

# The exponential statistical manifold: mean parameters, orthogonality and space transformations

GIOVANNI PISTONE<sup>1</sup> and MARIA PIERA ROGANTIN<sup>2</sup>

<sup>1</sup>*Department of Mathematics, Politecnico di Torino, Corso Duca degli Abruzzi 24, 10129 Torino, Italy. E-mail: pistone@polito.it*

<sup>2</sup>*Department of Mathematics, University of Genova, Via Dodecaneso 35, 16146 Genova, Italy. E-mail: rogantini@dima.unige.it*

Let  $(X, \mathcal{X}, \mu)$  be a measure space, and let  $\mathcal{M}(X, \mathcal{X}, \mu)$  denote the set of the  $\mu$ -almost surely strictly positive probability densities. It was shown by Pistone and Sempi in 1995 that the global geometry on  $\mathcal{M}(X, \mathcal{X}, \mu)$  can be realized by an affine atlas whose charts are defined locally by the mappings  $\mathcal{M}(X, \mathcal{X}, \mu) \supset \mathcal{U}_p \ni q \mapsto \log(q/p) + K(p, q) \in B_p$ , where  $\mathcal{U}_p$  is a suitable open set containing  $p$ ,  $K(p, q)$  is the Kullback–Leibler relative information and  $B_p$  is the vector space of centred and exponentially  $(p \cdot \mu)$ -integrable random variables. In the present paper we study the transformation of such an atlas and the related manifold structure under basic transformations, i.e. measurable transformation of the sample space. A generalization of the mixed parametrization method for exponential models is also presented.

*Keywords:* exponential families; exponential statistical manifolds; information; mean parameters; Orlicz spaces; orthogonality

## 1. Introduction

The present paper is devoted to mathematical developments connected to the so-called theory of statistical manifolds. As general references on the subject in book form we mention Amari (1985), Amari *et al.* (1987), Murray and Rice (1993) and Barndorff-Nielsen and Cox (1994). Other relevant references for the present paper are Rao (1945, 1949), Jeffreys (1946), Dawid (1975, 1977), Efron (1975, 1978), Madsen (1979), Amari (1982), Barndorff-Nielsen and Jupp (1989) and Kass (1989).

In the literature on statistical manifolds, the question of finding a suitable functional setting to a nonparametric extension of the geometric construction has been mentioned by many workers (see, for example, Dawid (1975, 1977), Amari (1982) and Murray and Rice (1993)). In those papers and books some fundamental ideas of the nonparametric theory

have been sketched but, as far as we know, no detailed formal construction has been published before that of Pistone and Sempi (1995).

Starting with unpublished seminars held in Lecce University by Pistone and Sempi in 1989, the following idea has been developed. The statistical object that induces the geometry is the exponential model with its particular form of the Fisher information (Efron 1975); so the starting point has to be a nonparametric definition of the exponential model. This definition in turn is related to the class of *exponentially integrable* random variables whose natural topology is given by the notion of Orlicz space for the exponential function. The present paper is dedicated to further developments of these ideas. Some of the results developed here in full detail were announced by Pistone and Rogantin (1994).

The content of the paper is as follows.

Sections 2 and 3 are mainly devoted to a new presentation of the previous results of Pistone and Rogantin (1990), Pistone and Rogantin (1994) and Pistone and Sempi (1995), where the proofs of the basic propositions have been given in detail. Some of those proofs rely on quite straightforward arguments from functional analysis, but nevertheless the programme sketched above is systematically developed in these references. A few new results are added here and also the presentation has been improved; we give a definition of tangent space, show the relevance of the Kullback information and recall the notion of submanifold.

The main results of the present paper are given in Sections 6 and 7, where we give a nonparametric version of the concept of *mixed parametrization* in exponential models, and finally the effect of space transformation on the exponential manifold. A paper by Rogantin (1996) has given further developments in the case of finite sample spaces, discussing the derivation of parameters' orthogonality for various classes of finite-state-space stochastic processes.

## 2. The exponential manifold

The definition of statistical manifold can be given in the framework of the theory of manifolds modelled on Banach spaces, as introduced for example by Lang (1995).

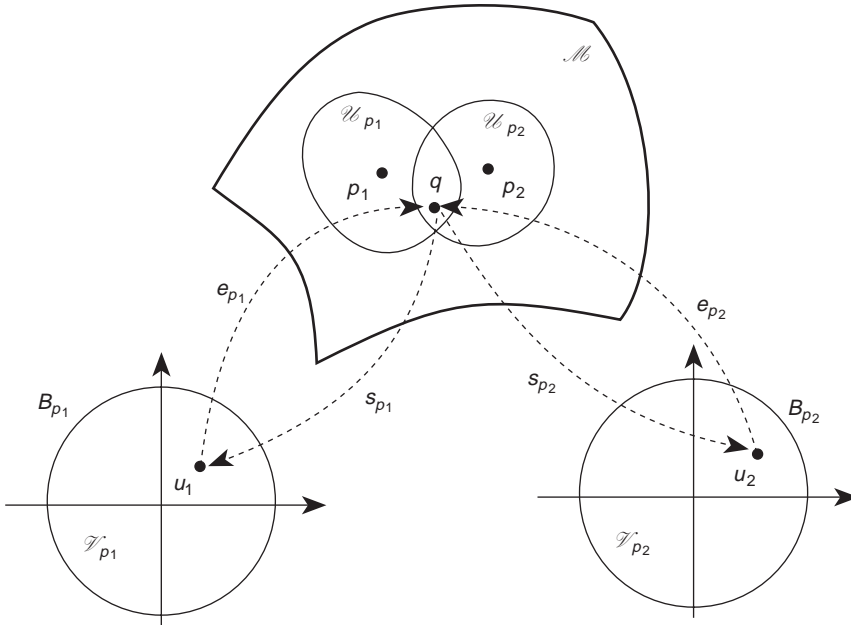
We consider a measure space  $(X, \mathcal{X}, \mu)$ , where  $\mu$  is a reference measure, and the set  $\mathcal{M}(X, \mathcal{X}, \mu)$  of the  $\mu$ -almost surely strictly positive probability densities. We shall define on the set  $\mathcal{M}(X, \mathcal{X}, \mu)$  a topology such that  $\mathcal{M}(X, \mathcal{X}, \mu)$  is an Hausdorff space (i.e. points can be separated by open sets). Then we shall construct a covering of  $\mathcal{M}(X, \mathcal{X}, \mu)$  with open sets  $\mathcal{U}_p$ ,  $p \in \mathcal{U}_p$ ,  $p \in \mathcal{M}(X, \mathcal{X}, \mu)$ , and a corresponding family of Banach spaces  $B_p$ , with norms  $\|\cdot\|_p$ ,  $p \in \mathcal{M}(X, \mathcal{X}, \mu)$ , such that each density  $q \in \mathcal{U}_p$  is represented with respect to  $p$  by a coordinate  $s_p(q) \in B_p$ .

We shall use the notation

$$s_p: \mathcal{U}_p \rightarrow \mathcal{F}_p \subset B_p, \quad (1)$$

$$e_p: \mathcal{F}_p \rightarrow \mathcal{U}_p \subset \mathcal{M}(X, \mathcal{X}, \mu), \quad (2)$$

to denote respectively the *charts*, i.e. the mappings from points to coordinates, and the *patches*, i.e. the mappings from coordinates to points (Figure 1).



**Figure 1.** The charts and the patches of the atlas.

The idea of a chart on  $U_p \subset \mathcal{M}(X, \mathcal{X}, \mu)$  is an abstraction of the parametrization for a family of probability densities (“statistical model”). Below we shall see that the condition of being a chart is actually more stringent than the condition of being a regular parametrization. Our model is nonparametric; so the coordinate mapping  $s_p$  cannot take values in a finite-dimensional vector space unless the sample space has a finite number of atoms.

As in differential geometry, we say that  $\{(\mathcal{U}_p, s_p) : p \in \mathcal{M}(X, \mathcal{X}, \mu)\}$  is an *atlas* if all the space is covered by its charts. If each of the *change in coordinates* (see Figure 1)

$$s_{p_2} \circ e_{p_1} : s_{p_1}(\mathcal{U}_{p_1} \cap \mathcal{U}_{p_2}) \rightarrow s_{p_2}(\mathcal{U}_{p_1} \cap \mathcal{U}_{p_2})$$

is a diffeomorphism of some regularity, the atlas is said to have that regularity. In such a case the atlas, possibly augmented by all the compatible charts, defines the manifold. Here a new chart is said to be compatible if all the corresponding changes of coordinate are regular (Lang 1995).

In our case we shall introduce a very special manifold, such that the change of coordinates are actually affine functions (i.e. they differ from a linear function by a additive constant), but we shall keep a weaker regularity, namely the  $C^\infty$  regularity (differentiability of any order) for compatible charts.

We shall denote by  $E_{p \cdot \mu}[\cdot]$  the expectation with respect to the probability measure  $p \cdot \mu$  (where  $(p \cdot \mu) dx = p(x)\mu dx$ ); if there is no ambiguity we shall use the notation  $E_p[\cdot]$ .

### 2.1. Functional framework

First we define the topology on  $\mathcal{M}(X, \mathcal{X}, \mu)$  as follows. For simplicity we give only the definition of convergence of sequences.

**Definition 1 (Exponential convergence).** *The sequence  $(p_n)_{n \in \mathbb{N}}$  in  $\mathcal{M}(X, \mathcal{X}, \mu)$  is e convergent (exponentially convergent) to  $p$  if  $(p_n)_{n \in \mathbb{N}}$  tends to  $p$  in  $\mu$  probability as  $n \rightarrow \infty$  and moreover the sequences  $(p_n/p)_{n \in \mathbb{N}}$  and  $(p/p_n)_{n \in \mathbb{N}}$  are eventually bounded in each  $L^\alpha(p)$ ,  $\alpha > 1$ , i.e.*

$$\forall \alpha > 1 \quad \limsup_{n \rightarrow \infty} E_p \left[ \left( \frac{p_n}{p} \right)^\alpha \right] < +\infty, \quad \limsup_{n \rightarrow \infty} E_p \left[ \left( \frac{p}{p_n} \right)^\alpha \right] < +\infty.$$

Some properties of the topology associated with this notion of convergence have been given by Pistone and Sempi (1995).

Now we shall introduce the Banach spaces on which the statistical manifold is modelled. We give a definition that shows how they are connected with well-known statistical objects (Barndorff-Nielsen 1978a; Letac 1992).

**Definition 2 (Cramér class).** *For each density  $p \in \mathcal{M}(X, \mathcal{X}, \mu)$ , the Cramér class at  $p$  is the set of all random variables  $u$  on  $(X, \mathcal{X}, \mu)$  such that the moment generating function of  $u$  with respect to the probability measure  $p \cdot \mu$  given by*

$$\hat{u}_p(t) = \int e^{tu} p \, d\mu = E_p[e^{tu}], \quad t \in \mathbb{R},$$

is finite in a neighbourhood of the origin 0.

If moreover the expectation of  $u$  is zero (the previous condition implies the existence of a finite expectation), then we shall call the set the centred Cramér class at  $p$ .

Note that, if  $p \in \mathcal{M}(X, \mathcal{X}, \mu)$  and  $p(\theta)$  is a one-dimensional exponential model,

$$p(\theta) = e^{\theta u - \psi(\theta)} p, \quad \theta \in I \text{ open real interval, } 0 \in I,$$

then the sufficient statistic  $u$  belongs to the Cramér class at each density  $p(\theta)$  in the model. In fact, for  $t \in I$ ,

$$\hat{u}_{p(\theta)}(t) = E_{p(\theta)}[e^{tu}] = E_p[e^{tu} e^{\theta u - \psi(\theta)}] = \frac{E_p[e^{(t+\theta)u}]}{e^{\psi(\theta)}} = \frac{e^{\psi(t+\theta)}}{e^{\psi(\theta)}}$$

and  $E_{p(\theta)}[e^{tu}]$  is finite for  $t$  in a neighbourhood of 0.

The following construction was presented by Pistone and Sempi (1995); we repeat it here for ease of presentation.

**Proposition 3 (A norm on the Cramér class).** *The Cramér class at  $p$  of Definition 2 is a vector space and a Banach space with the norm defined by*

$$\|u\|_p = \inf \left\{ r > 0: E_p \left[ \cosh \left( \frac{u}{r} \right) - 1 \right] \leq 1 \right\}. \tag{3}$$

The centred Cramér class, denoted by  $B_p$  and given by

$$B_p = \{u \in L^1(p \cdot \mu): 0 \in \text{dom}(\hat{u}_p)^\circ, E_p[u] = 0\},$$

is a closed subspace.

**Proof.** It is clear that  $E_p[e^{tu}] < +\infty$  in a neighbourhood of 0 if and only if  $E_p[e^{u/r}] < +\infty$  and  $E_p[e^{-u/r}] < +\infty$ , i.e.  $E_p[\cosh(u/r)] < +\infty$ . Moreover  $E_p[\cosh(|u|/r)] \rightarrow 1$  if  $r \rightarrow \infty$ . Consequently the Cramér class at  $p$  is defined by  $\|u\|_p < +\infty$ .

The Banach space property follows from general arguments on Orlicz spaces (Rao and Ren 1991; Krasnosel'skii and Rutickii 1961). □

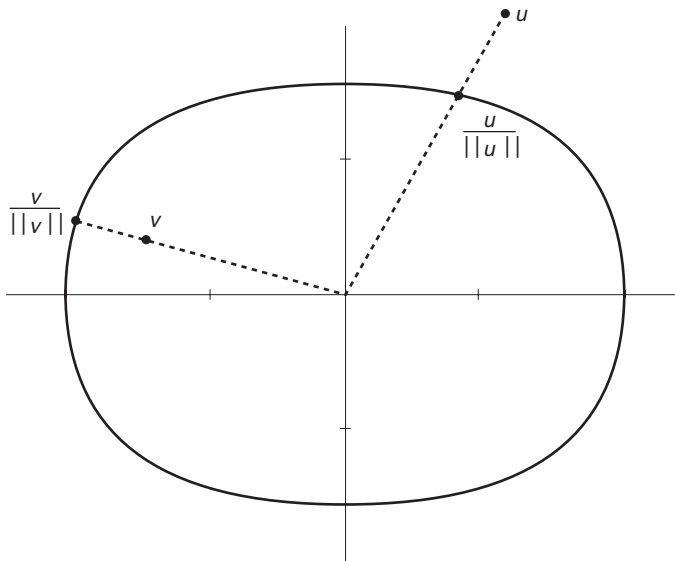
Note that  $\|u\|_p < 1$  if and only if there exists a real  $\alpha > 1$  such that

$$E_p[\cosh(\alpha u) - 1] \leq 1, \quad \text{i.e. } E_p[\cosh(\alpha u)] \leq 2.$$

Note also that, if  $u_n$  is a sequence of random variables, then  $\|u_n - u\|_p \rightarrow 0$  if and only if  $\forall \epsilon > 0, E_p[\cosh\{(u_n - u)/\epsilon\}] < 2$  eventually,  $n \rightarrow \infty$ .

**Example 4.** Figure 2 illustrates the construction of the norm  $\|\cdot\|_p$  when  $X = \{x_1, x_2\}$  and  $\mu$  is the counting measure. In such a case the space of the random variables  $u$  is  $\mathbb{R}^2$ . We consider the convex set  $\overline{\mathcal{F}}_p$  of all the vectors  $u = (u_1, u_2) \in \mathbb{R}^2$  such that

$$\overline{\mathcal{F}}_p = \{(u_1, u_2): \{\cosh(u_1) - 1\}\theta + \{\cosh(u_2) - 1\}(1 - \theta) \leq 1\},$$



**Figure 2.** Construction of the norm  $\|\cdot\|_p$ .

where  $\theta = p(x_1)$  and  $1 - \theta = p(x_2)$  (in Figure 2 we have assumed that  $\theta = \frac{1}{3}$ ). For any vector  $u \in \mathbb{R}^2$  this norm is the unique positive value  $r$  such that

$$\left\{ \cosh\left(\frac{u_1}{r}\right) - 1 \right\} \theta + \left\{ \cosh\left(\frac{u_2}{r}\right) - 1 \right\} (1 - \theta) = 1,$$

that is  $(u_1/r, u_2/r)$  lies on the boundary of  $\overline{\mathcal{F}}_p$ .

This example shows the construction of the norm defined in (3) in a finite sample space, where actually all norms define the same topology because the vector space of random variables has finite dimension. The construction that we present is really needed only when the underlying space of random variables is infinite dimensional.

