{i=1}^n Y_i^2, \quad \theta_1 = \frac{\mu}{\sigma^2}, \quad \theta_2 = \frac{1}{\sigma^2}.$$

Then the probability density (with respect to Lebesgue measure in  $R^n$ )

of  $Y=(Y_1, \dots, Y_n)$  given  $\theta$  can be written as

$$(2.7) \quad f(y|\mu, \sigma^2) = f^*(x_1, x_2|\theta_1, \theta_2) = \left(\frac{\theta_2}{2\pi}\right)^{n/2} \exp\left(-\frac{n\theta_1^2}{4\theta_2}\right) \exp(\theta_1 x_1 + \theta_2 x_2)$$

so that in (2.1)  $b(x) \equiv 1$ . Thus, by Corollary 4 of [1], the mixture of (2.7) under some prior  $\pi(\theta)$  of  $\theta$  is identifiable. For a concrete application of Theorem 2.1 in terms of  $y=(y_1, \dots, y_n)$ , we use the following result (see e.g., Berger [2], p. 158).

LEMMA 1. Let  $Y_1, \dots, Y_n$ ,  $n \geq 2$ , be a random sample from  $N(\mu, \sigma^2)$  and suppose the joint prior density of  $\mu$  and  $\sigma^2$  is given by

$$(2.8) \quad \pi(\mu, \sigma^2) = \pi_1(\mu|\sigma^2)\pi_2(\sigma^2),$$

where  $\pi_1(\mu|\sigma^2)$  is  $N(\mu_0, \lambda\sigma^2)$  and  $\pi_2(\sigma^2)$  is an inverse gamma,  $IG(\alpha_0, \beta_0)$  with parameters  $\alpha_0 > 0$ ,  $\beta_0 > 0$ , i.e., with density

$$g^*(t|\alpha_0, \beta_0) = [\Gamma(\alpha_0)\beta_0^\alpha]^{-1} t^{-(\alpha_0+1)} e^{-1/\beta_0 t}.$$

Then

(a) the joint posterior density of  $\mu$  and  $\sigma^2$  given  $Y=y$  is

$$(2.9) \quad \pi(\mu, \sigma^2|y) = \pi_1(\mu|\sigma^2, y)\pi_2(\sigma^2|y)$$

where  $\pi_1(\mu|\sigma^2, y)$  is  $N(\mu(y), \sigma_1^2)$ , with

$$(2.10) \quad \mu(y) = (\mu_0 + n\lambda\bar{y})/(n\lambda + 1), \quad \sigma_1^2 = (\lambda^{-1} + n)\sigma^2,$$

and  $\pi_2(\sigma^2|y)$  is an  $IG(\alpha_1, \beta_1)$  with

$$(2.11) \quad \alpha_1 = \alpha_0 + \frac{n-1}{2}, \quad \beta_1^{-1} = \beta_0^{-1} + \frac{1}{2} \sum_{i=1}^n (y_i - \bar{y})^2 + \frac{n(\bar{y} - \mu_0)^2}{2(1+n\lambda)}.$$

(b) The marginal posterior density of  $\mu$  given  $y$  is a  $t$ -distribution with  $2\alpha_0 + n - 1$  degrees of freedom, location parameter  $\mu(y)$  and scale parameter  $[(\lambda^{-1} + n)(\alpha_0 + (n-1)/2)\beta_1]^{-1}$ , where  $\mu(y)$  and  $\beta_1$  are defined in (2.10) and (2.11).

(c) The marginal posterior density of  $\sigma^2$  given  $y$  is  $\pi_2(\sigma^2|y)$ , as previously defined.

