

Article

Relative Entropy Derivative Bounds

Pablo Zegers *, Alexis Fuentes and Carlos Alarcón

Universidad de los Andes, Facultad de Ingeniería y Ciencias Aplicadas, Monseñor Álvaro del Portillo 12455, Las Condes, Santiago, Chile; E-Mails: afuentes@miuandes.cl (A.F.); ceag.201175@gmail.com (C.A.)

* Author to whom correspondence should be addressed; E-Mail: pzegers@miuandes.cl; Tel.: +56-226-189-324.

Received: 24 May 2013; in revised form: 12 July 2013 / Accepted: 16 July 2013 / Published: 23 July 2013

Abstract: We show that the derivative of the relative entropy with respect to its parameters is lower and upper bounded. We characterize the conditions under which this derivative can reach zero. We use these results to explain when the minimum relative entropy and the maximum log likelihood approaches can be valid. We show that these approaches naturally activate in the presence of large data sets and that they are inherent properties of any density estimation process involving large numbers of random variables.

Keywords: relative entropy; Kullback-Leibler divergence; Shannon differential entropy; asymptotic equipartition principle; typical set; Fisher information; maximum log likelihood

1. Introduction

Given a large ensemble of i.i.d.random variables, all of them generated according to some common density function, the asymptotic equipartition principle [1] guarantees that only a fraction of the possible ensembles, which is called the typical set, gathers almost all the probability ([2] [p. 226]). This fact opens interesting possibilities: What if the properties of the typical set impose conditions on any estimation process, such that the parameter search is focused on just a small subset of the available parameter set? We explore this using generalizations of the Shannon differential entropy [1] and the Fisher information [3–5] under typical set considerations. There are other works that take into account both concepts [6–8]. However, this work focuses on something new: defining algebraic bounds on the behavior of the derivative of the relative entropy with respect to their parameters. Furthermore, we

characterize the conditions under which seeking for the extreme of this expression is a valid approach. Most importantly, we prove that these conditions automatically activate if large data sets are used. We finish this work discussing the relation between these bounds and the density function estimation problem. Under the conditions that activate these bounds, we characterize when the minimum relative entropy and maximum log likelihood approaches render the same solution. Furthermore, we show that these equivalences become true by the sole fact of having a large number of random variables.

2. Known Density Functions Case

Let us assume for the ensuing discussion that there exists a known source density function, f_x , with support, S. We also know a set of i.i.d. samples, $\{x\}_n \equiv \{x_1, \ldots, x_n\}$, that was generated by f_x . There is also another density function, $f_{Y|\theta}(y)$, with the same support, S. This density function is indexed by the parameter vector, $\theta \in \Theta \subseteq \mathbb{R}^m$.

The first part of this work deals with this case, where all the density functions are known. The second part of this work studies the density function estimation problem, where the source density function, f_x , is not known.

3. Mixed Entropy Derivative: A Proxy for the Relative Entropy Derivative

We present the following definition:

Definition 1 The relative entropy ([2][p. 231]), also called Kullback-Leibler divergence, is defined by:

$$D_{RE}\left(f_{\boldsymbol{X}} \parallel f_{\boldsymbol{Y}|\boldsymbol{\theta}}\right) \equiv \int_{\mathcal{S}} f_{\boldsymbol{X}}(\boldsymbol{u}) \ln \frac{f_{\boldsymbol{X}}(\boldsymbol{u})}{f_{\boldsymbol{Y}|\boldsymbol{\theta}}(\boldsymbol{u})} d\boldsymbol{u}$$
(1)

Notice that $\lim_{u\to 0^+} u \ln u = 0$. Even though the definition is valid for any pair of density functions, we are using those useful for the ensuing analysis. The domain of this integral, and that of all the integrals in this work, corresponds to the support, S.

Definition 2 The mixed entropy, $h_M(f_X, f_{Y|\theta})$ of two density functions is defined by:

$$h_{M}\left(f_{\boldsymbol{X}}, f_{\boldsymbol{Y}|\boldsymbol{\theta}}\right) \equiv -\int_{\mathcal{S}} f_{\boldsymbol{X}}(\boldsymbol{u}) \ln f_{\boldsymbol{Y}|\boldsymbol{\theta}}(\boldsymbol{u}) d\boldsymbol{u}$$
(2)

when the integral exists.