In the previous proposition the function  $x \mapsto \cosh(x) - 1$  is a convex function that plays in the theory of the spaces  $B_p$  the same role as the function  $x \mapsto |x|^\alpha/\alpha$  in the theory of Lebesgue spaces  $L^\alpha$ ,  $\alpha > 1$ . We cite Krasnosel'skii and Rutickii (1961) and Rao and Ren (1991) as general references.

We shall use various types of such convex functions:

$$\phi_1: x \mapsto \cosh(|x|) - 1, \tag{4}$$

$$\phi_2: x \mapsto \exp(|x|) - |x| - 1, \tag{5}$$

$$\phi_3: x \mapsto (1 + |x|) \log(1 + |x|) - |x|. \tag{6}$$

For each of these functions it is possible to define a norm as in (3); we shall denote, for  $i = 1, 2, 3$ , the following.

- (a)  $\overline{\mathcal{F}}_{\phi_i, p}$  the convex set  $\{u \in L^1(p \cdot \mu): E_p[\phi_i(u)] \leq 1\}$ ;
- (b)  $\|\cdot\|_{\phi_i, p}$  the norm associated with  $\phi_i$ :

$$\|u\|_{\phi_i, p} = \inf \left\{ r > 0: E_p \left[ \phi_i \left( \frac{u}{r} \right) \right] \leq 1 \right\};$$

- (c)  $L^{\phi_i}(p \cdot \mu)$  (or  $L^{\phi_i}(p)$  if there is no ambiguity) the corresponding Banach spaces of non-centred random variables:

$$L^{\phi_i}(p) = \{u: \exists \alpha > 0 \text{ such that } E_p[\phi_i(\alpha u)] < +\infty\} = \{u: \|u\|_{\phi_i, p} < +\infty\};$$

- (d)  $L_0^{\phi_i}(p)$  the corresponding space of centred random variables.

The classes of centred random variables are closed subspaces. The function  $\phi_1$  is the most important for us; above, we have denoted  $\|\cdot\|_{\phi_1, p}$  by  $\|\cdot\|_p$  and  $L_0^{\phi_1}$  by  $B_p$ . We say that two  $\phi$  functions are equivalent if the corresponding norms are equivalent, i.e. there exist two real constants  $a$  and  $b$  such that  $a\|u\|_{\phi_i} \leq \|u\|_{\phi_j} \leq b\|u\|_{\phi_i}$ .

As the  $\phi$  functions in (4)–(6) are strictly convex and differentiable, it is possible to consider the inverse function  $[\phi']^{-1}(x)$  and to define the *conjugate function* of  $\phi$  as the function  $\psi$ , such that, for any  $y = \phi'(x)$ ,  $\psi'(y) = [\phi']^{-1}(x)$ . This implies that  $xy \leq \phi(x) + \psi(y)$ ,  $x, y \in \mathbb{R}_+$ , with equality if and only if  $y = \phi'(x)$ . If  $\phi$  and  $\psi$  are conjugate, then  $u \mapsto \sup_v \{E_p[uv]: E_p[\psi(v)] \leq 1\}$  is a norm equivalent to  $\|\cdot\|_{\phi, p}$  (Rao and Ren 1991, p. 61).

It can be shown that  $\phi_1$  and  $\phi_2$  are equivalent (and the corresponding Banach spaces coincide with the Cramer class at  $p$ ) and that  $\phi_2$  and  $\phi_3$  are conjugate.

Conjugacy implies that the Banach spaces are in a duality relation almost in the same way the Lebesgue spaces  $L^\alpha$  and  $L^\beta$  are in duality if  $\alpha^{-1} + \beta^{-1} = 1$ . Precisely the bilinear form

$$L^\phi(p) \times L^\psi(p) \ni (u, v) \mapsto E_p[uv] \in \mathbb{R}$$

is continuous, but in general  $L^\phi(p)$  and  $L^\psi(p)$  are not dual.

The following proposition follows immediately from the main result of Pistone and Sempi (1995). Because of its importance, we give here a new direct proof.

**Proposition 5.** *Let  $p$  and  $q$  be two probability densities in  $\mathcal{M}(X, \mathcal{X}, \mu)$  connected by a one-dimensional exponential model. Then*

$$L^{\phi_1}(p) = L^{\phi_1}(q).$$

**Proof.** Let  $r \in \mathcal{M}(X, \mathcal{X}, \mu)$  be given and let  $u \in L^{\phi_1}(r)$ . Let  $p(t) = e^{tu - \psi(t)} r$  be the one-dimensional exponential model associated with  $r$  and  $u$ , where  $t$  belongs to the real interval  $I$  with  $0 \in I$ . Let  $p(t_0) = p$  and  $p(t_1) = q$  be two densities in the given model. We can assume that  $t_0 < t_1$  because otherwise we change  $u$  with  $-u$ .

It is enough to show that  $L^{\phi_1}(p) \subset L^{\phi_1}(q)$  because the relation that connects  $p$  and  $q$  is symmetric.

Let  $w \in L^{\phi_1}(p)$  be given; we have to prove that there exists a  $\beta > 0$  such that

$$E_q[\cosh(\beta w)] < +\infty.$$

We have

$$\begin{aligned} E_q[\cosh(\beta w)] &= E_r[\cosh(\beta w) e^{t_1 u - \psi(t_1)}] \\ &= \frac{E_r[\cosh(\beta w) e^{t_1 u}]}{E_r[e^{t_1 u}]} \end{aligned}$$

The previous equation involves the convex function

$$g: (\theta, t) \mapsto E_r[\cosh(\theta w) e^{tu}].$$

We study the domain in which  $g$  is finite. If  $\theta = 0$  the value of the function is  $g(0, t) = E_r[e^{tu}]$  which is finite for  $t \in I$ . If  $t = t_0$  the value is

$$g(\theta, t_0) = E_r[\cosh(\theta w) e^{t_0 u}] = E_p[\cosh(\theta w)] E_r[e^{t_0 u}]$$

and is finite for  $\theta$  in some interval  $]-\bar{\theta}, \bar{\theta}[$ , because we have assumed that  $w \in L^{\phi_1}(p)$ . We chose  $b > t_1$ , such that  $b \in I$  and  $\gamma \in ]0, \bar{\theta}[$ . The function  $t \mapsto g(0, t)$  is finite in the interval  $[t_0, b]$  and the function  $\theta \mapsto g(\theta, t_0)$  is finite in the interval  $[-\gamma, \gamma]$ . Then, because of the convexity,  $(\theta, t) \mapsto g(\theta, t)$  is finite in the triangle with vertices at the points  $(t_0, \gamma)$ ,  $(b, 0)$  and  $(t_0, -\gamma)$ .

We consider the straight line between the points  $(t_0, \gamma)$  and  $(b, 0)$  and we denote by  $\beta$  the value of  $\theta$  at the intersection point of the previous straight line and the straight line  $\theta = t_1$ , namely  $\beta/\gamma = (b - t_1)/(b - t_0)$ . It follows that

$$E_q[\cosh(\beta w)] E_r[e^{t_1 u}] \leq \frac{t_1 - t_0}{b - t_0} E_p[\cosh(\gamma w)] E_r[e^{t_0 u}] + \frac{b - t_1}{b - t_0} E_r[e^{t_0 u}],$$

$$E_q[\cosh(\beta w)] \leq \frac{t_1 - t_0}{b - t_0} E_p[\cosh(\gamma w)] \frac{E_r[e^{t_0 u}]}{E_r[e^{t_1 u}]} + \frac{b - t_1}{b - t_0} \frac{E_r[e^{t_0 u}]}{[e^{t_1 u}]},$$

so that  $E_q[\cosh(\beta w)] < +\infty$ . □

Note that the equality between the spaces  $L^{\phi_1}(p)$  and  $L^{\phi_1}(q)$  holds true for each pair of points  $p, q$  which are connected by a one-dimensional exponential model. This establishes an equivalence relation, as we show below.

If  $p(t), t \in I$ , is a one-dimensional exponential model connecting  $p$  and  $q$ , and  $r(t), t \in J$ , is a one-dimensional exponential model connecting  $q$  and  $r$  we can assume those models to be of the form

$$p(t) = e^{tu - \psi_1(t)} q,$$

$$r(t) = e^{tv - \psi_2(t)} q,$$

with  $u, v \in L^{\phi_1}(q)$ . Then by convexity all the densities in the two-dimensional model

$$q(t_1, t_2) = e^{t_1 u + t_2 v - \phi_3(t_1, t_2)} q$$

are connected by a one-dimensional model; in particular  $p$  and  $r$  are connected.

**Definition 6 (x log x class).** We shall denote by  $*B_p$  the Banach space of centred random variables in  $L^{\phi_3}(p \cdot \mu)$ , i.e. the centred random variable of the so-called  $x \log x$  class.

**Proposition 7.** A  $(p \cdot \mu)$  integrable random variable  $u$  belongs to the  $x \log x$  class  $*B_p$  if and only if it is centred and  $(1 + |u|) \log(1 + |u|)$  is  $(p \cdot \mu)$  integrable.

**Proof.** If  $u \in L^{\phi_3}(p)$ , there exists a constant  $\alpha$  such that  $E_p[\phi_3(\alpha u)] < +\infty$ . This implies that  $E_p[\phi_3(u)] < +\infty$ ; in fact there exists a constant  $k$  such that, for any  $n$ ,  $\phi_3(2^n u) < k^n \phi_3(u)$  (Rao and Ren 1991, p. 22). For any  $\alpha$  there exists an  $n$  such that  $1/\alpha \leq 2^n$  and, because  $\phi_3$  is even and increasing in the positive real values,  $\phi_3(u) \leq \phi_3(2^n \alpha u) \leq k^n \phi_3(\alpha u)$ . So  $E_p[\phi_3(\alpha u)] < +\infty$  implies that  $E_p[\phi_3(u)] < +\infty$ . Finally both  $u \in L^{\phi_3}(p)$  and the integrability of  $(1 + |u|) \log(1 + |u|)$  imply that  $u \in L^1(p)$ , and the conclusion follows from  $E_p[u] = 0$ . □

Now we give some details about the Banach spaces  $B_p$  and  $*B_p$  which will be useful in the construction of the statistical manifold.



**Proposition 8.**

(a) All the elements  ${}^*u$  in  ${}^*B_p$  are identified with an element  $u^*$  of the dual space  $B_p^*$  of  $B_p$  by the formula:  $u^*(u) = E_p[{}^*uu]$ , with  $u \in B_p$ . In general,  ${}^*B_p$  is identified with a proper subset of  $B_p^*$ . The injection of  ${}^*B_p$  into  $B_p^*$  is continuous; we write

$${}^*B_p \subseteq B_p^*.$$

(b) All the elements  $u$  in  $B_p$  are identified with an element  $\bar{u}$  of the dual space  $({}^*B_p)^*$  of  ${}^*B_p$  by the formula  $\bar{u}({}^*u) = E_p[u {}^*u]$ , with  ${}^*u \in {}^*B_p$ . This identification is onto, i.e.  $B_p$  is identified with  $({}^*B_p)^*$ ; we write

$$({}^*B_p)^* \simeq B_p.$$

(c) The following continuous injections hold true:

$$L^\infty(p \cdot \mu) \subseteq B_p \subseteq \bigcap_{\alpha > 1} L_0^\alpha(p \cdot \mu) \subseteq {}^*B_p \subseteq B_p^*.$$

**Proof.** We recall some results about the Orlicz spaces; for further details see, for instance, Rao and Ren (1991).

Let  $\phi$  and  $\psi$  be conjugate functions. If  $f \in L^\phi(p)$  and  $g \in L^\psi(p)$ , then there exists a constant  $k$  such that  $\int |fg| p \, d\mu \leq k \|f\|_{\phi,p} \|g\|_{\psi,p}$  (Rao and Ren 1991, p. 58).

If  $u \in L^{\phi_2}(p)$  and  ${}^*u \in L^{\phi_3}(p)$ , then  $E_p[u {}^*u] \leq k \|u\|_{\phi_2,p} \|{}^*u\|_{\phi_3,p}$ . Thus for all  ${}^*u \in L^{\phi_3}(p)$  the mapping  $u \mapsto E_p[u {}^*u]$  is always defined, linear and continuous; it is an element of  $L^{\phi_2}(p)^*$ . Similarly for all  $u \in L^{\phi_2}(p)$  the mapping  ${}^*u \mapsto E_p[u {}^*u]$  is an element of  $L^{\phi_3}(p)^*$ .

$\phi_1$  and  $\phi_2$  are equivalent. Then, for any  ${}^*u \in L^{\phi_3}(p)$ ,  $(u \mapsto E_p[u {}^*u]) \in L^{\phi_1}(p)^*$  and for any  $u \in L^{\phi_1}(p)$ ,  $({}^*u \mapsto E_p[u {}^*u]) \in L^{\phi_3}(p)^*$ ,

$$L^{\phi_3}(p) \subseteq L^{\phi_1}(p)^* \quad \text{and} \quad L^{\phi_1}(p) \subseteq L^{\phi_3}(p)^*.$$

Moreover, if  $\phi$  is such that there exists  $k$  such that  $\phi(2x) \leq k\phi(x)$ , then  $L^\phi(p)^*$  is isometric to  $L^\phi(p)$  (Rao and Ren 1991, p. 111). We apply this to  $\phi_3$ ; then

$$L^{\phi_1}(p) \simeq L^{\phi_2}(p) = L^{\phi_3}(p)^*.$$

Now we have to show that the same properties extend to the centred spaces, i.e.

$${}^*B_p \subseteq B_p^* \quad \text{and} \quad B_p \simeq ({}^*B_p)^*.$$

Let  ${}^*u \in {}^*B_p$ . Then  ${}^*u \in L^{\phi_3}(p)$  and there exists  $u^* \in L^{\phi_1}(p)^*$  such that  $u^*(u) = E_p[{}^*uu]$ , for all  $u \in L^{\phi_1}(p)$ ; if  $u^*$  is restricted to  $B_p$ , then it is an element of  $B_p^*$  such that  $u^*(u) = E_p[{}^*uu]$ ,  $u \in B_p$ . The mapping  ${}^*u \rightarrow u^*$  is continuous from  ${}^*B_p$  to  $B_p^*$  because the restriction is a contraction.

The same argument applies to show that  $B_p \subseteq ({}^*B_p)^*$ .

Now let  $\bar{u} \in ({}^*B_p)^*$ .  $\bar{u}$  extends to a continuous linear operator  $\tilde{u}$  on  $L^{\phi_3}(p)$  defined by  $\tilde{u}({}^*u) = \bar{u}({}^*u - E_p[{}^*u])$ . Let  $u \in L^{\phi_1}(p)$  be the representation of  $\tilde{u}$ ; then  $\tilde{u}({}^*u) = E_p[u {}^*u]$ , with  ${}^*u \in L^{\phi_3}(p)$ . In particular, for  ${}^*u \in {}^*B_p$ ,  $\tilde{u}({}^*u) = \bar{u}({}^*u) = E_p[u {}^*u]$ .

Proposition 8(c) follows from the existence of constants  $k_1, k_2$  and  $k_3$  such that

$$\|\cdot\|_\infty \geq k_1 \|\cdot\|_{\phi_1} \geq k_2 \|\cdot\|_\alpha \geq k_3 \|\cdot\|_{\phi_3}.$$

Such inequalities depend on the different order of growth at  $\infty$  of the  $\phi$  functions involved (Rao and Ren 1991, p. 155). □

**Proposition 9.** *The multilinear mappings  $(u_1, \dots, u_n) \mapsto E_p[u_1 \dots u_n]$  with  $u_i \in B_p$ , are continuous; in particular the moments  $u \mapsto E_p[u^n]$  are continuous.*

**Proof.** This follows from the inclusion  $B_p \subseteq L^n(p)$ ; see Proposition 8(c). □

The Banach space  $B_p$  is an algebraic subspace of the Hilbert space  $L_0^2(p)$ , and its topology is stronger than the induced topology; in particular the scalar product of  $L_0^2(p)$  is defined on it and it is continuous.

**Definition 10 (Orthogonality in  $B_p$ ).** *The covariance induces a continuous scalar product on  $B_p$  defined as*

$$\langle u, v \rangle_p = E_p[uv] = \text{cov}_p[u, v] \text{ for all } u, v \in B_p.$$

*We say that  $u$  and  $v$  in  $B_p$  are orthogonal if  $\langle u, v \rangle_p = 0$ .*

The scalar product  $\langle \cdot, \cdot \rangle_p$  extends to the usual scalar product in  $L_0^2(p)$ . A different extension is possible, as is shown in the following definition.