As a corollary of Lemma 1 and Theorem 1, we obtain

PROPOSITION 1. In the notation of (2.7)–(2.11), for some real  $\mu_0$ ,  $\lambda > 0$ ,  $\alpha_0 > 0$ ,  $\beta_0 > 0$ , we have

$$(2.12) \quad m_1(y) = E(\mu|y) = \frac{\mu_0 + n\lambda\bar{y}}{n\lambda + 1}, \quad m_2(y) = E(\sigma^2|y) = \frac{1}{(\alpha_1 - 1)\beta_1},$$

if, and only if, the prior density  $\pi(\mu, \sigma^2)$  of  $(\mu, \sigma^2)$  is given by (2.8) and

the absolute density of  $Y$  is

$$(2.13) \quad f(y) = \frac{f(y|\mu, \sigma^2)\pi(\mu, \sigma^2)}{\pi(\mu, \sigma^2|y)},$$

which is easily seen to be a multivariate  $t$ -distribution.

*Remark 2.* Observe that the posterior mean vector  $m(y) = (m_1(y), m_2(y))$  is not linear in  $y$ ; nor is it so in terms of the sufficient statistic  $(x_1, x_2)$  (see (2.6)); thus, it was possible to characterize the prior for  $\mu$  and  $\sigma^2$  (see [5], Theorem 3) without assuming a linear  $m(y)$  (it so happened here that  $\pi(\mu, \sigma^2|y)$  and  $\pi(\mu, \sigma^2)$  are conjugate).

In terms of the canonical sufficient statistic  $X = (X_1, X_2)$  and the natural parameter  $\theta = (\theta_1, \theta_2)$ , as defined in (2.6), we have the following corollary.

**PROPOSITION 2.** *Under the assumptions of Proposition 1,*

$$(2.14) \quad m_1(x) = E(\theta_1|x) = \frac{\mu_0 + \lambda x_1}{(n\lambda + 1)} \alpha_1 \beta_1, \quad m_2(x) = E(\theta_2|x) = \alpha_1 \beta_1$$

where  $\alpha_1$  and  $\beta_1$  are given by (2.11), that is, in terms of  $x = (x_1, x_2)$

$$\beta_1^{-1} = \frac{1}{\beta_0} - x_2 - \frac{1}{2} \frac{x_1^2}{n} + \frac{n(n/x_1 - \mu_0)^2}{2(1+n\lambda)},$$

if, and only if, the prior density of  $\theta = (\theta_1, \theta_2)$  is

$$\pi^*(\theta_1, \theta_2) = \pi_1^*(\theta_1|\theta_2)\pi_2^*(\theta_2)$$

where  $\pi_1^*(\theta_1|\theta_2)$  is  $N(\mu_0\theta_2, \lambda\theta_2)$  and  $\pi_2^*(\theta_2)$  is a gamma with shape parameter  $\alpha_0 > 0$  and scale parameter  $\beta_0$ .

The specification of the marginal density of  $X$  can be given as for  $Y$  in (2.13); the posterior density  $\pi^*(\theta_1, \theta_2|x)$  is easily found from (2.9) and clearly the posterior density  $\pi_2^*(\theta_2|x)$  of  $\theta_2$  given  $x$  is a gamma with parameters  $\alpha_1$  and  $\beta_1$ . The posterior  $\pi_1^*(\theta_1|x)$  is rather complicated. The same is true, of course, for the marginal distribution  $f^*(x)$  of  $X$ . Nonetheless, given  $f(x|\theta)$  and  $m_1(x)$  and  $m_2(x)$  as in (2.14), it is possible to identify both  $\pi^*(\theta_1, \theta_2)$  and  $f^*(x)$ .

### 3. The discrete case

The treatment of the multiparameter discrete exponential family of distributions is analogous to the continuous case; the operation of differentiation (or  $\mathcal{F}$ ) is replaced by the shift operator  $E = (E_1, \dots, E_d)$  with  $E_i$  the partial shift operator, i.e.,  $E_i f(x_1, \dots, x_i, \dots, x_d) = f(x_1, \dots,$

$x_i+1, \dots, x_a)$ .

We can now state the analogue of Theorem 1.

**THEOREM 2.** *With the notation and assumptions of Theorem 1 except that now  $x \in I_0$ , the set of all lattice points in  $R^d$  with nonnegative components, the posterior mean vector*

$$m(x) = E[\theta | X=x], \quad x \in I_0$$

characterizes both  $\pi(\theta)$  and  $f(x)$ .

