



**HAL**  
open science

# Minimum divergence estimators, maximum likelihood and exponential families

Michel Broniatowski

► **To cite this version:**

Michel Broniatowski. Minimum divergence estimators, maximum likelihood and exponential families. *Statistics and Probability Letters*, 2014, 93 (6), pp.27-33. hal-00613126v5

**HAL Id: hal-00613126**

<https://hal.sorbonne-universite.fr/hal-00613126v5>

Submitted on 19 Aug 2011

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Minimum divergence estimators, maximum likelihood and exponential families

Michel Broniatowski  
LSTA Université Pierre et Marie Curie  
5 Place Jussieu, 75005 Paris, France  
e-mail: michel.broniatowski@upmc.fr

## Abstract

In this note we prove the dual representation formula of the divergence between two distributions in a parametric model. Resulting estimators for the divergence as for the parameter are derived. These estimators do not make use of any grouping nor smoothing. It is proved that all differentiable divergences induce the same estimator of the parameter on any regular exponential family, which is nothing else but the MLE.

Key words: statistical divergence; minimum divergence estimator; maximum likelihood; exponential family

## 1 Introduction

### 1.1 Context and scope of this note

This note presents a short proof of the duality formula for  $\varphi$ -divergences defined through differentiable convex functions  $\varphi$  in parametric models and discusses some unexpected phenomenon in the context of exponential families. First versions of this formula appear in [8] p 33, in [1] in the context of the Kullback-Leibler divergence and in [7] in a general form. The paper [3] introduces this form in the context of minimal  $\chi^2$ -estimation; a global approach to this formulation is presented in Broniatowski and Kéziou (2006)[2]. Independently Liese and Vajda (2006)[9] have obtained a similar expression based on a much simpler argument as presented in all the above mentioned papers (formula (118) in their paper); however the proof of their result is merely sketched and we have found it useful to present a complete treatment of this interesting result in the parametric setting, in contrast with the aforementioned approaches.

The main interest of the resulting expression is that it leads to a wide variety of estimators, by a plug in method of the empirical measure evaluated on the current data set; so, for any type of sampling its estimators and inference procedures, for any  $\varphi$ -divergence criterion. In the case of the simple i.i.d. sampling resulting properties of those estimators and subsequent inferential procedures are studied in [4].

A striking fact is that all minimum divergence estimators defined through this dual formula coincide with the MLE in exponential families. They henceforth enjoy strong optimality under the standard exponential models, leading to estimators different from the MLE under different models than the exponential one. Also this result proves that MLE 's of parameters of exponential families are strongly motivated by being generated by the whole continuum of  $\varphi$ -divergences.

This note results from joint cooperation with late Igor Vajda.

## 1.2 Notation

Let  $\mathcal{P} := \{P_\theta, \theta \in \Theta\}$  an identifiable parametric model on  $\mathbb{R}^d$  where  $\Theta$  is a subset of  $\mathbb{R}^s$ . All measures in  $\mathcal{P}$  will be assumed to be measure equivalent sharing therefore the same support. The parameter space  $\Theta$  need not be open in the present setting. It may even happen that the model includes measures which would not be probability distributions; cases of interest cover models including mixtures of probability distributions; see [4]. Let  $\varphi$  be a proper closed convex function from  $] - \infty, +\infty[$  to  $[0, +\infty[$  with  $\varphi(1) = 0$  and such that its domain  $\text{dom}\varphi := \{x \in \mathbb{R} \text{ such that } \varphi(x) < \infty\}$  is an interval with endpoints  $a_\varphi < 1 < b_\varphi$  (which may be finite or infinite). For two measures  $P_\alpha$  and  $P_\theta$  in  $\mathcal{P}$  the  $\varphi$ -divergence between  $Q$  and  $P$  is defined by

$$\phi(\alpha, \theta) := \int_{\mathcal{X}} \varphi \left( \frac{dP_\alpha}{dP_\theta}(x) \right) dP_\theta(x).$$

In a broader context, the  $\varphi$ -divergences were introduced by [5] as “ $f$ -divergences”. The basic property of  $\varphi$ -divergences states that when  $\varphi$  is strictly convex on a neighborhood of  $x = 1$ , then

$$\phi(\alpha, \theta) = 0 \text{ if and only if } \alpha = \theta.$$

We refer to [8] chapter 1 for a complete study of those properties. Let us simply quote that in general  $\phi(\alpha, \theta)$  and  $\phi(\theta, \alpha)$  are not equal. Hence,  $\varphi$ -divergences usually are not distances, but they merely measure some difference between two measures. A main feature of divergences between distributions of random variables  $X$  and  $Y$  is the invariance property with respect to common smooth change of variables.