**Definition 11 (Orthogonality on  $B_p^* \times B_p$ ).** *We shall denote by  $\langle \cdot, \cdot \rangle_{*,p}$  the bilinear form between the Banach space  $B_p$  and its dual Banach space  $B_p^*$ :*

$$B_p^* \times B_p \ni (u^*, u) \mapsto u^*(u) = \langle u^*, u \rangle_{*,p}.$$

*We say that  $u^* \in B_p^*$  is orthogonal to  $u \in B_p$  if  $\langle u^*, u \rangle_{*,p} = 0$ .*

As  $\langle \cdot, \cdot \rangle_{*,p}$  is not a scalar product, being defined on a product of different space, the previous definition is not consistent with the usual mathematical terminology. Nevertheless we suggest using the term ‘‘orthogonality’’ in this case because this notion is exactly what we need to discuss the notion of ‘‘orthogonal parametrization’’ in our framework; see below in Section 6.2.

We have seen in Proposition 8 that the Banach space  ${}^*B_p$  can be identified with a subspace of the dual space  $B_p^*$ . Thus, if  $u^* \in B_p^*$  is identified with  ${}^*u \in {}^*B_p$ , one has

$$\langle u^*, u \rangle_{*,p} = E_p[{}^*uu].$$

## 2.2. Analytical prerequisites for the construction of the atlas

The patches of the atlas will be defined on the open ball of radius 1:

$$\mathcal{I}_p = \{u \in B_p: \|u\|_p < 1\},$$

where we have denoted  $\|u\|_{\phi_1,p}$  by  $\|u\|_p$ .

**Proposition 12.** *If  $u \in \mathcal{V}_p$  and  $q = e^u p / E_p[e^u]$  then*

- (a) *the random variable  $e^u$  is  $(p \cdot \mu)$  integrable and  $q$  is a probability density in  $\mathcal{M}(X, \mathcal{X}, \mu)$  and*
- (b)  $L^{\phi_1}(p \cdot \mu) = L^{\phi_1}(q \cdot \mu)$ .

**Proof.**

(a) We remark that the condition  $\|u\|_p < 1$  is equivalent to the existence of an  $\alpha > 1$  such that  $E_p[\cosh(\alpha u) - 1] \leq 1$ , which in turn implies that  $E_p[e^u] \leq 4$ .

(b) The hypothesis implies that  $p$  and  $q$  are connected by a one-dimensional exponential model and the conclusion follows from Proposition 5. □

**Definition 13 (Moment generating functional).** *The moment generating functional*

$$G_p: L^{\phi_1}(p \cdot \mu) \rightarrow \bar{\mathbb{R}}_+ = [0, +\infty]$$

*is defined by*

$$G_p(u) = E_p[e^u].$$

**Proposition 14 (Properties of the moment generating functional).** *The moment generating functional  $G_p$*

(a) *takes the value 1 at 0; otherwise is strictly greater than 1, is convex and its proper domain  $\text{dom}(G_p) = \{u \in L^{\phi_1}(p \cdot \mu): G_p(u) < \infty\}$  is a convex set which contains the open unit ball of  $L^{\phi_1}(p \cdot \mu)$ ;*

(b) *is bounded and infinitely Fréchet differentiable on the open unit ball  $\mathcal{V}_p$  with differential*

$$D^n G_p(u)(v_1, \dots, v_n) = E_p[v_1 \dots v_n e^u].$$

**Proof.** See Pistone and Sempì (1995, Proposition 2.4). □

The previous abstract definition includes the usual definition of a (multivariate) moment generating function. In fact, let  $u = (u_1, \dots, u_n)$  be an  $n$ -dimensional random variable of the Cramer class under the density  $p$ . Then for each real vector  $t = (t_1, \dots, t_n)$  the linear combination  $\sum_{i=1}^n t_i u_i$  belongs to the Cramer class and its multivariate moment generating function  $G_p(t_1 \dots t_n)$  is equal to  $G_p(\sum_{i=1}^n t_i u_i)$ . We remark that the multivariate moments of order  $r_1, \dots, r_n$  of the random variables  $u_1, \dots, u_n$  are the value of the derivative at 0 of order  $r_1, \dots, r_n$  of  $G_p$  in the direction  $u_1, \dots, u_n$ :

$$E_p[u_1^{r_1} \dots u_n^{r_n}] = D^{r_1 + \dots + r_n} G_p(0) (u_1^{\circ r_1}, \dots, u_n^{\circ r_n})$$

where  $u^{\circ r}$  means  $\underbrace{u, \dots, u}_{r \text{ times}}$ .

**Definition 15 (Cumulant generating functional).** *The cumulant generating functional  $K_p: B_p \rightarrow [0, +\infty]$  is defined by*

$$K_p(u) = \log G_p(u).$$

We remark that we restrict the cumulant generating functional to be defined on *centred* random variables of the Cramér class at  $p$ .

**Proposition 16 (Cumulant).** *The cumulant generating functional  $K_p$  has proper domain  $\text{dom}(G_p) \cap B_p$ . If  $\mathcal{V}_p$  denotes the open ball of  $B_p$  of radius 1 then  $\mathcal{V}_p \subset \text{dom}(G_p) \cap B_p$ . Moreover  $K_p$  satisfies the following properties.*

(a)  $K_p$  is 0 at 0; otherwise is strictly positive, is convex and infinitely Fréchet differentiable on  $\mathcal{V}_p$ .

(b)  $\forall u \in \mathcal{V}_p$ ,  $q = e^{u-K_p(u)}$   $p$  is a probability density in  $\mathcal{M}(X, \mathcal{B}, \mu)$ . The value of the  $n$ th differential at  $u$  in the direction  $v$  ( $\in B_p$ ) of  $K_p$ , that is the  $n$ -linear continuous form

$D^n K_p(u)$  applied to  $\overbrace{(v, \dots, v)}^{n \text{ times}}$ , is the  $n$ th cumulant of  $v$  under the probability density  $q$ :

$$D^n K_p(u)v^n = \left. \frac{d^n}{dt^n} \log E_q[e^{tv}] \right|_{t=0}.$$

(c) For  $v, v_1$  and  $v_2$  in  $B_p$ , one has

$$DK_p(u)v = E_q[v], \tag{7}$$

$$D^2 K_p(u)(v_1, v_2) = E_q[v_1 v_2] - E_q[v_1] E_q[v_2] = \text{cov}_q[v_1, v_2].$$

(d)  $\forall u \in \mathcal{V}_p$  and  $q = e^{u-K_p(u)}$   $p$ , the random variable  $q/p - 1$  belongs to  ${}^*B_p$  and

$$DK_p(u)v = E_p \left[ \left( \frac{q}{p} - 1 \right) v \right], \quad v \in B_p.$$

In other words the differential of  $K_p$  at  $u$ ,  $DK_p(u)$ , is in  $B_p^*$  but actually is identified with an element of  ${}^*B_p$ , denoted by  $\nabla K_p(u)$ :

$$\nabla K_p(u) = e^{u-K_p(u)} - 1 = \frac{q}{p} - 1.$$

(e) The mapping  $B_p \ni u \mapsto \nabla K_p(u) \in {}^*B_p$  is monotonic, and in particular one to one.

(f) The weak derivative of the map  $B_p \ni u \mapsto \nabla K_p(u) \in {}^*B_p$  at  $u$  applied to  $w \in B_p$  is given by

$$D(\nabla K_p(u))w = \frac{q}{p}(w - E_q[w]),$$

and it is one to one at each point.

**Proof.**

- (a) For the first point see Pistone and Sempi (1995, Proposition 2.5(a)).
- (b) It is a consequence of the definition of  $K_p(u)$  and its properties, see Proposition 14.
- (c) Same proof as (b).
- (d) We note that the random variable  $q/p - 1$  is centred:  $E_p[q/p] - 1 = E_q[1] - 1 = 0$ .

To prove that  $q/p - 1 \in {}^*B_p$  we show, using the Proposition 7, that  $q/p \in L^{\phi_3}(p)$ , i.e.

$$E_p \left[ \left( 1 + \frac{q}{p} \right) \log \left( 1 + \frac{q}{p} \right) \right] < +\infty.$$

We have

$$\begin{aligned} E_p \left[ \left( 1 + \frac{q}{p} \right) \log \left( 1 + \frac{q}{p} \right) \right] &= E_p \left[ \log \left( 1 + \frac{q}{p} \right) \right] + E_p \left[ \frac{q}{p} \log \left( 1 + \frac{q}{p} \right) \right] \\ &\leq E_p \left[ \frac{q}{p} \right] + E_p \left[ \frac{q}{p} \log \left( 1 + \frac{q}{p} \right) \right] \\ &= 1 + E_p \left[ \frac{q}{p} \log \left( 1 + \frac{q}{p} \right) \right]. \end{aligned}$$

So we have to show that

$$E_p \left[ \frac{q}{p} \log \left( 1 + \frac{q}{p} \right) \right] < +\infty.$$

By the inequality  $x \log(1+x) \leq (1+x) \log^+(x) + 1$ , for  $x > 0$ , where  $\log^+(x) = \max\{\log(x), 0\}$ , we have

$$\begin{aligned} E_p \left[ \frac{q}{p} \log \left( 1 + \frac{q}{p} \right) \right] &\leq E_p \left[ \left( 1 + \frac{q}{p} \right) \log^+ \left( \frac{q}{p} \right) \right] + 1 \\ &= E_p \left[ \log^+ \left( \frac{q}{p} \right) \right] + E_p \left[ \frac{q}{p} \log^+ \left( \frac{q}{p} \right) \right] + 1. \end{aligned} \tag{8}$$

We know that

$$\log \left( \frac{q}{p} \right) = u - K_p(u) \in L^{\phi_1}(p) \subseteq L^1(p)$$

and also (see Proposition 12)

$$\log \left( \frac{q}{p} \right) \in L^{\phi_1}(q) \subseteq L^1(q).$$

Then all terms on the right-hand side of (8) are bounded. Since  $v \in B_p$  implies that  $E_p[v] = 0$ , the first equality in (7) may be written as

$$DK_p(u)v = E_p \left[ \left( \frac{q}{p} - 1 \right) v \right].$$

By definition the gradient  $\nabla K_p(u)$  is an element of  $B_p^*$  such that

$$DK_p(u)v = \langle \nabla K_p(u), v \rangle_{*,p}.$$

By comparing the two previous equations, we conclude that the gradient can be identified with  $q/p - 1 \in {}^*B_p$ . By a change in notation we write  $\nabla K_p(u) = q/p - 1 \in {}^*B_p$ .

(e) Note that, if  $u$  and  $\bar{u}$  are in  $B_p$  and  $\theta \in [0, 1]$ , the following relation holds:

$$\begin{aligned}
 E_p[(\nabla K_p(\bar{u}) - \nabla K_p(u))(\bar{u} - u)] &= \langle \nabla K_p(\bar{u}) - \nabla K_p(u), (\bar{u} - u) \rangle_{*,p} \\
 &= \left\langle \int_0^1 d\theta \frac{d}{d\theta} \nabla K_p((1 - \theta)u + \theta\bar{u}), (\bar{u} - u) \right\rangle_{*,p} \\
 &= \int_0^1 d\theta D^2 K_p((1 - \theta)u + \theta\bar{u})(\bar{u} - u, \bar{u} - u) \\
 &= \int_0^1 d\theta \text{var}_\theta(\bar{u} - u)
 \end{aligned} \tag{9}$$

where  $\text{var}_\theta$  is the variance with respect to the density  $p_\theta = e_p((1 - \theta)u + \theta\bar{u})$ . As the mapping  $\theta \mapsto D^2 K_p((1 - \theta)u + \theta\bar{u})$  is continuous and positive definite, then  $\theta \mapsto \text{var}_\theta(\bar{u} - u)$  is continuous and non-negative. If  $\nabla K_p(\bar{u}) = \nabla K_p(u)$ , then  $\text{var}_\theta(\bar{u} - u) = 0$  for any  $\theta$ . This implies that  $\bar{u} - u = \text{constant}$  and  $\bar{u} = u$ , because they are  $p$ -centred random variables.

(f) The map  $f: B_p \ni u \mapsto \nabla K_p(u) \in {}^* B_p$  is weakly differentiable if,  $\forall v \in B_p, u \mapsto E_p[f(u)v]$  is differentiable. We have

$$E_p[f(u)v] = E_p[\nabla K_p(u)v] = DK_p(u)v.$$

Now we consider the increment of  $\nabla K_p(u)$  in the direction  $w \in B_p$ :

$$\begin{aligned}
 E_p[(f(u + w) - f(u))v] &= DK_p(u + w)v - DK_p(u)v \\
 &= D^2 K_p(u)(w, v) + R_p(u, v, w),
 \end{aligned}$$

where  $|R_p(u, v, w)| = o(\|w\|_p)$ . We have shown that

$$D^2 K_p(u)(w, v) = \text{cov}_q[v, w] = E_p \left[ \frac{q}{p} (w - E_q[w])v \right], \tag{10}$$

so that  $E_p[DF(u)wv] = E_p[(q/p)(w - E_q[w])v]$ . To check that  $(q/p)(w - E_q[w])$  is an element of  ${}^* B_p$  we use (10) and the properties of the norms with respect to two conjugate functions. We want to prove that

$$\sup_{\|v\|_p < 1} E_p \left[ \frac{q}{p} (w - E_q[w])v \right] < +\infty.$$

Now  $\text{cov}_q[v, w] = D^2 K_p(u)(v, w) \leq k\|v\|_{\phi_{2,p}}\|w\|_{\phi_{2,p}}$  because  $K_p$  is twice differentiable at  $u$ . Finally  $\sup_{\|v\|_{\phi_{2,p}} \leq 1} \|v\|_p < +\infty$  because the norms  $\|\cdot\|_{\phi_{2,p}}$  and  $\|\cdot\|_p$  are equivalent. The weak derivative is one to one; in fact, if there exists  $v_1$  such that

$$\frac{q}{p}(v - E_q[v]) = \frac{q}{p}(v_1 - E_q[v_1]),$$

then  $v - v_1 = \text{constant}$  and  $v = v_1$  because they are centred random variables. □

Using the previous definitions and properties, it is possible to give a definition of the nonparametric exponential model as follows.

**Definition 17 (Maximal exponential model).** For each  $p$  in  $\mathcal{M}(X, \mathcal{X}, \mu)$  the maximal exponential model at  $p$  is the statistical model

$$\mathcal{E}_p = \{e^{u-K_p(u)} p : u \in \text{dom}(K_p)^\circ, E_p[u] = 0\}.$$

The function

$$B_p \supset \text{dom}(K_p)^\circ \ni u \mapsto e^{u-K_p(u)} p \in \mathcal{M}(X, \mathcal{X}, \mu)$$

is the likelihood function of the maximal exponential model;  $u$  plays the role of the “model parameter”.

Any parametric exponential model generated by  $p$  is embedded into the maximal exponential model as follows. Let  $(u_1, \dots, u_d)$  be the sufficient statistic of the parametric exponential model, with  $u_i \in B_p$ , and let

$$\Theta = \left\{ \theta = (\theta_1, \dots, \theta_d) \in \mathbb{R}^d : \sum \theta_i u_i \in \text{dom}(K_p)^\circ \right\}. \tag{11}$$

Then the parametric exponential model is

$$p_\theta = e^{\sum \theta_i u_i - \psi(\theta)} p \tag{12}$$

with  $\psi(\theta) = K_p(\sum \theta_i u_i)$ .

Note that the parametric exponential model is uniquely characterized by the linear subspace spanned by  $(u_1, \dots, u_d)$ . Note also that

$$\frac{\partial}{\partial \theta_j} \psi(\theta) = DK_p \left( \sum \theta_i u_i \right) u_j = E_{p_\theta}[u_j] = E_p \left[ \frac{p_\theta}{p} u_j \right] = E_p \left[ \left( \frac{p_\theta}{p} - 1 \right) u_j \right] = E_p[{}^* u_\theta u_j],$$

where  ${}^* u_\theta = p_\theta/p - 1 = \nabla K_p(u_\theta)$ ; the second equality follows from Proposition 16 and the last but one holds because the  $u_j$  are  $p$ -centred random variables.

We denote

$$\frac{\partial}{\partial \theta_j} \psi(\theta)$$

by  $\eta_j$ . The parameters  $(\eta_1, \dots, \eta_d)$  are the mean parameters of Barndorff-Nielsen and Cox (1994) or the mixture coordinates of Amari (1982).