Definition 3 *The Shannon differential entropy* [1] *is defined by:*

$$h_{s}(f_{\mathbf{x}}) \equiv -\int_{\mathcal{S}} f_{\mathbf{x}}(\mathbf{u}) \ln f_{\mathbf{x}}(\mathbf{u}) d\mathbf{u}$$
(3)

From these definitions, when $f_{Y|\theta} = f_x$, from a functional point of view, the mixed entropy, $h_M(f_x, f_{Y|\theta})$, equals the differential entropy, $h_S(f_x)$

From examination of the relative entropy definition:

$$D_{KL} (f_{\mathbf{X}} \parallel f_{\mathbf{Y}|\boldsymbol{\theta}}) = \int f_{\mathbf{X}}(\boldsymbol{u}) \ln f_{\mathbf{X}}(\boldsymbol{u}) d\boldsymbol{u} - \int f_{\mathbf{X}}(\boldsymbol{u}) \ln f_{\mathbf{Y}|\boldsymbol{\theta}}(\boldsymbol{u}) d\boldsymbol{u}$$

= $-h_{S}(f_{\mathbf{X}}) + h_{M} (f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}})$ (4)

Hence:

$$\frac{\partial D_{KL}\left(f_{\mathbf{X}} \parallel f_{\mathbf{Y}|\boldsymbol{\theta}}\right)}{\partial \theta_{k}} = \frac{h_{M}\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}}\right)}{\partial \theta_{k}}$$
(5)

In the previous equation it was assumed that the $\frac{\partial}{\partial \theta_k}$ are defined in the interior of Θ for $k \in \{1, \ldots, m\}$. In this work, where we assume that f_x does not depend on any θ_k , we use extensively the fact that the relative entropy derivative is equal to the mixed entropy derivative. Hence, in the following, we focus on studying the properties of the mixed entropy in order to be able to say something about the properties of the relative entropy derivatives.

4. Mixed Entropy Typical Set

An interesting insight related to sequences of i.i.d. random variables, discovered by Shannon [1], is the usefulness of the weak law of large numbers to characterize large ensembles of random variables. In the context of this work, this law implies:

Defining

Definition 4 The likelihood, $f_{\{\mathbf{Y}\}_n|\theta}$, is defined by:

$$f_{\{\mathbf{Y}\}_n|\boldsymbol{\theta}} \equiv \prod_{k=1}^n f_{\mathbf{Y}|\boldsymbol{\theta}}(\boldsymbol{x}_k)$$
(6)

we can state:

Theorem 1 For any positive real number, ε , and fixed parameter vector, θ :

$$\mathbb{P}\left\{\left|-\frac{1}{n}\ln f_{\{\mathbf{Y}\}_{n}\mid\boldsymbol{\theta}}-h_{M}\left(f_{\mathbf{X}},f_{\mathbf{Y}\mid\boldsymbol{\theta}}\right)\right|\leq\varepsilon\right\}\xrightarrow[n\to\infty]{}1$$
(7)

Proof: Use that the random variables are i.i.d., and the weak law of large numbers. \Box

Hence, it makes sense to define the following set ([2][p. 226]):

Definition 5 For any positive real number, ε , fixed parameter vector, θ , and any $n \ge 1$, the mixed entropy typical set, $\mathcal{M}_{\varepsilon}^{n}$, with respect to the density function, $f_{\mathbf{x}}$, is defined by:

$$\mathcal{M}_{\varepsilon}^{n} = \left\{ \{ \boldsymbol{x} \}_{n} \in \mathcal{S}^{n} : \left| -\frac{1}{n} \ln f_{\{\boldsymbol{Y}\}_{n}|\boldsymbol{\theta}} - h_{M} \left(f_{\boldsymbol{X}}, f_{\boldsymbol{Y}|\boldsymbol{\theta}} \right) \right| \leq \varepsilon \right\}$$
(8)

Furthermore, using the weak law of large numbers, it is possible to prove:

Lemma 1 For any positive real number, ε , and any parameter vector, θ , it is true that:

$$\mathbb{P}\left\{\left\{\boldsymbol{x}\right\}_{n}\in\mathcal{M}_{\varepsilon}^{n}\right\}\xrightarrow[n\to\infty]{}1$$
(9)

Proof: Use the weak law of large numbers in the mixed entropy typical set definition. \Box

This proves that for large values of n, almost all the probability is contained by the mixed entropy typical set and that the probability of being outside this set becomes negligible.