### 1.3 Examples of $\varphi$ -divergences

The Kullback-Leibler ( $KL$ ), modified Kullback-Leibler ( $KL_m$ ),  $\chi^2$ , modified  $\chi^2$  ( $\chi_m^2$ ), Hellinger ( $H$ ), and  $L_1$  divergences are respectively associated to the convex functions  $\varphi(x) = x \log x - x + 1$ ,  $\varphi(x) = -\log x + x - 1$ ,  $\varphi(x) = \frac{1}{2}(x-1)^2$ ,  $\varphi(x) = \frac{1}{2}(x-1)^2/x$ ,  $\varphi(x) = 2(\sqrt{x}-1)^2$  and  $\varphi(x) = |x-1|$ . All these divergences except the  $L_1$  one, belong to the class of the so called “power divergences” introduced in [6] (see also [8] chapter 2), a class which takes its origin from Rényi [10]. They are defined through the class of convex functions

$$x \in ]0, +\infty[ \mapsto \varphi_\gamma(x) := \frac{x^\gamma - \gamma x + \gamma - 1}{\gamma(\gamma - 1)} \quad (1)$$

if  $\gamma \in \mathbb{R} \setminus \{0, 1\}$ ,  $\varphi_0(x) := -\log x + x - 1$  and  $\varphi_1(x) := x \log x - x + 1$ . So, the  $KL$ -divergence is associated to  $\varphi_1$ , the  $KL_m$  to  $\varphi_0$ , the  $\chi^2$  to  $\varphi_2$ , the  $\chi_m^2$  to  $\varphi_{-1}$  and the Hellinger distance to  $\varphi_{1/2}$ .

It may be convenient to extend the definition of the power divergences in such a way that  $\phi(\alpha, \theta)$  may be defined (possibly infinite) even when  $P_\alpha$  or  $P_\theta$  is not a probability measure. This is achieved setting

$$x \in ]-\infty, +\infty[ \mapsto \begin{cases} \varphi_\gamma(x) & \text{if } x \in [0, +\infty[, \\ +\infty & \text{if } x \in ]-\infty, 0[. \end{cases} \quad (2)$$

when  $\text{dom} \varphi = \mathbb{R}^+ / \{0\}$ . Note that for the  $\chi^2$ -divergence, the corresponding  $\varphi$  function  $\phi_2(x) := \frac{1}{2}(x-1)^2$  is defined and convex on whole  $\mathbb{R}$ .

We will only consider divergences defined through differentiable functions  $\varphi$ , which we assume to satisfy

- There exists a positive  $\delta$  such that for all  $c$  in  $[1 - \delta, 1 + \delta]$ ,
- (RC)** we can find numbers  $c_1, c_2, c_3$  such that
- $$\varphi(cx) \leq c_1 \varphi(x) + c_2 |x| + c_3, \text{ for all real } x.$$

Condition **(RC)** holds for all power divergences including  $KL$  and  $KL_m$  divergences.

## 2 Dual form of the divergence and dual estimators in parametric models

Let  $\theta$  and  $\theta_T$  be any parameters in  $\Theta$ . We intend to provide a new expression for  $\phi(\theta, \theta_T)$ .

By strict convexity, for all  $a$  and  $b$  the domain of  $\varphi$  it holds

$$\varphi(b) \geq \varphi(a) + \varphi'(a)(b - a) \quad (3)$$

with equality if and only if  $a = b$ .