### 2.3. The atlas

We now have all the elements for the definition of the atlas (see Figure 1). Let us consider the following map defined on a subset  $\mathcal{T}_p$  of the proper domain of  $K_p$ :

$$e_p : \mathcal{T}_p \ni u \mapsto q = e^{u-K_p(u)} p \in \mathcal{M}(X, \mathcal{X}, \mu), \tag{13}$$

where  $K_p(u) = \log E_p[e^u] = \log G_p(u)$  is the cumulant generating functional computed at  $u$ .

This mapping is one to one because  $u$  is centred; in fact, if  $u_1, u_2 \in \mathcal{T}_p$  and  $e^{u_1 - K_p(u_1)} = e^{u_2 - K_p(u_2)}$ , then  $u_1 - K_p(u_1) = u_2 - K_p(u_2)$ , and  $u_1 - u_2$  is constant and this constant has to be 0.

According to (1) and (2) we shall denote by  $\mathcal{U}_p$  the image of  $\mathcal{T}_p$  by the mapping  $e_p$  and by  $s_p$  the inverse of  $e_p$  on  $\mathcal{U}_p$ . Such an inverse,  $s_p: \mathcal{U}_p \rightarrow \mathcal{T}_p$ , is easily computed as

$$s_p: \mathcal{U}_p \ni q \mapsto \log\left(\frac{q}{p}\right) - E_p\left[\log\left(\frac{q}{p}\right)\right] \in \mathcal{T}_p. \tag{14}$$

The functions  $s_p, p \in \mathcal{M}(X, \mathcal{X}, \mu)$ , will be the coordinate mappings of our manifold in the sense that, locally around each  $p \in \mathcal{M}(X, \mathcal{X}, \mu)$ , each  $q \in \mathcal{U}_p$  will be “parametrized” by its *centred log-likelihood*.

Let us now compute the change-in-coordinates formula; if  $p_1$  and  $p_2$  are two points in  $\mathcal{M}(X, \mathcal{X}, \mu)$  such that  $\mathcal{U}_{p_1} \cap \mathcal{U}_{p_2} \neq \emptyset$ , then for all  $q$  in that intersection

$$\log\left(\frac{p_1}{p_2}\right) = \log\left(\frac{p_1}{q}\right) + \log\left(\frac{q}{p_2}\right)$$

belongs to  $L^{\phi_1}(p_1) = L^{\phi_1}(q) = L^{\phi_1}(p_2)$ . The composite transition mapping

$$s_{p_2} \circ e_{p_1}: s_{p_1}(\mathcal{U}_{p_1} \cap \mathcal{U}_{p_2}) \rightarrow s_{p_2}(\mathcal{U}_{p_1} \cap \mathcal{U}_{p_2})$$

simplifies to

$$s_{p_2} \circ e_{p_1}(u) = u + \log\left(\frac{p_1}{p_2}\right) - E_{p_2}\left[u + \log\left(\frac{p_1}{p_2}\right)\right], \tag{15}$$

where the algebraic computations are done in the space of  $\mu$  classes of measurable functions and the expectation is well defined as long as  $\mathcal{U}_{p_1} \cap \mathcal{U}_{p_2} \neq \emptyset$ .

**Theorem 18.** *The collection of pairs  $\{(\mathcal{U}_p, s_p): p \in \mathcal{M}(X, \mathcal{X}, \mu)\}$  is an affine  $C^\infty$  atlas on  $\mathcal{M}(X, \mathcal{X}, \mu)$ . The induced topology on sequences is equivalent to  $e$  convergence and the transition mappings are those defined in (15).*

**Proof.** This is the main result of Pistone and Sempi (1995, Theorem 3.6). □

Note that the derivative of the transition mapping defined in (15) is

$$B_{p_1} \ni u \mapsto u - E_{p_2}[u] \in B_{p_2},$$

and this is an isomorphism between  $B_{p_1}$  and  $B_{p_2}$  as topological linear spaces (Lang 1995). This extends Proposition 12 (b).

**Definition 19 (Manifold).** *The exponential (statistical) manifold is the manifold defined by the property in Theorem 18 on the set  $\mathcal{M}(X, \mathcal{X}, \mu)$ .*

The maximal exponential model defined in Definition 17 has a precise place in the general framework, as the following theorem shows.



**Theorem 20.** *The maximal exponential model  $\mathcal{E}_p$  is the connected component containing  $p$  of the exponential manifold  $\mathcal{M}(X, \mathcal{X}, \mu)$ .*

**Proof.** See Pistone and Sempi (1995, Theorem 4.1). □

The manifold structure that we have defined is special; many other types of atlas have been suggested in the literature, in particular the coordinates based on the mean parameters of Barndorff-Nielsen and Cox (1994) and the so-called Amari embeddings described by Amari (1982). In the infinite-dimensional case those different geometric structures are not equivalent to the exponential manifold, but in some restricted sense they are, because they induce the same manifold structure on finite-dimensional submanifolds (i.e. parametric statistical manifolds) (see again Amari (1982, 1985) and Murray and Rice (1993)).

### 3. Tangent space

A basic object of the theory of manifolds is the *tangent bundle*. In the case of the exponential statistical manifold it has been remarked from the very beginning (Dawid 1975) that there is a very natural identification between the tangent vectors and the exponential one-dimensional models around a point  $p$ . In fact each differentiable curve in  $\mathcal{M}(X, \mathcal{X}, \mu)$ , i.e. each one-dimensional statistical model  $p(t)$ ,  $t \in I \subseteq R$ , such that  $p(0) = p$ , has a tangent model of the exponential form  $e^{tu - K_p(tu)}$ . This prompts the identification of the tangent space with the set of one-dimensional exponential models.

Let  $p_1(t)$  and  $p_2(t)$  be two regular curves such that  $p_1(t_0) = p_2(t_0) = p$  and let  $u_1(t)$  and  $u_2(t)$  be the corresponding representations by a chart  $s_q$ , i.e.  $u_1(t) = s_q(p_1(t))$  and  $u_2(t) = s_q(p_2(t))$ . Then the curves  $p_1$  and  $p_2$  are equivalent at  $p$  if  $\dot{u}_1(t_0) = \dot{u}_2(t_0)$ . We denote by  $T_p \mathcal{M}$  the set of equivalence classes of regular curves through  $p$ . In local coordinates determined by the chart  $s_q$  the equivalence class corresponding to  $p(t)$  may be represented by the tangent vector to  $u(t)$  at  $t_0$ , that is by  $\dot{u}(t_0)$ , if  $p(t_0) = p$  and  $u(t) = s_q(p(t))$ . This definition does not depend on the chart chosen (Lang 1995).

**Proposition 21 (Tangent space).** *Let  $p(t)$  be a regular curve in  $\mathcal{M}(X, \mathcal{X}, \mu)$  with  $p(t_0) = p$ , and let  $u(t) \in B_q$  be its representation by a chart  $s_q$ , where  $t \in \{t: u(t) \in \text{dom}(K_q)^\circ\}$ . Then  $p(t) = e^{u(t) - K_q(u(t))} q$ .*

(a) *The relation between  $\dot{u}(t_0)$ , the tangent to  $u(t)$  at  $t_0$ , and the score function of  $p(t)$  with respect to the density  $p$  is*

$$\dot{u}(t_0) - E_q[\dot{u}(t_0)] = \left. \frac{d}{dt} \log \left( \frac{p(t)}{p} \right) \right|_{t=0}.$$

*If  $q = p$ , i.e. if the chart is centred at the same point where the log-likelihood is calculated, we have*

$$\dot{u}(t_0) = \left. \frac{d}{dt} \log \left( \frac{p(t)}{p} \right) \right|_{t=t_0}.$$

Then the space of the score function is a representation of the tangent space  $T_p \mathcal{M}$ .

(b) The curve

$$t \mapsto \frac{p(t)}{p} - 1$$

is in  ${}^*B_p$  and its weak derivative at  $t_0$  is  $\dot{u}(t_0)$ .

(c) The score function of any one-dimensional exponential model through  $p$ , i.e.

$$e^{tu - K_p(tu)} p,$$

at  $t = t_0$ , is  $u$  and vice versa any  $u \in B_p$  has such a corresponding one-dimensional exponential model. Then the space of the one-dimensional exponential models is another representation of the tangent space  $T_p \mathcal{M}$ .

**Proof.**

(a) We have

$$\begin{aligned} \left. \frac{d}{dt} \log \left( \frac{p(t)}{p} \right) \right|_{t=t_0} &= \left. \frac{d}{dt} \{u(t) - u(t_0) - K_q(u(t)) + K_q(u(t_0))\} \right|_{t=t_0} \\ &= \dot{u}(t_0) - E_p[\dot{u}(t_0)]. \end{aligned}$$

If  $p = q$ , then  $E_p[\dot{u}(t_0)] = E_q[\dot{u}(t_0)] = 0$ .

(b) From Proposition 16 (c) we have

$$\frac{p(t)}{p} - 1 = \nabla K_p(u(t)).$$

Its weak derivative (see Proposition 16 (f)) calculated at  $t = t_0$ , is

$$\left. \frac{\dot{p}(t_0)}{p} = D\nabla K_p(u(t)) \dot{u}(t) \right|_{t=t_0} = \left. \frac{p(t)}{p} (\dot{u}(t) - E_{p(t)}[\dot{u}(t)]) \right|_{t=t_0} = \dot{u}(t_0)$$

and it is in  ${}^*B_p$ .

(c) It follows by direct computation that

$$\left. \frac{d}{dt} \{tu - K_p(tu)\} \right|_{t=t_0} = u. \quad \square$$

The tangent space inherits the structure of vector space and the topology from  $B_p$ .

As a scalar product is defined on  $B_p$  (see Definition 10), a scalar product is defined on  $T_p \mathcal{M}$  together with the definition of orthogonality.

Let  $v \in T_p \mathcal{M}$  and let  $p(t)$  a regular curve whose tangent vector is  $v$ , i.e. a curve tangent at  $p$  to  $p(t) = e^{tv - K_p(tv)} p$ . Let  $\varphi: \mathcal{M}(X, \mathcal{X}, \mu) \rightarrow \mathbb{R}$ . The differential of  $\varphi$  at  $p$ , denoted by  $d_p \varphi$ , is a linear form on the tangent space:

$$d_p\varphi(v) = \left. \frac{d}{dt}\varphi(p(t)) \right|_{t=0}.$$

If  $p(t) = e_q(u(t))$ , then

$$d_p\varphi(v) = D\{\varphi \circ e_q(u(t_0))\}v.$$

### 3.1. Regular parametrization

Now we give a definition of a parametrization and we shall show an example of a parametrization that is not a chart in our sense.

**Definition 22.** Let  $A$  be an open set of the exponential statistical manifold  $\mathcal{M}(X, \mathcal{X}, \mu)$ , and let  $B$  be a Banach space. We shall say that  $F: A \rightarrow B$  is a  $C^k$  parametrization of  $A$ ,  $k = 1, 2, \dots, \infty$ , if  $F$  is a one-to-one,  $k$ -times continuously differentiable parameter and the tangent mapping  $d_pF$  is one to one at each  $p$ .

Note that the condition of being a  $C^\infty$  parametrization is weaker than the condition of being a chart (essentially because we do not require  $(d_pF)^{-1}$  to be continuous), unless the manifold is finite dimensional. In fact the regularity of the inverse mapping is not ensured by the conditions in Definition 22.

As an example, we consider the parameter shown in Proposition 16 based on the likelihood:  $q \mapsto q/p - 1$ . For any  $u \in \mathcal{U}_p$  there exists an  ${}^*u \in {}^*B_p$  such that  ${}^*u = \nabla K_p(u) = q/p - 1$ . The map  $u \mapsto {}^*u = q/p - 1$  is one to one and  ${}^*u$  belongs to  $x - \log x$  class (see Proposition 16); then this is a reasonable parametrization but it does not define a chart.

To see this, we denote by  $\mathcal{L}_p$  the image of  $\mathcal{U}_p$  in  ${}^*B_p$  under the mapping  $q \mapsto q/p - 1$ . If  $u \mapsto {}^*u = q/p - 1$ , a chart  $\mathcal{L}_p$  would be an open set of  ${}^*B_p$ . This is false; in fact a base of the neighbourhood of  ${}^*B_p$  is of the form  $\{{}^*u: E_p[\phi_3({}^*u)] < k\}$  while  ${}^*u \in \mathcal{L}_p$  satisfies the condition  ${}^*u \geq -1$  and, if the space does not have a finite number of atoms,  $\int ({}^*u + 1) \log ({}^*u + 1) p \, d\mu < +\infty$  does not imply that  ${}^*u \geq -1$ . So in the nonparametric theory the mean parameters do not define a system of charts.

Another example is the global parametrization  $p \mapsto p^{1/2} \in L^2(\mu)$ . This case has been discussed in more detail by Brigo and Pistone (1996).

## 4. Information

We now describe briefly how the notion of information (or entropy) is connected to the notion of the exponential statistical manifold.

**Definition 23 (Kullback–Leibler information).** Let the probability densities  $p$  and  $q$  in  $\mathcal{M}(X, \mathcal{X}, \mu)$  be given. If  $(q/p) \log(q/p)$  is  $(p \cdot \mu)$  integrable, then the Kullback–Leibler relative information of  $q$  with respect to  $p$  is the number

$$K(q, p) = \int \frac{q}{p} \log\left(\frac{q}{p}\right) p \, d\mu = E_p \left[ \frac{q}{p} \log\left(\frac{q}{p}\right) \right] = E_q \left[ \log\left(\frac{q}{p}\right) \right].$$

**Proposition 24.** Let  $p \in \mathcal{M}(X, \mathcal{X}, \mu)$ , let  $q \in \mathcal{U}_p$  and let  $u$  be the  $s_p$  coordinate of  $q$ , i.e.  $q = e^{u-K_p(u)} p$ . Let  ${}^*u \in {}^*B_p$  defined as  ${}^*u = q/p - 1$  and let

$$H_p({}^*u) = E_p[(1 + {}^*u) \log(1 + {}^*u)].$$

Then

- (a)  $K_p(u) = K(p, q),$
- (b)  $H_p({}^*u) = K(q, p),$
- (c)  $E_q[u] = K(p, q) + K(q, p) = K_p(u) + H_p({}^*u).$

**Proof.** We have

- (a)  $K(p, q) = -E_p \left[ \log\left(\frac{q}{p}\right) \right] = -E_p[u] + K_p(u) = K_p(u),$
- (b)  $H_p({}^*u) = E_p[(1 + {}^*u) \log(1 + {}^*u)] = E_p \left[ \frac{q}{p} \log\left(\frac{q}{p}\right) \right] = K(q, p),$
- (c)  $K(q, p) = E_q[u - K_p(u)] = E_q[u] - K_p(u) = E_q[u] - K(p, q). \quad \square$

Note that the value  $s_p(q)$  of the mapping  $s_p(\cdot)$  defined in (14) is the log-likelihood of  $q$  with respect to  $p$  plus the Kullback–Leibler relative information  $K(p, q)$  whose value is  $E_p[\log(p/q)]$

$$s_p(q) = \log\left(\frac{q}{p}\right) + K(p, q).$$

The mapping  $s_p$  is connected to the maximum-likelihood estimator as follows.

**Proposition 25 (Maximum expected log-likelihood).** Let  $p \in \mathcal{M}(X, \mathcal{X}, \mu)$ ,  $u \in \mathcal{V}_p$  and  $e_p(u) = e^{u-K_p(u)} p$ . Let  $\hat{u} \in \mathcal{V}_p$  and  $q = e^{\hat{u}-K_p(\hat{u})} p$ .