Then, it is possible to prove that:

Lemma 2 For any positive real number, ε , and any fixed parameter vector, θ , it exists $n_0 \equiv n_0(\varepsilon, \theta)$, such that for all $n > n_0$, the total probability of the typical set is lower bounded by:

$$1 - \varepsilon \le \mathbb{P}\left\{\mathcal{M}^n_\varepsilon\right\} \tag{10}$$

Proof: Check ([2][p. 226], [9][p. 118]) for details. □

5. Micro-Differences in Mixed Entropy Typical Sets

Assuming that $\{x\}_n \in \mathcal{M}_{\varepsilon}^n$ and from the definition of mixed entropy typical set:

$$\ln f_{\{\mathbf{Y}\}_n|\boldsymbol{\theta}} \in \left] - nh_M\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}}\right) - n\varepsilon, -nh_M\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}}\right) + n\varepsilon\right[$$
(11)

The width of this interval is $2\varepsilon n$. This allows us to compare $\ln f_{\{Y\}_n|\theta}$ with $-nh_M(f_Y, f_{Y|\theta})$, the value that is obtained in the limit thanks to the weak law of large numbers. Their ratio is:

$$\left|\frac{2\varepsilon n}{-nh_{M}\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}}\right)}\right| = \left|\frac{2\varepsilon}{h_{M}\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}}\right)}\right|$$
(12)

In other words, given that ε can be chosen as small as needed, the range of values that $\ln f_{\{Y\}_n|\theta}$ can take may be indistinguishable from $-nh_M(f_X, f_{Y|\theta})$, with high probability. This is another way of stating the weak law of large numbers.

The previous facts allow us to present the following definition:

Definition 6 When $\{x\}_n \in \mathcal{M}^n_{\varepsilon}$, its associated micro-difference is defined by the following function:

$$\delta \equiv \delta(f_{\mathbf{x}}, n, \{\mathbf{x}\}_n, f_{\mathbf{Y}|\boldsymbol{\theta}}, \boldsymbol{\theta}, \varepsilon) \in]0, 1[$$
(13)

such that:

$$\ln f_{\{\mathbf{Y}\}_n|\boldsymbol{\theta}} = -n(h_M\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}}\right) + \varepsilon) + 2\varepsilon n\delta$$
(14)

hence:

$$\delta = \frac{1}{2\varepsilon n} \left(\ln f_{\{\mathbf{Y}\}_n \mid \boldsymbol{\theta}} + n(h_M \left(f_{\mathbf{X}}, f_{\mathbf{Y} \mid \boldsymbol{\theta}} \right) + \varepsilon) \right)$$
(15)

The analysis performed in the preceding paragraphs shows that the behavior of the micro-differences might be negligible in the case of large values of n. However, the following sections of this work deal with the the derivative of the micro-difference value with respect to the parameters; so, we cannot neglect it, and we need to study its behavior.

From Equation (14):

$$\frac{\partial \ln f_{\{\mathbf{Y}\}_n|\boldsymbol{\theta}}}{\partial \theta_k} = -n \frac{\partial h_M \left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}} \right)}{\partial \theta_k} + 2\varepsilon n \frac{\partial \delta}{\partial \theta_k} \tag{16}$$

for all $k \in [1, \ldots, m]$. Thus:

$$2\varepsilon \frac{\partial \delta}{\partial \theta_k} = \frac{1}{n} \frac{\partial \ln f_{\{\mathbf{Y}\}_n \mid \boldsymbol{\theta}}}{\partial \theta_k} + \frac{\partial h_M \left(f_{\mathbf{X}}, f_{\mathbf{Y} \mid \boldsymbol{\theta}} \right)}{\partial \theta_k}$$
(17)

$$= \frac{\partial}{\partial \theta_k} \left(\frac{1}{n} \ln f_{\{\mathbf{Y}\}_n \mid \boldsymbol{\theta}} + h_M \left(f_{\mathbf{X}}, f_{\mathbf{Y} \mid \boldsymbol{\theta}} \right) \right)$$
(18)

$$= -\frac{\partial}{\partial \theta_k} \left(-\frac{1}{n} \ln f_{\{\mathbf{Y}\}_n | \boldsymbol{\theta}} - h_M \left(f_{\mathbf{X}}, f_{\mathbf{Y} | \boldsymbol{\theta}} \right) \right)$$
(19)

This last expression helps us to understand that the derivative of δ with respect to the parameters is related to changes experienced by quantities upper bounded by ε .

6. Mixed Information

We define:

Definition 7 The mixed information, $i_F (f_{\{\mathbf{X}\}_n}, f_{\{\mathbf{Y}\}_n \mid \theta})_k$, is defined by:

$$i_{F}\left(f_{\{\boldsymbol{X}\}_{n}},f_{\{\boldsymbol{Y}\}_{n}|\boldsymbol{\theta}}\right)_{k}\equiv\int_{\mathcal{S}^{n}}f_{\{\boldsymbol{X}\}_{n}}\left(\frac{\partial\ln f_{\{\boldsymbol{Y}\}_{n}|\boldsymbol{\theta}}}{\partial\theta_{k}}\right)^{2}d\boldsymbol{u}$$
(20)

with:

$$f_{\{\boldsymbol{x}\}_n} \equiv \prod_{k=1}^n f_{\boldsymbol{x}}(\boldsymbol{x}_k) \tag{21}$$

when the integral exists.