Denote

$$\varphi^\#(x) := x\varphi'(x) - \varphi(x).$$

For any  $\alpha$  in  $\Theta$  denote

$$a := \frac{dP_\theta}{dP_\alpha}(x).$$

Define

$$b := \frac{dP_\theta}{dP_{\theta_T}}(x).$$

Inserting these values in (3) and integrating with respect to  $P_{\theta_T}$  yields

$$\phi(\theta, \theta_T) \geq \int \left[ \varphi' \left( \frac{dP_\theta}{dP_\alpha} \right) dP_\theta - \varphi^\# \left( \frac{dP_\theta}{dP_\alpha} \right) \right] dP_{\theta_T}.$$

Assume at present that this entails

$$\phi(\theta, \theta_T) \geq \int \varphi' \left( \frac{dP_\theta}{dP_\alpha} \right) dP_\theta - \int \varphi^\# \left( \frac{dP_\theta}{dP_\alpha} \right) dP_{\theta_T} \quad (4)$$

for suitable  $\alpha$ 's in some set  $\mathcal{F}_\theta$  included in  $\Theta$ .

When  $\alpha = \theta_T$  the inequality in (4) turns to equality, which yields

$$\phi(\theta, \theta_T) = \sup_{\alpha \in \mathcal{F}_\theta} \int \varphi' \left( \frac{dP_\theta}{dP_\alpha} \right) dP_\theta - \int \varphi^\# \left( \frac{dP_\theta}{dP_\alpha} \right) dP_{\theta_T} \quad (5)$$

Denote

$$h(\theta, \alpha, x) := \int \varphi' \left( \frac{dP_\theta}{dP_\alpha} \right) dP_\theta - \varphi^\# \left( \frac{dP_\theta}{dP_\alpha} \right) \quad (6)$$

from which

$$\phi(\theta, \theta_T) = \sup_{\alpha \in \mathcal{F}_\theta} \int h(\theta, \alpha, x) dP_{\theta_T}. \quad (7)$$

Furthermore by (4), for all suitable  $\alpha$

$$\begin{aligned} & \phi(\theta, \theta_T) - \int h(\theta, \alpha, x) dP_{\theta_T} \\ &= \int h(\theta, \theta_T, x) dP_{\theta_T} - \int h(\theta, \alpha, x) dP_{\theta_T} \geq 0 \end{aligned}$$

and the function  $x \rightarrow h(\theta, \theta_T, x) - h(\theta, \alpha, x)$  is non negative, due to (3). It follows that  $\phi(\theta, \theta_T) - \int h(\theta, \alpha, x) dP_{\theta_T}$  is zero only if  $h(\theta, \alpha, x) = h(\theta, \theta_T, x) - P_{\theta_T}$  a.e. Therefore for any  $x$  in the support of  $P_{\theta_T}$

$$\left[ \int \varphi' \left( \frac{dP_\theta}{dP_{\theta_T}} \right) dP_\theta - \int \varphi' \left( \frac{dP_\theta}{dP_\alpha} \right) dP_\theta \right] - \varphi^\# \left( \frac{dP_\theta}{dP_\alpha}(x) \right) + \varphi^\# \left( \frac{dP_\theta}{dP_{\theta_T}}(x) \right) = 0$$

which cannot hold for all  $x$  when the functions  $\varphi^\# \left( \frac{dP_\theta}{dP_\alpha}(x) \right)$ ,  $\varphi^\# \left( \frac{dP_\theta}{dP_{\theta_T}}(x) \right)$  and 1 are linearly independent, unless  $\alpha = \theta_T$ . We have proved that  $\theta_T$  is the unique optimizer in (5).

We have skipped some sufficient conditions which ensure that (4) holds.

Assume that

$$\int \left| \varphi' \left( \frac{dP_\theta}{dP_\alpha} \right) \right| dP_\theta < \infty. \quad (8)$$