The maximum expected (at  $q$ ) log-likelihood, i.e. the maximum of the function,

$$\mathcal{V}_p \ni u \mapsto E_q \left[ \log\left(\frac{e_p(u)}{p}\right) \right]$$

is obtained at the point  $\hat{u}$  and

$$\max_u \left\{ E_q \left[ \log\left(\frac{e_p(u)}{p}\right) \right] \right\} = K(q, p).$$

**Proof.** We have

$$E_q \left[ \log \left( \frac{e_p(u)}{p} \right) \right] = E_q[u] - K_p(u).$$

The function  $u \mapsto E_q[u] - K_p(u)$  is concave in  $u$ . Its derivative at  $u$  in the direction  $v \in B_p$  is  $E_q[v] - DK_p(u)v = E_q[v] - E_{e_p(u)}[v]$ . Such a derivative is zero for all  $v$  if and only if  $q = e_p(u)$ ; then the maximum is obtained at  $\hat{u} = s_p(q)$  and consequently

$$\max_u \left\{ E_q \left[ \log \left( \frac{e_p(u)}{p} \right) \right] \right\} = E_q[\hat{u}] - K_p(\hat{u}) = K(q, p). \quad \square$$

**Proposition 26.** Let  $u \in B_p$ . Let  ${}^*u$  and  $H_p({}^*u)$  be as defined in Proposition 24.

The functions  $K_p(u)$  and  $H_p({}^*u)$  are conjugate convex functions. The relations of conjugacy are

$$H_p({}^*u) = \max_u \{ E_p[{}^*uu] - K_p(u) \},$$

$$K_p(u) = \max_{{}^*u} \{ E_p[{}^*uu] - H_p({}^*u) \}.$$

**Proof.** The first relation follows from Proposition 24 (b) and from Proposition 25:

$$\begin{aligned} H_p({}^*u) &= K(q, p) \\ &= \max_u \{ E_q[u] - K_p(u) \} \\ &= \max_u \left\{ E_p \left[ \frac{q}{p} u \right] - K_p(u) \right\} \\ &= \max_u \left\{ E_p \left[ \left( \frac{q}{p} - 1 \right) u \right] - K_p(u) \right\}. \end{aligned}$$

The last relation follows from the general theory of convex analysis (Ekeland and Temam 1974). □

The results of this section develop standard arguments on the entropy function (Kullback and Leibler 1951; Donsker and Varadhan 1975; Amari 1985; Kullback 1997 (some 38 years ago)).

## 5. Submanifold

**Definition 27 (Submanifold, submodel).** Let  $\mathcal{N}$  be a subset of the exponential manifold  $\mathcal{M}(X, \mathcal{X}, \mu)$  and, for each density  $p \in \mathcal{N}$ , let  $V_p^1$  and  $V_p^2$  be closed subspaces of  $B_p$ , such that there exist

(a) a homeomorphism (i.e. a linear invertible and bi-continuous mapping) between  $B_p$  and the direct product  $V_p^1 \times V_p^2$  (we say that  $V_p^1$  and  $V_p^2$  are split in  $B_p$ ) and

(b) a chart on a neighbourhood  $\mathcal{W}_p$  of  $p$ :

$$\sigma_p: \mathcal{W}_p \rightarrow B_p \simeq V_p^1 \times V_p^2,$$

where  $\sigma_p$  maps  $\mathcal{W}_p$  onto the product of two open sets  $\mathcal{V}_p^1 \times \mathcal{V}_p^2$  (with  $\mathcal{V}_p^1 \subset V_p^1$  and  $\mathcal{V}_p^2 \subset V_p^2$ ) and maps  $\mathcal{N} \cap \mathcal{W}_p$  onto  $\mathcal{V}_p^1 \times \{0\}$ .

We shall say that  $\mathcal{N}$  is a submodel or a submanifold of the exponential statistical manifold  $\mathcal{M}(X, \mathcal{X}, \mu)$ .

A submanifold  $\mathcal{N}$  is a manifold whose charts are the restriction of the charts  $\sigma_p$  to  $\mathcal{N}$  (Figure 3).

Frequently below we shall use this particular splitting;  $V_p^1$  and  $V_p^2$  are closed subspaces of  $B_p$  such that  $V_p^1 \cap V_p^2 = \{0\}$  and  $B_p = V_p^1 + V_p^2$ . Then any element  $u \in B_p$  can be written uniquely as  $u = u_1 + u_2$ , with  $u_i \in V_p^i$ ,  $i = 1, 2$ .

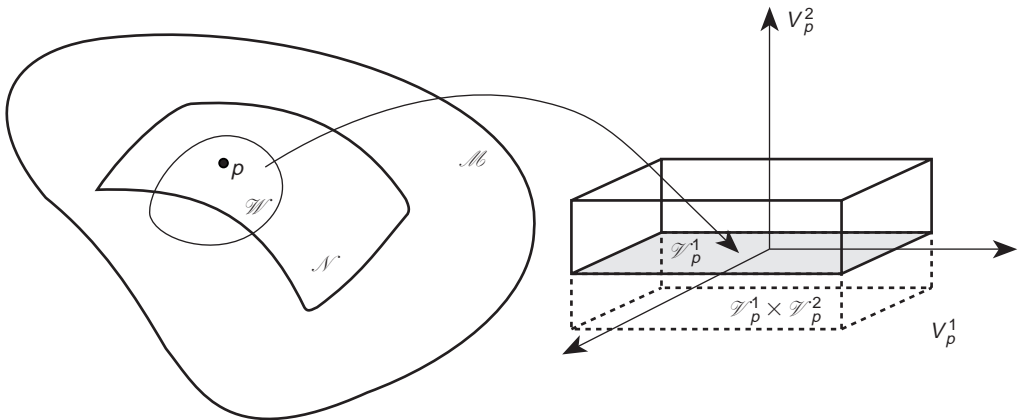
In this case  $\mathcal{N} = \{q = e^{u_1 - K_p(u_1)} : u_1 \in V_p^1 \cap \mathcal{V}_p^1\}$  is a submanifold, which we call “exponential submodel” of the maximal exponential model.

### 5.1. Some examples of submanifolds

We discuss two relevant examples of submanifolds.

#### 5.1.1. (m, d)-curved exponential models

We consider the parametric exponential model defined above in (12). Let  $u_1, \dots, u_d \in B_p$  and  $V_p^1 = \{\sum_{i=1}^d \theta_i u_i; \theta_i \in \mathbb{R}\}$  be given. Let  $P_{u_1, \dots, u_d}$  be the orthogonal projector from  $B_p$  (as a subspace of  $L_0^2(p)$ ) onto the finite-dimensional Hilbert space  $V_p^1$ . For each  $u \in B_p$ ,  $u = P_{u_1, \dots, u_d} u + (I - P_{u_1, \dots, u_d})u$ . The orthogonal projector  $P$  is characterized by



**Figure 3.** Submanifold: in this example the parallelepiped is  $\mathcal{V}_p^1 \times \mathcal{V}_p^2$  and the emphasized rectangle is  $\mathcal{V}_p^1$  (subset of horizontal plane).

$\langle u, u_j \rangle_p = \sum \hat{\theta}_i \langle u_i, u_j \rangle_p$  for  $j = 1, \dots, d$  and  $\hat{\theta} \in \Theta$ , with  $\Theta$  as in (11).  $V_p^1$  is closed in  $B_p$  because it is finite dimensional.  $P_{u_1, \dots, u_d}$  and  $I - P_{u_1, \dots, u_d}$  are continuous functions in  $B_p$ ; then  $V_p^2 = \ker(P_{u_1, \dots, u_d})$  is a closed subset of  $B_p$ . This implies the existence of a splitting.

Then the  $d$ -dimensional parametric exponential model

$$\mathcal{N} = \left\{ e^{\sum_i \theta_i u_i - \psi(\theta)} p \right\} \text{ with } \psi(\theta) = K_p \left( \sum_i \theta_i u_i \right)$$

is a submanifold of the nonparametric model  $\mathcal{M}(X, \mathcal{X}, \mu)$ . The set  $\Theta$ , as in (11), is an open set of  $\mathbb{R}^d$  and  $\theta \leftrightarrow u = \sum \theta_i u_i$  is a  $C^\infty$  chart of  $\mathcal{N}$  into  $\Theta$ . If  $A$  is an  $m$ -dimensional submanifold of  $\Theta$  the corresponding  $m$ -dimensional parametric submodel  $\mathcal{S}$  is a submanifold of  $\mathcal{N}$  and thus of  $\mathcal{M}(X, \mathcal{X}, \mu)$ .  $\mathcal{S}$  is the  $(m, d)$ -curved exponential model of Barndorff-Nielsen and Cox (1994, p. 65).

### 5.1.2. Conditional expectation

Let  $\mathcal{A}$  a sub- $\sigma$ -algebra of  $\mathcal{X}$ . We consider  $V_p^1 = L_0^{\phi_1}(X, \mathcal{A}, \mu)$  and the conditional expectation

$$E_p[\cdot | \mathcal{A}]: B_p \rightarrow V_p^1.$$

This mapping is well defined; in fact any element  $u \in B_p$  maps to an element of  $V_p^1$ :  $E_p[\phi_1(\alpha E_p[u | \mathcal{A}])] \leq E_p[E_p[\phi_1(\alpha u | \mathcal{A})]] = E_p[\phi_1(\alpha u)]$  (the Jensen inequality). It is surjective (any element of  $V_p^1$  maps to itself) and continuous (Neveu 1972). The subspace  $V_p^1$  is closed. We consider  $V_p^2 = \{u: E_p[u | \mathcal{A}] = 0\}$ . This subspace is closed because it is the kernel of a linear and continuous map.

$V_p^1$  and  $V_p^2$  split in  $B_p$  because  $u = E_p[u | \mathcal{A}] + (u - E_p[u | \mathcal{A}])$  is a unique decomposition. Then

$$\mathcal{M}(\mathcal{A}) = \{p \in \mathcal{M}(X, \mathcal{X}, \mu): p \text{ is } \mathcal{A} \text{ measurable}\}$$

is a submanifold. In particular this applies to invariance with respect to measurable transformation of the sample space.

## 6. Splitting and orthogonality

What follows is a nonparametric version of results taken from Amari (1982); we extend to the statistical exponential manifold the method of mixed parametrization for exponential models (Barndorff-Nielsen 1978a,b).

### 6.1. Splitting

Let  $V_p^1$  and  $V_p^2$  be two closed subspaces of  $B_p$ , such that  $V_p^1 \cap V_p^2 = \{0\}$  and  $B_p = V_p^1 + V_p^2$ , i.e. any element  $u \in B_p$  can be uniquely written as  $u = u_1 + u_2$ , with  $u_i \in V_p^i$ ,  $i = 1, 2$ . We consider the linear projectors  $P_i: B_p \rightarrow V_p^i$  defined by the splitting as  $P_i(u) = u_i$

with  $i = 1, 2$ . The projectors are both continuous because of the Banach closed-graph theorem (Lang 1995, p. 4).

The mixed parametrization in exponential models is based on analytical computations that involve partial derivatives. Our next step is to describe how to compute the partial derivatives when the two components are not Cartesian components, but the two projections induced by a splitting. Such an extension is straightforward; nevertheless it is useful to have precise notation for the nonparametric case.

We consider the *annihilating subspaces*  $(V_p^1)^0$  and  $(V_p^2)^0$  of  $B_p^*$ , and  ${}^0(V_p^1)$  and  ${}^0(V_p^2)$  of  ${}^*B_p$  defined by

$$(V_p^i)^0 = \{u^* \in B_p^*: \langle u^*, u_j \rangle_{*,p} = 0, \forall u_j \in V_p^j, j \neq i\},$$

$${}^0(V_p^i) = \{^*u \in {}^*B_p: E_p[^*uu_j] = 0, \forall u_j \in V_p^j, j \neq i\},$$

with  $i, j \in \{1, 2\}$ . Such spaces can be considered the “orthogonal” spaces to the spaces of the splitting. This means, for example, that, if  ${}^*u \in {}^0(V_p^2) \cap L_0^2(p)$ , then  ${}^*u$  belongs to  $(V_p^1)^\perp$  in the Hilbert sense. Note that the numbering is such that sup-1 is “orthogonal” to sup-2 and vice versa for consistency with the natural  $L^2$  notation.

The following proposition shows that the annihilating subspaces split the dual space and gives a characterization of the dual projections.

**Proposition 28.**  *$(V_p^1)^0$  and  $(V_p^2)^0$  split in  $B_p^*$  and any element of  $B_p^*$  can be written as  $u^* = u_1^* + u_2^*$ , with  $u_i^* \in (V_p^i)^0$ ,  $i = 1, 2$ . If  $v \in B_p$ ,  $v = v_1 + v_2$ ,  $v_i \in V_p^i$ , then  $u_i^*(v) = u^* \circ P_i(v)$ ; so  $u_i^* = u^* \circ P_i$ , for  $i = 1, 2$ .*

**Proof.** This follows from the following properties of the subspaces.

(a)  $(V_p^1)^0$  and  $(V_p^2)^0$  are closed in  $B_p^*$ ; in fact, if  $u_j$  is fixed in  $V_p^j$ , then the space  $\{u^* \in B_p^*: \langle u^*, u_j \rangle_{*,p} = 0\}$  is closed because it is the kernel of a continuous map, and  $(V_p^i)^0$ ,  $i \neq j$ , is the intersection of such closed subsets.

(b)  $(V_p^1)^0 \cap (V_p^2)^0 = \{0\}$ ; in fact, if  $u^* \in (V_p^1)^0 \cap (V_p^2)^0$ , then, for any  $v = v_1 + v_2 \in B_p$ ,  $\langle u^*, v_1 \rangle_{*,p} = 0$  and  $\langle u^*, v_2 \rangle_{*,p} = 0$ ; so  $u^* = 0$  because the bilinear form of the duality is separating.

(c)  $B_p^* = (V_p^1)^0 + (V_p^2)^0$ ; in fact, if  $u^* \in B_p^*$ , then  $u^* = u^* \circ P_1 + u^* \circ P_2$ ; then, if  $v \in B_p$ ,  $v = v_1 + v_2$ , we have  $u^*(v) = u^* \circ P_1(v) + u^* \circ P_2(v) = u^*(v_1) + u^*(v_2)$  and  $u_i^*$  (with  $u_i^* = u^*(P_i)$ )  $\in (V_p^i)^0$ . The last assertion follows from the definition of  $(V_p^i)^0$ ; in fact, if  $v_j \in V_p^j$  and  $i \neq j$ , then

$$\langle u_i^*, v_j \rangle_{*,p} = u_i^*(v_j) = u^* \circ P_i(v_j) = u^*(0) = 0. \quad \square$$

Note that properties (a) and (b) hold also for  ${}^0(V_p^1)$  and  ${}^0(V_p^2)$ .

### 6.1.1. Partial derivatives of $K_p$

Given a splitting of  $B_p^*$ , we now look for a characterization of  $u_1^*$  and  $u_2^*$  for those elements



$u^*$  of  $B_p^*$  (with  $u^* = DK_p(u)$ ) identified with an element  ${}^*u$  of  ${}^*B_p$  such that  ${}^*u = \nabla K_p(u)$  (see Proposition 8 (a) and Proposition 16 (c)).

If  ${}^*u = \nabla K_p(u) \in {}^*B_p$  is identified with  $u^* = DK_p(u) \in B_p^*$ ,  $v = v_1 + v_2 \in B_p$  and  $u^* = u_1^* + u_2^* \in B_p^*$ , then

$$u_i^*(v) = u^*(v_i) = DK_p(u)v_i = E_p[\nabla K_p(u)v_i], \quad i = 1, 2. \tag{16}$$

For  $i = 1, 2$  we define the two partial derivative of  $K_p(u)$  in the direction  $v$  as follows:

$$\partial_i K_p(u)v = \frac{\partial}{\partial t} K_p(u + P_i(tv))_{t=0} = DK_p(u)P_i(v) = DK_p(u)v_i = E_p[\nabla K_p(u)v_i]. \tag{17}$$

Note that  $\partial_i K_p(u)$  is an extension to  $B_p$  of the partial derivative of the function  $(u_1, u_2) \mapsto K_p(u_1, u_2) = K_p(u_1 + u_2)$  defined on  $V_p^1 \times V_p^2$ ; such a partial derivative takes the value 0 on  $V_p^j$ , with  $i \neq j$  and  $i, j \in \{1, 2\}$ .

By (16) and (17), we have that  $\partial_i K_p(u) = u_i^*$ . In general, this element of the dual space will not be in  ${}^*B_p$ , even if the gradient  ${}^*u = \nabla K_p(u)$  of  $K_p$  is; in the following we shall study special cases where such inclusion takes place; see Propositions 35 and 37 below.

The problem of finding general conditions which ensure that  $u_i^* \in B_p^*$  may be identified with an element of  ${}^*B_p$  remains open.

### 6.2. Orthogonality and mean parameters

In this section we show how to extend the notion of the Fisher information matrix and the corresponding notion of orthogonality (see, for example, Barndorff-Nielsen and Cox (1994)) to the exponential statistical manifold.