**PROOF.**  $E$  operating on  $f(x|\theta)$  with respect to  $x$  gives

$$E f(x|\theta) = \frac{E b(x)}{b(x)} f(x|\theta) + \theta f(x|\theta),$$

and hence the analogue of (2.5):

$$(3.1) \quad \frac{E f(x)}{f(x)} = \frac{E b(x)}{b(x)} + m(x).$$

Hence, as in Theorem 1, the assertion follows.

A similar remark holds as for Theorem 1 and the analogue of (2.5) takes the form

$$(3.2) \quad E_i m_j(x) = E_j m_i(x) \quad \text{for all } i \neq j.$$

*Characterization of the Dirichlet distribution as a prior for the negative multinomial.* For an illustration of Theorem 2, consider the negative multinomial distribution with p.f.

$$(3.3) \quad f(x|p) = \frac{\Gamma(r+x_0)}{\Gamma(r) \prod_{i=1}^d \Gamma(x_i+1)} p_0^r \prod_{i=1}^d p_i^{x_i}$$

where we set

$$p_0 = 1 - \sum_{i=1}^d p_i, \quad x_0 = \sum_{i=1}^d x_i, \quad p = (p_1, p_2, \dots, p_d).$$

The mixture of  $f(x|\theta)$  w.r.t. a Dirichlet prior  $\pi(\theta)$  is identifiable by Proposition 1 of [1] and the assumptions of Theorem 2 are satisfied. Hence a Dirichlet prior  $\pi(\cdot)$ , with

$$(3.4) \quad \pi(p) = \frac{\Gamma(\alpha)}{\prod_{i=0}^k \Gamma(\alpha_i)} \prod_{j=0}^d p_j^{\alpha_j-1}, \quad \alpha = \sum_{j=0}^d \alpha_j \quad (\alpha_i > 0)$$

can be characterized by the posterior expectation

$$m(x) = E(p|x), \quad x \in I_0;$$

this is given by

$$(3.5) \quad m_i(x) = E(p_i|x) = \frac{x_i + \alpha_i}{x_0 + r + \alpha}, \quad i = 1, \dots, d.$$

The result is summarized in

PROPOSITION 3. Let  $X = (X_1, \dots, X_d)$  be an observation from (3.3). Then  $m(x)$  is given by (3.5) if, and only if, the prior for  $p$  is the Dirichlet (3.4) and the (marginal) density of  $X$  is a compound negative multinomial distribution, Mosimann [6], defined by

$$(3.6) \quad f(x_1, \dots, x_k) = \frac{\Gamma(\alpha)\Gamma(r+x_0)\Gamma(r+\alpha_0) \prod_{i=1}^d \Gamma(\alpha_i+x_i)}{\Gamma(\alpha_0)\Gamma(r)\Gamma(x_0+\alpha+r) \prod_{i=1}^d \Gamma(\alpha_i)x_i!}.$$

If the  $\alpha_i$  are integers, then (3.6) becomes a multiple generalized hypergeometric distribution.

UNIVERSITY OF ATHENS

#### REFERENCES

- [1] Barndorff-Nielsen, O. (1965). Identifiability of mixtures of exponential families, *J. Math. Anal. Appl.*, **12**, 115-121.
- [2] Berger, J. (1980). *Statistical Decision Theory*, Springer-Verlag.
- [3] Cacoullos, T. and Papageorgiou, H. (1983). Characterizations of discrete distributions by a conditional distribution and a regression function, *Ann. Inst. Statist. Math.*, **35**, 95-103.
- [4] Cacoullos, T. and Papageorgiou, H. (1984). Characterizations of mixtures of continuous distributions by their posterior means, *Scand. Actuarial J.*, 23-30.
- [5] Diaconis, P. and Ylvisaker, D. (1979). Conjugate priors for exponential families, *Ann. Statist.*, **7**, 269-281.
- [6] Mosimann, J. E. (1963). On the compound negative multinomial distribution and correlations among inversely sampled pollen counts, *Biometrika*, **50**, 47-54.
- [7] Teicher, H. (1960). On the mixture of distributions, *Ann. Math. Statist.*, **31**, 55-73.