Assume further that  $\phi(\theta, \theta_T)$  is finite. Since

$$\begin{aligned} - \int \varphi^\# \left( \frac{dP_\theta}{dP_\alpha}(x) \right) dP_{\theta_T} &\leq \phi(\theta, \theta_T) - \int \varphi' \left( \frac{dP_\theta}{dP_{\theta_T}} \right) dP_\theta \\ &\leq \phi(\theta, \theta_T) + \int \left| \varphi' \left( \frac{dP_\theta}{dP_{\theta_T}} \right) \right| dP_\theta < +\infty \end{aligned}$$

we obtain

$$\int \varphi^\# \left( \frac{dP_\theta}{dP_\alpha} \right) dP_{\theta_T} > -\infty$$

which entails (4). When  $\int \varphi^\# \left( \frac{dP_\theta}{dP_\alpha} \right) dP_{\theta_T} = +\infty$  then clearly, under (8)

$$\phi(\theta, \theta_T) > \int \varphi' \left( \frac{dP_\theta}{dP_\alpha} \right) dP_\theta - \int \varphi^\# \left( \frac{dP_\theta}{dP_\alpha} \right) dP_{\theta_T} = -\infty.$$

We have proved that (5) holds when  $\alpha$  satisfies (8).

Sufficient and simple conditions encompassing (8) can be assessed under standard requirements for nearly all divergences. We state the following Lemma (see Liese and Vajda (1987)[8]) and Broniatowski and Kéziou (2006) [2], Lemma 3.2).

**Lemma 1** *Assume that **RC** holds and  $\phi(\theta, \alpha)$  is finite. Then (8) holds.*

Summing up, we state

**Theorem 2** *Let  $\theta$  belong to  $\Theta$  and let  $\phi(\theta, \theta_T)$  be finite. Assume that **RC** holds. Let  $\mathcal{F}_\theta$  be the subset of all  $\alpha$ 's in  $\Theta$  such that  $\phi(\theta, \alpha)$  is finite. Then*

$$\phi(\theta, \theta_T) = \sup_{\alpha \in \mathcal{F}_\theta} \int \varphi' \left( \frac{dP_\theta}{dP_\alpha} \right) dP_\theta - \int \varphi^\# \left( \frac{dP_\theta}{dP_\alpha} \right) dP_{\theta_T}.$$

*Furthermore the sup is reached at  $\theta_T$  and uniqueness holds.*

For the Cressie-Read family of divergences with  $\gamma \neq 0, 1$  this representation writes

$$\phi_\gamma(\theta, \theta_T) = \sup_{\alpha \in \mathcal{F}_\theta} \left\{ \frac{1}{\gamma - 1} \int \left( \frac{dP_\theta}{dP_\alpha} \right)^{\gamma-1} dP_\theta - \frac{1}{\gamma} \int \left( \frac{dP_\theta}{dP_\alpha} \right)^\gamma dP_{\theta_T} - \frac{1}{\gamma(\gamma - 1)} \right\}.$$

The set  $\mathcal{F}_\theta$  may depend on the choice of the parameter  $\theta$ . Such is the case for the  $\chi^2$  divergence i.e.  $\varphi(x) = (x - 1)^2/2$ , when  $p_\theta(x) = \theta \exp(-\theta x) 1_{[0, \infty)}(x)$ . In most cases the difficulty of dealing with a specific set  $\mathcal{F}_\theta$  depending on  $\theta$  can be encompassed when

$$\begin{aligned} &\text{There exists a neighborhood } \mathcal{U} \text{ of } \theta_T \text{ for which} & (A) \\ &\phi(\theta, \theta') \text{ is finite whatever } \theta \text{ and } \theta' \text{ in } \mathcal{U} \end{aligned}$$

which for example holds in the above case for any  $\theta_T$ . This simplification deserves to be stated in the next result

**Theorem 3** *When  $\phi(\theta, \theta_T)$  is finite and **RC** holds, then under condition (A)*

$$\phi(\theta, \theta_T) = \sup_{\alpha \in \mathcal{U}} \int \varphi' \left( \frac{dP_\theta}{dP_\alpha} \right) dP_\theta - \int \varphi^\# \left( \frac{dP_\theta}{dP_\alpha} \right) dP_{\theta_T}.$$

Furthermore the sup is reached at  $\theta_T$  and uniqueness holds.

**Remark 4** *Identifying  $\mathcal{F}_\theta$  might be cumbersome. This difficulty also appears in the classical MLE case, a special case of the above statement with divergence function  $\varphi_0$ , for which it is assumed that*

$$\int \log p_\theta(x) p_{\theta_T}(x) d\lambda(x) \text{ is finite}$$

for  $\theta$  in a neighborhood of  $\theta_T$ .