We assume that there exists a splitting  $V_p^1, V_p^2$  of  $B_p$ . Let  $q = e^{u-K_p(u)} p \in \mathcal{U}_p$  be given. From the remark immediately after the proof of Theorem 18 it is seen that the derivative of the transition mapping, i.e.

$$B_p \ni v \mapsto \tilde{v} = v - E_q[v] \in B_q$$

is a top-linear isomorphism (Lang 1995, Chapter 3). It follows also directly from Proposition 5. Here and in the following the  $\tilde{v}$  and the  $v$  differ only by a constant. The subspaces  $V_q^1$  and  $V_q^2$  defined as

$$V_q^i = \{\tilde{v}_i = v_i - E_q[v_i]: v_i \in V_p^i\}, \quad i = 1, 2,$$

split in  $B_q$ . In such a case, for all  $w_1, w_2 \in B_p$  and  $\tilde{w}_i = w_i - E_q[w_i]$ ,  $i = 1, 2$ , by Proposition 16 (b), we have

$$\begin{aligned} D^2 K_p(u)(w_1, w_2) &= \text{cov}_q[w_1, w_2] \\ &= E_q[(w_1 - E_q[w_1])(w_2 - E_q[w_2])] \\ &= \langle \tilde{w}_1, \tilde{w}_2 \rangle_q. \end{aligned} \tag{18}$$

As in the previous section the cumulant generating functional  $K_p$  can be considered as a function of two variables:

$$K_p(\cdot, \cdot): V_p^1 \times V_p^2 \ni (w_1, w_2) \mapsto K_p(w_1, w_2) = K_p(w_1 + w_2).$$

We can take partial derivatives that we shall denote by  $\partial_1, \partial_2, \partial_{11}, \partial_{12}, \partial_{21}, \partial_{22}$  by composing with the relevant projections; see (17).

**Definition 29 (Fisher information operator).** *The value at  $q = e_p(u)$  of the Hessian linear operator from  $B_p$  to  $B_p^*$  of  $K_p$  at  $u$  will be denoted by  $I(p, q)$ :*

$$\langle I(p, q)w, v \rangle_{*,p} = D^2 K_p(u)(w, v) \text{ with } w, v \in B_p.$$

and it is called the Fisher information operator.

Given a splitting of  $B_p$ , if  $v = v_1 + v_2$  and  $w = w_1 + w_2$ , we consider the partitioned operators  $I_{ij}$ , with  $i, j \in \{1, 2\}$ , restricted from  $V_p^j$  to  $(V_p^k)^\circ$ ,  $k \neq i$ , and such that

$$\langle I_{ij}(p, q)w_j, v_i \rangle_{*,p} = \langle I(p, q)w_j, v_i \rangle_{*,p}.$$

In matrix form,

$$I(p, q) = \begin{bmatrix} I_{11}(p, q) & I_{12}(p, q) \\ I_{21}(p, q) & I_{22}(p, q) \end{bmatrix}: B_p \simeq V_p^1 \times V_p^2 \rightarrow (V_p^1)^\circ \times (V_p^2)^\circ \simeq B_p^*.$$

From (18) we get, for  $\tilde{w} = w - E_q[w]$ ,  $\tilde{v} = v - E_q[v]$ ,

$$\text{cov}_q[w, v] = \langle \tilde{w}, \tilde{v} \rangle_q = \langle I(p, q)w, v \rangle_{*,p}.$$

We now consider the problem of finding an orthogonal projection of  $B_q$  onto  $V_q^1$ .

The space  $B_q$  is a pre-Hilbert space for the scalar product  $\langle \cdot, \cdot \rangle_q$  as in Definition 10, i.e. the space is in general not complete under the norm induced by the scalar product. The existence of the orthogonal projection of  $w \in B_q$  onto  $V_q^1$  is not ensured but, if the projection  $w_{|1}$  exists, then it is unique; moreover, given a splitting such that  $w = w_1 + w_2$ , then  $w_{2|1}$  defined as

$$w_{2|1} = w_2 - w_1$$

is the orthogonal projection of  $w_2$  onto  $V_q^1$ .

The following proposition shows how to characterize  $w_{2|1}$  with the partitioned Fisher information operator.

**Proposition 30.** *If  $w = w_1 + w_2$  and  $\tilde{w}_i = w_i - E_q[w_i]$ , for  $i = 1, 2$ , then  $\tilde{w} = \tilde{w}_1 + \tilde{w}_2$ . If there exists  $w_{2|1} \in V_p^1$  such that the normal equation*

$$I_{11}(p, q)w_{2|1} = I_{12}(p, q)w_2 \tag{19}$$

is satisfied, then  $\tilde{w}_{2|1} = w_{2|1} - E_q[w_{2|1}]$  is the orthogonal projection of  $\tilde{w}_2 \in V_q^2$  onto  $V_q^1$ .

**Proof.** If (19) holds true, then for any  $w_1 \in V_p^1$  we have, from (18) and using the Fisher information operator,

$$\begin{aligned} \langle \tilde{w}_2 - \tilde{w}_{2|1}, \tilde{w}_1 \rangle_q &= \langle \tilde{w}_2, \tilde{w}_1 \rangle_q - \langle \tilde{w}_{2|1}, \tilde{w}_1 \rangle_q \\ &= \langle I_{12}(p, q)w_2, w_1 \rangle_{*,p} - \langle I_{11}(p, q)w_{2|1}, w_1 \rangle_{*,p} \\ &= 0. \end{aligned}$$

This ends the proof. □

**Definition 31.** Let  $B_1, B_2$  be Banach spaces, and let  $F$  be a  $C^k$  parametrization of  $A$ , an open subset of  $\mathcal{M}(X, \mathcal{X}, \mu)$ , on  $B_1 \times B_2$ :

$$A \ni q \mapsto F(q) = (F_1(q), F_2(q)) \in B_1 \times B_2.$$

We consider the pair of regular curves

$$]a, b[ \ni t \mapsto q_i(t) \in A, \quad i = 1, 2,$$

such that  $0 \in ]a, b[$  and

$$q_1(0) = q_2(0) = q \quad \text{and} \quad F_i(q_j(t)) = \text{constant}, \quad i \neq j, t \in ]a, b[. \quad (20)$$

We shall say that the two parameters  $F_1$  and  $F_2$  are orthogonal at  $q \in A$  if each such pair of curves is orthogonal at  $q$ , i.e.  $\dot{q}_1(0)$  is orthogonal to  $\dot{q}_2(0)$ . If this is true for all  $q \in A$ , then we shall say that the two parameters are orthogonal.

**Theorem 32 (Mixed parametrization).** Given a splitting  $V_p^1, V_p^2$  of  $B_p$ , the mapping

$$F: \mathcal{U}_p \ni q \mapsto (\eta_1, u_2) \in (V_p^1)^\circ \times V_p^2$$

with  $q = e^{u-K_p(u)} p$ ,  $u = u_1 + u_2$ ,  $u_i \in V_p^i$ ,  $i = 1, 2$ , and

$$\eta_1 = \partial_1 K_p(u)$$

is an orthogonal  $C^\infty$  parametrization of  $\mathcal{U}_p$ .

Note that, from (17) and (7),

$$\eta_1(v) = \langle \eta_1, v \rangle_{*,p} = DK_p(u)P_1(v) = DK_p(u)v_1 = E_q[v_1];$$

so the present parametrization will be called mixed parametrization, as in the parametric case (Barndorff-Nielsen and Cox 1989; 1994, p. 62).

**Proof.** First we show that the mapping  $q = e^{u-K_p(u)} p \mapsto (\partial_1 K_p(u), u_2)$  is one to one. In (9), if  $\bar{u} = \bar{u}_1 + \bar{u}_2$  and  $\bar{u}_2 = u_2$ , then  $\bar{u} - u$  equals  $\bar{u}_1 - u_1$  and, denoting  $e^{\theta u - K_p(\theta u)} p$  by  $p_\theta$ , we have

$$\begin{aligned} \int_0^1 d\theta \operatorname{var}_{p_\theta}(\bar{u} - u) &= \{DK_p(\bar{u}) - DK_p(u)\}(\bar{u}_1 - u_1) \\ &= \{\partial_1 K_p(\bar{u}) - \partial_1 K_p(u)\}(\bar{u} - u) \\ &= \langle \bar{\eta}_1 - \eta_1, \bar{u} - u \rangle_{*,p}. \end{aligned}$$

If  $\partial_1 K_p(\bar{u}) = \partial_1 K_p(u)$  and  $\bar{u}_2 = u_2$ , then  $\bar{u} - u = \text{constant}$  and  $\bar{u}_1 = u_1$ , because  $\bar{u}$  and  $u$  are centred random variables.

Let  $u^{(i)}$ ,  $i = 1, 2$ , be the representation in the chart  $s_p$  of the regular curves through  $q$  in Definition 31, corresponding to the mixed parameter  $(\eta_1, u_2)$ :

$$q_i(t) = e^{u^{(i)}(t) - K_p(u^{(i)}(t))} p, \quad i = 1, 2.$$

Because of the second condition of (20) the first curve leaves the second parameter constant,  $u_2^{(1)}(t) = u_2$ , and the second curve leaves the first parameter constant, i.e., for any  $w \in B_p$ ,  $\partial_1 K_p(u^{(2)}(t))w = \text{constant}$ . If we take the derivative of the last equation with respect to  $t$  we get

$$D^2 K_p(u^{(2)}(t))(P_1(w), \dot{u}^{(2)}(t)) = 0. \tag{21}$$

The condition on the first curve can be written as  $P_2(u^{(1)}(t)) = u_2$  and derivation gives  $P_2(\dot{u}^{(1)}(t)) = 0$ . Consequently,  $\dot{u}^{(1)}(t) = P_1(\dot{u}^{(1)}(t))$  and substituting  $w$  in (21) with  $\dot{u}^{(1)}(t)$  we get

$$D^2 K_p(u^{(2)}(t))(\dot{u}^{(1)}(t), \dot{u}^{(2)}(t)) = 0.$$

We denote by  $\tilde{u}^{(i)}(t)$  the coordinate of  $q_i(t)$  with respect to the chart  $s_q$ :

$$\tilde{u}^{(i)}(t) = u^{(i)}(t) - u^{(i)}(0) - E_q[u^{(i)}(t) - u^{(i)}(0)];$$

then  $\tilde{\dot{u}}^{(i)}(t) = \dot{u}^{(i)}(t) - E_q[\dot{u}^{(i)}(t)]$ . If  $t = 0$ , it follows from (17) and Definition 29 that

$$\begin{aligned} \langle \tilde{\dot{u}}^{(1)}(0), \tilde{\dot{u}}^{(2)}(0) \rangle_q &= \langle I(p, q)\dot{u}^{(1)}(0), \dot{u}^{(2)}(0) \rangle_{*,p} \\ &= D^2 K_p(u^{(2)}(0))(\dot{u}^{(1)}(0), \dot{u}^{(2)}(0)) \\ &= 0. \end{aligned}$$

So  $q_1(t)$  and  $q_2(t)$  are orthogonal at  $q$ . □

## 7. Transformation of the sample space

The problem that we consider in this section is the action of a transformation of the sample space on the manifold structure. Let us first introduce some notation. Let a measurable mapping on the sample space be given:

$$\varphi: (X, \mathcal{X}, \mu) \rightarrow (Y, \mathcal{Y}, \nu), \tag{22}$$

i.e.

$$\varphi: X \rightarrow Y, \quad \varphi^{-1}: \mathcal{Y} \rightarrow \mathcal{X}, \quad \nu = \mu \circ \varphi^{-1}.$$

If  $p \in \mathcal{M}(X, \mathcal{X}, \mu)$ , then the image under  $\varphi$  of the probability measure  $p \cdot \mu$  is characterized by

$$(p \cdot \mu) \circ \varphi^{-1}(B) = \int_{\varphi^{-1}(B)} p \, d\mu = \int_{\varphi^{-1}(B)} E_{\mu}[p|\varphi^{-1}(\mathcal{Y})] \, d\mu, \quad B \in \mathcal{Y}.$$

By the Doob lemma,

$$E_{\mu}[p|\varphi^{-1}(\mathcal{Y})] = \hat{p} \circ \varphi, \tag{23}$$

where  $\hat{p}$  is defined up to sets of  $\nu$ -measure 0. We shall write  $\hat{p} = E_{\nu}[p|\varphi]$ , so that

$$(p \cdot \mu) \circ \varphi^{-1}(B) = \int_{\varphi^{-1}(B)} \hat{p} \circ \varphi \, d\mu = \int_B E_{\nu}[p|\varphi] \, d\nu, \quad B \in \mathcal{Y},$$

and the sample space transformation  $\varphi$  induces a transformation  $\hat{\varphi}$  on  $\mathcal{M}(X, \mathcal{X}, \mu)$  given by

$$\hat{\varphi}: \mathcal{M}(X, \mathcal{X}, \mu) \ni p \mapsto \hat{p} = E_{\nu}[p|\varphi] \in \mathcal{M}(Y, \mathcal{Y}, \nu).$$

### 7.1. Transformation

The basic result about the kind of smoothness of the mapping  $\hat{\varphi}: p \mapsto \hat{p}$  that we are able to prove is shown by the following Proposition 33. We refer to the book by Lang (1995) for an introduction to the basic notions regarding the differentiability and regularity of mappings between differentiable manifolds. In particular we shall mention the notion of *submersion*. This notion is connected with the stability of the rank of a mapping. Note that there is misuse of the language because the derivative mapping of the transformation is a linear continuous mapping between the tangent spaces, but we are able to prove the continuous differentiability in a weaker sense only.

The given regularity results are not connected in any way with the regularity of the mapping  $\varphi$  itself, nor with any algebraic or topological property of the spaces  $X$  or  $Y$ . The possible relations between the differentiable structure of the exponential statistical manifold and the regularity of the densities themselves are not dealt with at all in the present paper and suggest different research work. A special *ad hoc* condition is introduced in (27) to prove the Gateaux regularity of the mapping between the exponential statistical manifolds. We do not discuss examples of application here. We just mention the fact that this condition is easily shown to be true in trivial cases, such as constant or injective transformation, or  $Y$  with a finite number of atoms.

Differentiability of higher order could be proved under the analogous conditions. We shall not discuss this topic further as we are not going to use it.

**Proposition 33 (Sample space transformation).** *Let us consider a space transformation  $\varphi$  as in (22) and let*

$$\hat{\varphi}: \mathcal{M}(X, \mathcal{X}, \mu) \ni p \mapsto \hat{\varphi}(p) = \hat{p} \in \mathcal{M}(Y, \mathcal{Y}, \nu)$$

denote the action of  $\varphi$  on the exponential statistical manifolds. Then we have the following.

(a)  $\hat{\varphi}$  is onto and it is injective if and only if the mapping  $\varphi$  is  $\mu$ -almost surely injective, i.e. it generates the  $\sigma$ -algebra  $\mathcal{X}$  up to  $\mu$ -null sets.

(b) The coordinate form

$$\hat{\varphi}_p = s_{\hat{p}} \circ \hat{\varphi} \circ e_p, \quad \hat{p} = \hat{\varphi}(p),$$

of the mapping  $\hat{\varphi}$  around the point  $p$  is given by

$$\hat{\varphi}_p(u) = K_{\hat{p}}(u|\varphi) - E_{\hat{p}}[K_{\hat{p}}(u|\varphi)], \tag{24}$$

where

$$K_{\hat{p}}(u|\varphi) = \log E_{\hat{p}}[e^u|\varphi] \tag{25}$$

in a neighbourhood of 0. Such a mapping is continuous from an open neighbourhood of 0 of  $B_p$  to  $B_{\hat{p}}$ ; so the mapping  $\hat{\varphi}$  is of class  $C^0$  from the manifold  $\mathcal{M}(X, \mathcal{X}, \mu)$  to the manifold  $\mathcal{M}(Y, \mathcal{Y}, \nu)$ . Moreover for all  $\alpha \geq 1$  the mapping  $\hat{\varphi}_p$  is continuously differentiable as a mapping in  $L^\alpha(Y, \mathcal{Y}, \hat{p} \cdot \nu)$  and its derivative at  $u \in \mathcal{T}_p$  in the direction  $v \in B_p$  is given by

$$D\hat{\varphi}_p(u)v = E_{\hat{q}}[v|\varphi] - E_{\hat{p}}[E_{\hat{q}}[v|\varphi]], \quad q = e_p(u). \tag{26}$$

(c) The derivative mapping in (26) is a linear continuous operator from  $B_p$  to  $B_{\hat{p}}$ . If moreover  $u \in B_p$  is such that, for all  $\alpha \geq 1$ ,

$$E_{\hat{p}}[e^{u - E_p[u|\varphi]}|\varphi] \in L^\alpha(Y, \mathcal{Y}, \hat{p} \cdot \nu), \tag{27}$$

then the mapping  $\hat{\varphi}_p$  is Gateaux differentiable from  $B_p$  to  $B_{\hat{p}}$  in the direction  $u$ .