Under the above notation and hypotheses define

$$T_\theta(P_{\theta_T}) := \arg \sup_{\alpha \in \mathcal{F}_\theta} \int h(\theta, \alpha, x) dP_{\theta_T}. \quad (9)$$

It then holds

$$T_\theta(P_{\theta_T}) = \theta_T$$

for all  $\theta_T$  in  $\Theta$ . Also let

$$S(P_{\theta_T}) := \arg \inf_{\theta \in \Theta} \sup_{\alpha \in \mathcal{F}_\theta} \int h(\theta, \alpha, x) dP_{\theta_T}. \quad (10)$$

which also satisfies

$$S(P_{\theta_T}) = \theta_T$$

for all  $\theta_T$  in  $\Theta$ . We thus state

**Theorem 5** *When  $\phi(\theta, \theta_T)$  is finite for all  $\theta$  in  $\Theta$  and **RC** holds, both functionals  $T_\theta$  and  $S$  are Fisher consistent for all  $\theta_T$  in  $\Theta$ .*

### 3 Plug in estimators

From (7) simple estimators for  $\theta_T$  can be defined, plugging any convergent empirical measure in place of  $P_{\theta_T}$  and taking the infimum in  $\theta$  in the resulting estimator of  $\phi(\theta, \theta_T)$ .

In the context of simple i.i.d. sampling, introducing the empirical measure

$$P_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$$

where the  $X_i$ 's are i.i.d. r.v's with common unknown distribution  $P_{\theta_T}$  in  $\mathcal{P}$ , the natural estimator of  $\phi(\theta, \theta_T)$  is

$$\begin{aligned} \phi_n(\theta, \theta_T) &:= \sup_{\alpha \in \mathcal{F}_\theta} \left\{ \int h(\theta, \alpha, x) dP_n(x) \right\} \\ &= \sup_{\alpha \in \mathcal{F}_\theta} \int \varphi' \left( \frac{dP_\theta}{dP_\alpha} \right) dP_\theta - \frac{1}{n} \sum_{i=1}^n \varphi^\# \left( \frac{dP_\theta}{dP_\alpha} (X_i) \right). \end{aligned} \quad (11)$$

Since

$$\inf_{\theta \in \Theta} \phi(\theta, \theta_T) = \phi(\theta_T, \theta_T) = 0$$

the resulting estimator of  $\phi(\theta_T, \theta_T)$  is

$$\phi_n(\theta_T, \theta_T) := \inf_{\theta \in \Theta} \phi_n(\theta, \theta_T) = \inf_{\theta \in \Theta} \sup_{\alpha \in \mathcal{F}_\theta} \left\{ \int h(\theta, \alpha, x) dP_n(x) \right\}. \quad (12)$$

Also the estimator of  $\theta_T$  is obtained as

$$\hat{\theta} := \arg \inf_{\theta \in \Theta} \sup_{\alpha \in \mathcal{F}_\theta} \left\{ \int h(\theta, \alpha, x) dP_n(x) \right\}. \quad (13)$$

When A holds then  $\mathcal{F}_\theta$  may be substituted by  $\mathcal{U}$  in the above definitions.

The resulting minimum dual divergence estimators (12) and (13) do not require any smoothing or grouping, in contrast with the classical approach which involves quantization. The paper [4] provides a complete study of those estimates and subsequent inference tools for the usual i.i.d. sample scheme. For all divergences considered here, these estimators are asymptotically efficient in the sense that they achieve the Cramer-Rao bound asymptotically. The case when  $\varphi = \varphi_0$  leads to  $\theta_{ML}$  defined as the celebrated Maximum Likelihood Estimator (MLE), in the context of the simple sampling.