(d)  $u \mapsto D\hat{\varphi}_p(u)$  is surjective and its kernel splits.

**Proof.**

(a) The surjectivity follows from the condition formula (23); in fact for each density  $\tilde{p} \in \mathcal{M}(Y, \mathcal{Y}, \nu)$ , the density  $\tilde{p} \circ \varphi$  in  $\mathcal{M}(X, \mathcal{X}, \mu)$  and it is such that  $\hat{\varphi}(\tilde{p} \circ \varphi) = \tilde{p}$ . If  $\varphi^{-1}(\mathcal{Y}) = \mathcal{X}$   $\mu$ -almost surely, then using (23)

$$\hat{p} \circ \varphi = E_\mu[p|\varphi^{-1}(\mathcal{Y})] = p.$$

If  $\hat{\varphi}(p_1) = \hat{\varphi}(p_2) = \hat{p}$ , then  $p_1 = p_2 = \hat{p} \circ \varphi$ . On the other hand, if  $\varphi^{-1}(\mathcal{Y}) \subset \mathcal{X}$  is a proper sub- $\sigma$ -algebra, then there exists a density  $p$  which is not  $\varphi^{-1}(\mathcal{Y})$  measurable, so that  $p$  and  $\hat{p} \circ \varphi$  are different and have the same image  $\hat{p}$ .

(b) Let us compute the coordinate form of the mapping  $\hat{\varphi}$ . Consider a density  $p$ , the chart domain  $\mathcal{U}_p$ , and let  $q$  be a density representable in such a chart,  $q = e^{u - K_p(u)} p$ , where  $\|u\|_p < 1$ . Then  $\hat{\varphi}(q) = \hat{q}$  is given by

$$\hat{q} \circ \varphi = E_{\mu}[e^{u-K_{\rho}(u)} p | \varphi^{-1}(\mathcal{Y})].$$

Now we compute the likelihood with respect to  $\hat{p}$ :

$$\begin{aligned} \left(\frac{\hat{q}}{\hat{p}}\right) \circ \varphi &= \frac{\hat{q} \circ \varphi}{\hat{p} \circ \varphi} \\ &= \frac{E_{\mu}[e^{u-K_{\rho}(u)} p | \varphi^{-1}(\mathcal{Y})]}{E_{\mu}[p | \varphi^{-1}(\mathcal{Y})]} \\ &= E_{\rho}[e^{u-K_{\rho}(u)} | \varphi^{-1}(\mathcal{Y})] \\ &= e^{-K_{\rho}(u)} E_{\rho}[e^u | \varphi^{-1}(\mathcal{Y})]. \end{aligned}$$

The previous formula simplifies to

$$\frac{\hat{q}}{\hat{p}} = e^{-K_{\rho}(u)} E_{\hat{p}}[e^u | \varphi],$$

and we can write

$$\hat{q} = \hat{\varphi}(q) = \hat{\varphi} \circ e_{\rho}(u) = e^{-K_{\rho}(u)} E_{\hat{p}}[e^u | \varphi] \hat{p}.$$

We have to find the coordinates of  $\hat{q}$  with respect to the chart on  $\mathcal{U}_{\hat{p}}$ . If  $\hat{q} \in \mathcal{U}_{\hat{p}}$ , then

$$\hat{u} = s_{\hat{p}}(\hat{q}) = \log\left(\frac{\hat{q}}{\hat{p}}\right) - E_{\hat{p}}\left[\log\left(\frac{\hat{q}}{\hat{p}}\right)\right] = \log E_{\hat{p}}[e^u | \varphi] - E_{\hat{p}}[\log E_{\hat{p}}[e^u | \varphi]]. \tag{28}$$

Using the notation in (25), we get (24).

We shall show that  $\hat{u}$  in (28) actually belongs to  $\mathcal{U}_{\hat{p}}$  if  $u$  belongs to a suitable neighbourhood of zero. We start by showing that the mapping

$$\mathcal{V}_{\rho} \ni u \mapsto K_{\hat{p}}(u | \varphi) \in L^{\phi_1}(\hat{p} \cdot \nu) \tag{29}$$

is defined and norm decreasing. In fact, using the definition of norm in (3), we take  $\alpha$  such that  $\alpha \|u\|_{\hat{p}} < 1$  and  $\alpha = 1/r > 1$ . From the convexity of the function

$$f(x) = \frac{1}{2}(x^{\alpha} + x^{-\alpha})$$

for  $\alpha > 1$ , and the Jensen inequality, it follows that the norm

$$\|K_{\hat{p}}(u | \varphi)\|_{\hat{p}} = \|\log E_{\hat{p}}[e^u | \varphi]\|_{\hat{p}}$$

can be bounded above. In fact,

$$\begin{aligned} \cosh(\alpha \log E_{\hat{p}}[e^u | \varphi]) - 1 &= \frac{1}{2}(E_{\hat{p}}[e^u | \varphi]^{\alpha} + E_{\hat{p}}[e^u | \varphi]^{-\alpha}) - 1 \\ &\leq E_{\hat{p}}[\cosh(\alpha u) - 1 | \varphi] \end{aligned}$$

and, taking the expected values, as  $1 > 1/\alpha = r > \|u\|_{\hat{p}}$ ,

$$E_{\hat{p}}[\cosh(\alpha \log E_{\hat{p}}[e^u | \varphi]) - 1] \leq E_{\hat{p}}[\cosh(\alpha u) - 1] \leq 1. \tag{30}$$

As  $\|u\|_p < 1$  and from (30), it follows that

$$\|K_{\hat{p}}(u|\varphi)\|_{\hat{p}} \leq \|u\|_p.$$

As the expectation is a contraction operator on the space of  $L^{\phi_1}(\hat{p} \cdot \nu)$ , then  $\hat{u}$  in (28) belongs to the domain  $\mathcal{U}_{\hat{p}}$  of the chart at  $\hat{p}$  for  $u$  in a suitable neighbourhood of  $p$ .

To prove the continuity of  $\hat{\varphi}_p$ , consider the difference of the values of the function  $K_{\hat{p}}(\cdot|\varphi)$  defined in (25) at two points  $u$  and  $v$  of the domain of  $\hat{\varphi}_p$ :

$$\begin{aligned} K_{\hat{p}}(v|\varphi) - K_{\hat{p}}(u|\varphi) &= \log E_{\hat{p}}[e^v|\varphi] - \log E_{\hat{p}}[e^u|\varphi] \\ &= \log \left( \frac{E_{\hat{p}}[e^v|\varphi]}{E_{\hat{p}}[e^u|\varphi]} \right) \\ &= \log E_{\hat{q}}[e^{v-u}|\varphi], \end{aligned}$$

where  $q = e_p(u)$ . The norm-decreasing property shown above, together with the equivalence of  $B_{\hat{q}}$  and  $B_{\hat{p} \cdot \nu}$  (see Proposition 5), now implies the continuity of  $K_{\hat{p}}(\cdot|\varphi)$  and in turn the continuity of  $\hat{\varphi}_p$ . Finally  $\hat{\varphi}$  is continuous because its chart representation is continuous.

For the differentiability the same argument as above shows that it is enough to show the differentiability of the mapping  $u \mapsto K_{\hat{p}}(u|\varphi)$  defined in (25).

Let us check first the directional derivative. For  $u, v \in B_p, u \in \mathcal{Z}'_p, q = e^{u-K_p(u)} \mu, \hat{q} = \hat{\varphi}(q)$ , and  $t \rightarrow 0$ , we want  $R(u, v, t)$  to go to 0, where

$$\begin{aligned} R(u, v, t) &= t^{-1}(K_{\hat{p}}(u + tv|\varphi) - K_{\hat{p}}(u|\varphi)) - E_{\hat{q}}[v|\varphi] \\ &= t^{-1}K_{\hat{q}}(tu|\varphi) - E_{\hat{q}}[v|\varphi] \\ &= t^{-1}K_{\hat{q}}(t(v - E_{\hat{q}}[v|\varphi])|\varphi). \end{aligned}$$

From the concavity of the log function we get

$$\begin{aligned} R(u, v, t) &= t^{-1} \log E_{\hat{q}}[e^{t(v - E_{\hat{q}}[v|\varphi])}|\varphi] \\ &\geq t^{-1} E_{\hat{q}}[\log e^{t(v - E_{\hat{q}}[v|\varphi])}|\varphi] \\ &= E_{\hat{q}}[v - E_{\hat{q}}[v|\varphi]|\varphi] \\ &= 0. \end{aligned}$$

We use now a classical argument from Hardy *et al.* (1952). When  $t$  decreases to 0, the function

$$t \mapsto E_{\hat{q}}[e^{t(v - E_{\hat{q}}[v|\varphi])}|\varphi]^{1/t}$$

is decreasing and the same is true for  $t \mapsto R(u, v, t)$ . From the inequality  $\log x \leq x - 1$  we get



$$\begin{aligned}
 R(u, v, t) &\leq t^{-1}(\mathbb{E}_{\hat{q}}[e^{(v - \mathbb{E}_{\hat{q}}[v|\varphi])|\varphi}] - 1) \\
 &= \mathbb{E}_{\hat{q}} \left[ \left( \frac{e^{t(v - \mathbb{E}_{\hat{q}}[v|\varphi])} - 1}{t} \right) \middle| \varphi \right].
 \end{aligned}
 \tag{31}$$

The right-hand side of (31) goes to 0 as  $t$  goes to 0, and it is eventually bounded in all  $L^\alpha$ . The same is true for negative  $t$ , as can be shown by changing  $v$  to  $-v$ . This and the addition of the expected value show directional (Gateaux) differentiability at any point of the domain of  $\hat{\varphi}_p$  and the form of the derivative.

Finally we consider the continuity of the derivative. We consider the difference

$$\mathbb{E}_{\hat{q}}[v|\varphi] - \mathbb{E}_{\hat{p}}[v|\varphi] = \mathbb{E}_{\hat{p}}[v(e^{u - K_{\hat{p}}(u|\varphi)} - 1)|\varphi].$$

The  $L^\alpha$  norm of this difference is bounded by the  $L^\alpha$  norm of

$$v(e^{u - K_p(u|\varphi^{-1}(\mathcal{Z}))} - 1),$$

where  $K_p(u|\varphi^{-1}(\mathcal{Z})) = \log \mathbb{E}_p[e^u|\varphi^{-1}(\mathcal{Z})]$ . Now

$$\begin{aligned}
 \|v e^{u - K_p(u|\varphi^{-1}(\mathcal{Z}))}\|_{L^\alpha(p)}^\alpha &= \mathbb{E}_p[|v|^\alpha e^{\alpha(u - K_p(u|\varphi^{-1}(\mathcal{Z})))}] \\
 &\leq \mathbb{E}_p[|v|^{2\alpha}]^{1/2} \mathbb{E}_p[e^{2\alpha(u - K_p(u|\varphi^{-1}(\mathcal{Z})))}]^{1/2}.
 \end{aligned}
 \tag{32}$$

However,  $K_p(\cdot|\varphi^{-1}(\mathcal{Z}))$  is a contraction; then on the open set of  $\mathcal{Z}_p$  where  $2\alpha\|u - K_p(u|\varphi^{-1}(\mathcal{Z}))\|_p < 1$  the second factor of (32) is bounded by 2. The first factor is bounded by a constant times  $\|v\|_p^\alpha$ .

(c) The linear mapping

$$v \mapsto \mathbb{E}_{\hat{q}}[v|\varphi]$$

is a contraction from  $L^{\phi_1}(q)$  to  $L^{\phi_1}(\hat{q})$ , because it is a conditional expectation. Then it is continuous from  $B_p$  to  $B_{\hat{p}}$  because of the equivalence of the  $B$  spaces.

Assume now the condition (27). Then, for all  $\alpha > 1$ ,

$$\cosh(\alpha R(u, v, t) - 1) = \frac{1}{2}(\mathbb{E}_{\hat{q}}[e^{t(v - \mathbb{E}_{\hat{q}}[v|\varphi])|\varphi}]^{\alpha/t} + \mathbb{E}_{\hat{q}}[e^{t(v - \mathbb{E}_{\hat{q}}[v|\varphi])|\varphi}]^{-\alpha/t}) - 1.$$

As  $t$  decreases to 0, this quantity has been shown to decrease to zero. If the expected value at  $t = 1$  is finite, then the bounded convergence theorem implies the convergence to zero of expected values. This is precisely what is implied by the condition (27).

(d) We have shown that the derivative of the local coordinate form of the transformation  $\hat{\varphi}_p$  at  $u \in \mathcal{Z}_p$  in the direction  $v \in B_p$  is given by (26). This linear mapping is surjective; in fact, for each  $\hat{v} \in B_{\hat{p}}$  we have  $v = \hat{v} \circ \varphi \in B_p$  and

$$\begin{aligned}
 D\hat{\phi}_p(u)v &= D\hat{\phi}_p(u)v \\
 &= E_{\hat{q}\cdot\nu}[\hat{v} \circ \phi|\phi] - E_{\hat{p}}[E_{\hat{q}\cdot\nu}[\hat{v} \circ \phi|\phi]] \\
 &= \hat{v} - E_p[v] \\
 &= \hat{v}.
 \end{aligned}$$

Let us now study the kernel of  $D\hat{\phi}_p(u)$ . This linear mapping can be written as  $D\hat{\phi}_p(u) = C \circ E(u)$ , where

$$E = E_{\hat{q}}[\cdot|\phi]: B_p \rightarrow L^{\phi_1}(\hat{p} \cdot \nu)$$

and

$$C: L^{\phi_1}(\hat{p} \cdot \nu) \ni \hat{w} \mapsto \hat{w} - E_{\hat{p}}[\hat{w}].$$

The kernel of  $C$  is the set of constant random variables in  $L^{\phi_1}(\hat{p} \cdot \nu)$ . The constant elements  $v \in B_p$ , such that this constant is 0, and we have shown that

$$\ker D\hat{\phi}_p(u) = \{v \in B_p: E_p[v|\phi^{-1}(\mathcal{Y})] = 0\}.$$

As we have shown in the previous paragraphs (see the second example of submanifold in Section 5.1), the kernel of the conditional expectation splits. □

**Remark 34.** If the mapping  $\hat{\phi}$  is of class  $C^n$  from  $\mathcal{M}(X, \mathcal{X}, \mu)$  to  $\mathcal{M}(Y, \mathcal{Y}, \nu)$ ,  $n \geq 1$ , then the previous splitting property implies that this mapping is a *submersion*, and, for all  $\hat{p} \in \mathcal{M}(Y, \mathcal{Y}, \nu)$ ,  $\hat{\phi}^{-1}(\hat{q})$  is a submanifold of  $\mathcal{M}(X, \mathcal{X}, \mu)$  (Lang 1995, p. 25). This general result cannot be applied directly, because we have not proved in general the  $C^1$  regularity. Nevertheless it can be applied to special submanifolds. If we assume that  $\mathcal{N}$  is a submanifold of  $\mathcal{M}(X, \mathcal{X}, \mu)$  such that the restriction of  $\hat{\phi}$  to  $\mathcal{N}$  is  $C^1$ , and the splitting property holds on the submanifold, then the splitting result implies that  $\hat{\phi}^{-1}(\hat{q}) \cap \mathcal{N}$  is a submanifold of  $\mathcal{M}(X, \mathcal{X}, \mu)$ .

The following sections show a different direct approach to the differential structure of  $\hat{\phi}^{-1}(\hat{q})$ , based on the idea of regular parametrization.

### 7.2. Bivariate densities with one marginal given

Let us assume in the remaining part of this section that the sample space  $(X, \mathcal{X}, \mu)$  is a product space  $X = Y \times Z$ ,  $\mathcal{X} = \mathcal{Y} \otimes \mathcal{Z}$ ,  $\mu = \nu \otimes \pi$ . Let us denote by  $\phi_1: Y \times Z \rightarrow Y$  and  $\phi_2: Y \times Z \rightarrow Z$  the two marginal mappings.

**Proposition 35 (One marginal given).**

(a) Let  $p_1 = \hat{\phi}_1(p)$  be given. For each  $q \in \mathcal{U}_p$ ,  $q = e^{u-K_p(u)} p$ , let  $q_1 = \hat{\phi}_1(q)$ . Then we have the following.

(i) The mapping  $\mathcal{U}_p \ni q \mapsto q_1/p_1 - 1 \in {}^*B_{p_1}$  is  $C^\infty$  from  $\mathcal{M}(X, \mathcal{X}, \mu)$  to  $B_{p_1}^*$ .