## 4 Minimum divergence estimators in exponential families

In this section we prove the following result

**Theorem 6** *For all divergence  $\phi$  defined through a differentiable function  $\varphi$  satisfying Condition **(RC)**, the minimum dual divergence estimator defined by (13) coincides with the MLE on any full exponential families such that  $\phi(\theta, \alpha)$  is finite for all  $\theta$  and  $\alpha$  in  $\Theta$ .*

Let  $\mathcal{P}$  be an exponential family on  $\mathbb{R}^s$  with canonical parameter in  $\mathbb{R}^d$

$$\mathcal{P} := \left\{ P_\theta \text{ such that } p_\theta(x) = \frac{dP_\theta}{d\lambda}(x) \right\} \\ = \left\{ \exp [T(x)' \theta - C(\theta)]; \theta \in \Theta \right\}$$

where  $x$  is in  $\mathbb{R}^s$  and  $\Theta$  is an open subset of  $\mathbb{R}^d$ , and  $\lambda$  is a dominating measure for  $\mathcal{P}$ . We assume  $\mathcal{P}$  to be full, namely that the Hessian matrix  $(\partial^2/\partial\theta^2) C(\theta)$  is definite positive for all  $\theta$  in  $\Theta$ .

Let  $X_1, \dots, X_n$  be  $n$  i.i.d. random variables with common distribution  $P_{\theta_T}$  with  $\theta_T$  in  $\Theta$ . Introduce

$$M_n(\theta, \alpha) := \int \varphi' \left( \frac{dP_\theta}{dP_\alpha} \right) dP_\theta - \frac{1}{n} \sum_{i=1}^n \varphi^\# \left( \frac{dP_\theta}{dP_\alpha}(X_i) \right)$$

We will prove that

$$\inf_{\theta} \sup_{\alpha} M_n(\theta, \alpha) = 0 \quad (14)$$

whatever the function  $\varphi$  satisfying the claim. In (14)  $\theta$  and  $\alpha$  run in  $\Theta$ . This result extends the maximum likelihood case for which  $\inf_{\theta} \sup_{\alpha} M_n(\theta, \alpha) = \sup_{\theta} \inf_{\alpha} \left[ \frac{1}{n} \sum_{i=1}^n \log p_\theta(X_i) - \frac{1}{n} \sum_{i=1}^n \log p_\alpha(X_i) \right] = 0$ .

Direct substitution shows that for any  $\theta$ ,

$$\sup_{\alpha} M_n(\theta, \alpha) \geq M_n(\theta, \theta) = 0$$

from which

$$\inf_{\theta} \sup_{\alpha} M_n(\theta, \alpha) \geq 0 \quad (15)$$

We prove that

$$\alpha = \theta_{ML} \text{ is the unique maximizer of } M_n(\theta_{ML}, \alpha) \quad (16)$$

which yields

$$\inf_{\theta} \sup_{\alpha} M_n(\theta, \alpha) \leq \sup_{\alpha} M_n(\theta_{ML}, \alpha) = M_n(\theta_{ML}, \theta_{ML}) = 0 \quad (17)$$

which together with (15) completes the proof.

Define

$$\begin{aligned} M_{n,1}(\theta, \alpha) &:= \int \varphi'(\exp A(\theta, \alpha, x)) \exp B(\theta, x) d\lambda(x) \\ M_{n,2}(\theta, \alpha) &:= \frac{1}{n} \sum_{i=1}^n \exp(A(\theta, \alpha, X_i)) \varphi'(\exp A(\theta, \alpha, X_i)) \\ M_{n,3}(\theta, \alpha) &:= \frac{1}{n} \sum_{i=1}^n \varphi(\exp A(\theta, \alpha, X_i)) \end{aligned}$$

with

$$\begin{aligned} A(\theta, \alpha, x) &:= T(x)'(\theta - \alpha) + C(\alpha) - C(\theta) \\ B(\theta, x) &:= T(x)'\theta - C(\theta). \end{aligned}$$