(ii) The mapping

$$\mathcal{U}_p \ni q \mapsto \left( \frac{q_1}{p_1} - 1, u - E_p[u|\varphi_1^{-1}(\mathcal{Y})] \right) \in {}^* B_{p_1} \times \ker E_p[\cdot|\varphi_1^{-1}(\mathcal{Y})]$$

is an orthogonal parametrization of  $\mathcal{U}_p$  on  $B_{p_1}^* \times \ker E_p[\cdot|\varphi_1^{-1}(\mathcal{Y})]$ .

(b) The set of all densities whose first marginal is  $p_1$ , i.e.

$$\mathcal{N}_{p_1} = \{p \in \mathcal{M}(X, \mathcal{X}, \mu): \hat{\varphi}_1(p) = p_1\},$$

has a  $C^\infty$  parametrization on  $\mathcal{U}_p$  given by

$$\mathcal{N}_{p_1} \ni p \mapsto u - E_p[u|\varphi_1^{-1}(\mathcal{Y})] \in \ker E_p[\cdot|\varphi_1^{-1}(\mathcal{Y})].$$

**Proof.**

(a) It is a particular case of Proposition 32 based on the splitting

$$B_p \ni u \leftrightarrow E_p[u|\varphi_1^{-1}(\mathcal{Y})] + (u - E_p[u|\varphi_1^{-1}(\mathcal{Y})]).$$

To identify the first component  $\eta_1$  of the orthogonal parametrization we compute, as in Proposition 32, the partial derivative with respect to the first component of the splitting. We have

$$\begin{aligned} \langle \eta_1, v \rangle_{*,p} &= \partial_1 K_p(u)v \\ &= DK_p(u) P_1(v) \\ &= DK_p(u) E_p[v|\varphi_1^{-1}(\mathcal{Y})] \\ &= E_q[E_p[v|\varphi_1^{-1}(\mathcal{Y})]] \\ &= E_q[E_{p_1}[v|\varphi_1^{-1}(\mathcal{Y})] \circ \varphi_1] \\ &= E_{q_1}[E_{p_1}[v|\varphi_1]] \\ &= E_{p_1} \left[ \frac{q_1}{p_1} E_{p_1}[v|\varphi_1] \right] \\ &= E_p \left[ \frac{q_1}{p_1} \circ \varphi_1 E_p[v|\varphi_1^{-1}(\mathcal{Y})] \right] \\ &= E_p \left[ \frac{q_1}{p_1} \circ \varphi_1 v \right] \\ &= E_p \left[ \left( \frac{q_1}{p_1} - 1 \right) \circ \varphi_1 v \right]. \end{aligned}$$

We have characterized  $\eta_1$  as  $(q_1/p_1 - 1) \circ \varphi_1$ .

(b)  $q \in \mathcal{N}_{p_1} \cap \mathcal{U}_p$  if and only if  $q_1/p_1 - 1 = 0$ . □

Note that there is a one-to-one mapping between the bivariate densities with a given

marginal and the set of conditional densities, i.e. the Bayes formula. This induces a parametrization on the set of the corresponding conditional densities  $\{p/p_1: p \in \mathcal{N}_{p_1}\}$ .

### 7.3. Bivariate densities with two given marginals

**Definition 36 (Fréchet class).** Using the previous notation, the set of all densities whose marginals are given, i.e.

$$\mathcal{N}_{p_1 p_2} = \{p \in \mathcal{M}(X, \mathcal{X}, \mu): \hat{\varphi}_i(p) = p_i, i = 1, 2\} \tag{33}$$

is called the Fréchet class with marginals  $p_1, p_2$ .

The following proposition shows that each Fréchet class has a regular parametrization and gives some details of its structure.

**Proposition 37 (Fréchet manifold).** Let two marginal densities  $p_1, p_2$  be given: we consider the Fréchet class  $\mathcal{N}_{p_1 p_2}$  as in (33).

(a) Let  $p = p_1 p_2$  be the product density, and let  $u \in B_p$  be given. The formulae

$$u_1 = \int u p_2 \, d\pi,$$

$$u_2 = \int u p_1 \, d\nu,$$

$$u_{12} = u - u_1 - u_2,$$

define a splitting of  $B_p$  into the spaces  $(B_{p_1} \times B_{p_2})$  and  $B_p^{12}$ , where

$$B_p^{12} = \left\{ w \in B_p: \int w p_2 \, d\pi = \int w p_1 \, d\nu = 0 \right\}.$$

(b) The mapping

$$\mathcal{U}_p \ni q \mapsto \left( \left( \frac{q_1}{p_1} - 1, \frac{q_2}{p_2} - 1 \right), u_{12} \right) \in (*B_{p_1} \times *B_{p_2}) \times B_p^{12} \tag{34}$$

is an orthogonal parametrization of  $\mathcal{U}_p$  on  $(B_{p_1}^* \times B_{p_2}^*) \times B_p^{12}$ .

(c)  $\mathcal{N}_{p_1 p_2} \cap \mathcal{U}_p$  has a  $C^\infty$  parametrization given by  $q \mapsto u_{12} \in B_p^{12}$ .

**Proof.** It is an application of the Proposition 32 based on the following projections:

$$\begin{aligned} B_p \ni u \mapsto & \left( \int u(\cdot, z) p_2(z) \, d\pi \, dz, \int u(y, \cdot) p_1(y) \, d\nu \, dy, u - \int u(\cdot, z) p_2(z) \pi \, dz \right. \\ & \left. + \int u(y, \cdot) p_1(y) \nu \, dy \right). \end{aligned} \quad \square$$

Note that the term  $u_{12}$  has the character of an interaction term.

### 7.4. An example

Rogantin (1996) developed in detail the application of the previous result to the case of finite spaces and to Markov chains. In the present paper we give a very simple example aimed at showing how each of the abstract results in the previous propositions correspond to the usual objects of the statistical analysis of statistical dependence (Cox and Reid 1987).

Let us consider the sample space  $\{0, 1\}^2$ ; let us take as a reference measure the uniform distribution, that is  $\mu_{ij} = \frac{1}{4}$  for  $i, j \in \{0, 1\}$ , and consider a neighbourhood of the unit density, i.e.  $p_{ij} = 1$  for  $i, j \in \{0, 1\}$ . The  $s_p$  coordinate of the generic density  $q_{ij}$ , with  $i, j \in \{0, 1\}$ , is given in matrix form by

$$u = \begin{bmatrix} \log q_{00} - \frac{1}{4} \sum_{ij} \log q_{ij} & \log q_{01} - \frac{1}{4} \sum_{ij} \log q_{ij} \\ \log q_{10} - \frac{1}{4} \sum_{ij} \log q_{ij} & \log q_{11} - \frac{1}{4} \sum_{ij} \log q_{ij} \end{bmatrix}$$

$$= \frac{1}{4} \begin{bmatrix} 3 \log q_{00} - \log q_{01} - \log q_{10} - \log q_{11} & -\log q_{00} + 3 \log q_{01} - \log q_{10} - \log q_{11} \\ -\log q_{00} - \log q_{01} + 3 \log q_{10} - \log q_{11} & -\log q_{00} - \log q_{01} - \log q_{10} + 3 \log q_{11} \end{bmatrix}.$$

We show in this example Proposition 34; the first marginal is given, i.e.  $\varphi_1(i, j) = i$ . The conditional mean value with respect to the first margin is

$$E_p[u|\varphi_1] \circ \varphi_1 = \frac{1}{4} \begin{bmatrix} \log q_{00} + \log q_{01} - \log q_{10} - \log q_{11} & \log q_{00} + \log q_{01} - \log q_{10} - \log q_{11} \\ \log q_{10} + \log q_{11} - \log q_{00} - \log q_{01} & \log q_{10} + \log q_{11} - \log q_{00} - \log q_{01} \end{bmatrix}$$

$$= \frac{1}{4} \log \left( \frac{q_{00}q_{01}}{q_{10}q_{11}} \right) \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix}$$

and

$$u - E_p[u|\varphi_1] \circ \varphi_1 = \frac{1}{2} \begin{bmatrix} \log \left( \frac{q_{01}}{q_{00}} \right) \\ \log \left( \frac{q_{10}}{q_{11}} \right) \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix}.$$

If we introduce the transition probabilities  $\alpha = q_{01}/(q_{00} + q_{01})$  and  $\beta = q_{10}/(q_{10} + q_{11})$ , then we find that

$$u - E_p[u|\varphi_1] \circ \varphi_1 = \frac{1}{2} \begin{bmatrix} \log \left( \frac{\alpha}{1 - \alpha} \right) \\ \log \left( \frac{\beta}{1 - \beta} \right) \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix}.$$

Thus this two-dimensional parameter coincides with the log-odds of the transition probabilities and it is orthogonal to mean parameter  $q_1 - 1$ , where

$$q_1 = E_\mu[q|\varphi_1] = \frac{1}{2} \begin{bmatrix} 4 - (q_{10} + q_{11}) \\ q_{10} + q_{11} \end{bmatrix}.$$

Now we show Proposition 36; the marginals are given, i.e.  $\varphi_1(i, j) = i$  and  $\varphi_2(i, j) = j$  and the splitting is  $u = u_1 + u_2 + u_{12}$ , where

$$u_1 = E_p[u|\varphi_1] \circ \varphi_1 = \frac{1}{4} \log \left( \frac{q_{00}q_{01}}{q_{10}q_{11}} \right) \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix},$$

$$u_2 = E_p[u|\varphi_2] \circ \varphi_2 = \frac{1}{4} \log \left( \frac{q_{00}q_{10}}{q_{01}q_{11}} \right) \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix};$$

so  $u_{12}$  is given by

$$u_{12} = u - u_1 - u_2 = \frac{1}{4} \log \left( \frac{q_{01}q_{10}}{q_{00}q_{11}} \right) \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix}.$$

In terms of transition probabilities,

$$\log \left( \frac{q_{01}q_{10}}{q_{00}q_{11}} \right) = \log \left( \frac{\alpha}{1-\alpha} \frac{\beta}{1-\beta} \right).$$

Then this parameter is orthogonal to the mean parameters  $q_1 - 1$  and  $q_2 - 1$ , where  $q_1$  is as above and

$$q_2 = E_\mu[q|\varphi_2] = \frac{1}{2} [4 - (q_{01} + q_{11}) \quad q_{01} + q_{11}].$$

## Acknowledgements

The present version of this paper has been greatly improved over the previous versions by working out many very useful constructive comments and meticulous remarks by the anonymous referees of the journal. The authors want to express to them and to the Editor their warmest thanks.

Some parts of the paper have been discussed with Damiano Brigo and Paolo Gibilisco who are currently working on further development of the theory of exponential statistical manifold (Brigo and Pistone 1996; Gibilisco and Pistone 1998). Their comments have contributed to the clarification of the presentation in many points.

## References

- Amari, S.-I. (1982) Differential geometry of curved exponential families—curvature and information loss. *Ann. Statist.*, **18**, 357–385.
- Amari, S.-I. (1985) *Differential-Geometrical Methods in Statistics*. Lecture Notes Statist., **28**. Berlin: Springer-Verlag.
- Amari, S.-I., Barndorff-Nielsen, O.E., Kass, R., Lauritzen, S.L. and Rao, C.R. (1987) *Differential Geometry and Statistical Inference*. Hayward, CA: Institute of Mathematical Statistics.

- Barndorff-Nielsen, O.E. (1978a) *Information and Exponential Families in Statistical Theory*. New York: Wiley.
- Barndorff-Nielsen, O.E. (1978b) *Parametric Statistical Models and Likelihood*. Lecture Notes Statist., **50**. Berlin: Springer-Verlag.
- Barndorff-Nielsen, O.E. and Cox, D.R. (1989) *Asymptotic Techniques for use in Statistics*. London: Chapman & Hall.
- Barndorff-Nielsen, O.E. and Cox, D.R. (1994) *Inference and Asymptotics*. London: Chapman & Hall.
- Barndorff-Nielsen, O.E. and Jupp, P.E. (1989) Approximating exponential models. *Ann. Inst. Statist. Math.*, **41**, 247–267.
- Brigo, D. and Pistone, G. (1996) Projecting the Fokker–Planck equation onto a finite dimensional exponential family. Preprint 4, Dipartimento di Matematica Pura e Applicata Università di Padova.
- Cox, D.R. and Reid, N. (1987) Parameter orthogonality and approximate conditional inference (with discussion). *J. Roy. Statist. Soc. B*, **49**, 1–39.
- Dawid, A.P. (1975) Discussion of a paper by B. Efron. *Ann. Statist.*, **3**, 1231–1234.
- Dawid, A.P. (1977) Further comments on a paper by Bradley Efron. *Ann. Statist.*, **5**, 1249.
- Donsker, M.D. and Varadhan, S.R.S. (1975) Asymptotic evaluation of certain Markov processes expectations for large time I. *Commun. Pure Appl. Math.*, **28**, 1–47.
- Efron, B. (1975) Defining the curvature of a statistical problem (with applications to second-order efficiency) (with discussion). *Ann. Statist.*, **3**, 1189–1242.
- Efron, B. (1978) The geometry of exponential families. *Ann. Statist.*, **6**, 362–376.
- Ekeland, I. and Temam, R. (1974) *Analyse Convexe et Problèmes Variationnels*. Paris: Dunod, Gauthier–Villars.
- Gibilisco, P. and Pistone, G. (1998) Connections on non-parametric statistical manifolds by Orlicz space geometry, infinite dimensional analysis. *Quantum Probab. Related Topics*, **1**, 325–347.
- Hardy, G.H., Littlewood, J.E. and Pólya, G. (1952) *Inequalities*, 2nd edn. London: Cambridge University Press.
- Jeffreys, H. (1946) An invariant form of the prior probability in estimation problems. *Proc. Roy. Soc. London A*, **196**, 453–461.
- Kass, R.E. (1989) The geometry of asymptotic inference (with discussion) *Statist. Sci.*, **4**, 188–234.
- Krasnosel'skii, M.A. and Rutickii, Ya.B. (1961) *Convex Functions and Orlicz Spaces*. Groningen: Noordhoff. (Russian original (1958) Moscow: Fizmatgiz.)
- Kullback, S. (1997) *Information Theory and Statistics*, reprint of 2nd (1968) edn., Mineola, NY: Dover.
- Kullback, S. and Leibler, R.A. (1951) On information and sufficiency. *Ann. Math. Statist.*, **22**, 79–86.
- Lang, S. (1995) *Differentiable Manifolds and Riemannian Manifolds*. New York: Springer-Verlag.
- Letac, G. (1992) *Lectures on Natural Exponential Families and their Variance Functions*, Monogr. Mat., **50**. Rio de Janeiro: Instituto de Matemática Pura e Aplicada.
- Madsen, L.T. (1979) The geometry of statistical model—a generalization of curvature. Research Report 79-1, Statistics Research Unit, Danish Medical Research Council.
- Murray, M.K. and Rice, J.W. (1993) *Differential Geometry and Statistics*, Monogr. Statist. Appl. Probab., **48**. London: Chapman & Hall.
- Neveu, J. (1972) *Martingales à Temps Discrets*. Paris: Masson.
- Pistone, G. and Rogantin, M.P. (1990) Gli strumenti della geometria differenziale nell'inferenza statistica. In P. Nastasi (ed.) *Memorial Beniamino Gulotta, Giornate di Lavoro di Probabilità e Statistica*, pp. 85–99. Palermo: Istituto Gramsci.
- Pistone, G. and Rogantin, M.P. (1994) *The Transformation of the Non-Parametric Statistical Manifold under Conditioning and Sampling*. Proceedings of the 57th IMS Annual Meeting and Third World Congress of the Bernoulli Society, Chapel Hill, NC, 20–25 June 1994.

- Pistone, G. and Sempi, C. (1995) An infinite dimensional geometric structure on the space of all the probability measures equivalent to a given one. *Ann. Statist.* **23**, 1543–1561.
- Rao, C.R. (1945) Information and accuracy attainable in the estimation of statistical parameters. *Bull. Calcutta Math. Soc.*, **37**, 81–89.
- Rao, C.R. (1949) On the distance between two populations. *Sankhyā*, **9**, 246–248.
- Rao, M.M. and Ren, Z.D. (1991) *Theory of Orlicz Spaces*. New York: Dekker.
- Rogantin, M.P. (1996) Geometrical modelling of Markovian dependence. *Metron*, **54**, 45–65.

Received April 1995 and revised October 1997