It holds

$$M_n(\theta, \alpha) = M_{n,1}(\theta, \alpha) - M_{n,2}(\theta, \alpha) + M_{n,3}(\theta, \alpha)$$

with

$$\frac{\partial}{\partial \alpha} M_{n,1}(\theta, \alpha)_{\alpha=\theta} = -\varphi^{(2)}(1) [\nabla C(\theta) - \nabla C(\alpha)_{\alpha=\theta}] = 0$$

for all  $\theta$ ,

$$\frac{\partial}{\partial \alpha} M_{n,2}(\theta_{ML}, \alpha)_{\alpha=\theta_{ML}} = \varphi^{(2)}(1) \frac{1}{n} \sum_{i=1}^n [-T(X_i) + \nabla C(\alpha)_{\alpha=\theta_{ML}}] = 0$$

and

$$\frac{\partial}{\partial \alpha} M_{n,3}(\theta_{ML}, \alpha) = \frac{1}{n} \sum_{i=1}^n [-T(X_i) + \nabla C(\alpha)_{\alpha=\theta_{ML}}] = 0$$

where the two last displays hold iff  $\alpha = \theta_{ML}$ . Now

$$\begin{aligned} \frac{\partial^2}{\partial \alpha^2} M_{n,1}(\theta_{ML}, \alpha)_{\alpha=\theta_{ML}} &= (\varphi^{(3)}(1) + 2\varphi^{(2)}(1)) (\partial^2 / \partial \theta^2) C(\theta_{ML}) \\ \frac{\partial^2}{\partial \alpha^2} M_{n,2}(\theta_{ML}, \alpha)_{\alpha=\theta_{ML}} &= (\varphi^{(3)}(1) + 4\varphi^{(2)}(1)) (\partial^2 / \partial \theta^2) C(\theta_{ML}) \\ \frac{\partial^2}{\partial \alpha^2} M_{n,3}(\theta_{ML}, \alpha)_{\alpha=\theta_{ML}} &= \varphi^{(2)}(1) (\partial^2 / \partial \theta^2) C(\theta_{ML}), \end{aligned}$$

whence

$$\begin{aligned}\frac{\partial}{\partial \alpha} M_n(\theta_{ML}, \alpha)_{\alpha=\theta_{ML}} &= 0 \\ \frac{\partial^2}{\partial \alpha^2} M_n(\theta_{ML}, \alpha)_{\alpha=\theta_{ML}} &= -\varphi^{(2)}(1) (\partial^2 / \partial \theta^2) C(\theta_{ML})\end{aligned}$$

which proves (16), and closes the proof.

## References

- [1] Broniatowski, M. Estimation of the Kullback-Leibler divergence. *Math. Methods Statist.* 12 (2003), no. 4, 391–409 .
- [2] Broniatowski, M.; Keziou, A. Minimization of  $\varphi$ -divergences on sets of signed measures. *Studia Sci. Math. Hungar.* 43 (2006), no. 4, 403–442.
- [3] Broniatowski, M.; Leorato, S. An estimation method for the Neyman chi-square divergence with application to test of hypotheses. *J. Multivariate Anal.* 97 (2006), no. 6, 1409–1436.
- [4] Broniatowski, M. Keziou, A. Parametric estimation and tests through divergences and the duality technique. *J. Multivariate Anal.* 100 (2009), no. 1, 16–36.
- [5] Csiszár, I. Eine informationstheoretische Ungleichung und ihre Anwendung auf den Beweis der Ergodizität von Markoffschen Ketten. (German) *Magyar Tud. Akad. Mat. Kutató Int. Közl.* 8 1963 85–108.
- [6] Read, T. R. C., Cressie, N. A. C. Goodness-of-fit statistics for discrete multivariate data. Springer Series in Statistics. Springer-Verlag, New York, 1988. xii+211 pp. ISBN: 0-387-96682-X
- [7] Keziou, A. Dual representation of  $\varphi$ -divergences and applications. *C. R. Math. Acad. Sci. Paris* 336 (2003), no. 10, 857–862
- [8] Liese, F., Vajda, I. Convex statistical distances. Teubner-Texte zur Mathematik [Teubner Texts in Mathematics], 95. BSB B. G. Teubner Verlagsgesellschaft, Leipzig, 1987, ISBN: 3-322-00428-7 .
- [9] Liese, F., Vajda, I. On divergences and informations in statistics and information theory. *IEEE Trans. Inform. Theory* 52 (2006), no. 10, 4394–4412
- [10] Rényi, A. On measures of entropy and information. 1961 Proc. 4th Berkeley Sympos. Math. Statist. and Prob., Vol. I pp. 547–561 Univ. California Press, Berkeley, Calif.