7. Mixed Entropy Derivative Bounds

The main contribution of this work is the following pair of inequalities, which are obtained thanks to a combination of the mixed information and the mixed entropy.

Theorem 2 Given a positive real number, ε , any parameter vector, θ , and an i.i.d. sequence with n elements, then the components of $\nabla_{\theta} h_M(f_{\mathbf{x}}, f_{\mathbf{y}|\theta})$ comply with:

$$\frac{\sqrt{I_{R}^{k}} - \frac{1}{n}\sqrt{i_{F}\left(f_{\{\mathbf{X}\}_{n}}, f_{\{\mathbf{Y}\}_{n}|\boldsymbol{\theta}}\right)_{k}}}{\sqrt{I_{L}}} \leq \left|\frac{\partial h_{M}\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}}\right)}{\partial \theta_{k}}\right| \leq \frac{\sqrt{I_{R}^{k}} + \frac{1}{n}\sqrt{i_{F}\left(f_{\{\mathbf{X}\}_{n}}, f_{\{\mathbf{Y}\}_{n}|\boldsymbol{\theta}}\right)_{k}}}{\sqrt{I_{L}}} \tag{22}$$

where:

$$I_{R}^{k} \equiv \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{x}\}_{n}} \left| 2\varepsilon \frac{\partial \delta}{\partial \theta_{k}} \right|^{2} d\boldsymbol{u}$$
(23)

$$I_{L} \equiv \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{X}\}_{n}} d\boldsymbol{u}$$
(24)

Proof: From the mixed information definition:

$$i_F\left(f_{\{\mathbf{X}\}_n}, f_{\{\mathbf{Y}\}_n|\boldsymbol{\theta}}\right)_k = I^k_{\mathcal{M}^n_{\varepsilon}} + I^k_{\sim \mathcal{M}^n_{\varepsilon}} \ge I^k_{\mathcal{M}^n_{\varepsilon}} > 0$$

$$\tag{25}$$

where:

$$I_{\mathcal{M}_{\varepsilon}^{n}}^{k} \equiv \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{x}\}_{n}} \left(\frac{\partial \ln f_{\{\mathbf{Y}\}_{n}|\boldsymbol{\theta}}}{\partial \theta_{k}}\right)^{2} d\boldsymbol{u}$$
(26)

and the symbol, \sim , denotes the complement set.

Replacing Equation (16) in Equation (26), it is obtained:

$$I_{\mathcal{M}_{\varepsilon}^{n}}^{k} = n^{2} \left(\frac{\partial h_{M} \left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}} \right)}{\partial \theta_{k}} \right)^{2} \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{X}\}_{n}} d\boldsymbol{u} - 2n \frac{\partial h_{M} \left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}} \right)}{\partial \theta_{k}} \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{X}\}_{n}} \left(2\varepsilon n \frac{\partial \delta}{\partial \theta_{k}} \right) d\boldsymbol{u} + \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{X}\}_{n}} \left(2\varepsilon n \frac{\partial \delta}{\partial \theta_{k}} \right)^{2} d\boldsymbol{u}$$

$$(27)$$

Using the Cauchy-Bunyakovsky-Schwarz inequality:

$$\left| \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{x}\}_{n}} \left(2\varepsilon n \frac{\partial \delta}{\partial \theta_{k}} \right) d\mathbf{u} \right|^{2} = \left| \int_{\mathcal{M}_{\varepsilon}^{n}} \sqrt{f_{\{\mathbf{x}\}_{n}}} \sqrt{f_{\{\mathbf{x}\}_{n}}} \left(2\varepsilon n \frac{\partial \delta}{\partial \theta_{k}} \right) d\mathbf{u} \right|^{2}$$

$$\leq \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{x}\}_{n}} d\mathbf{u} \cdot \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{x}\}_{n}} \left| 2\varepsilon n \frac{\partial \delta}{\partial \theta_{k}} \right|^{2} d\mathbf{u}$$

$$= I_{L} \cdot n^{2} I_{R}^{k}$$
(28)

Therefore:

$$I_{\mathcal{M}_{\varepsilon}^{n}}^{k} \geq n^{2} \left(\frac{\partial h_{M} \left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}} \right)}{\partial \theta_{k}} \right)^{2} I_{L} - 2n^{2} \left| \frac{\partial h_{M} \left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}} \right)}{\partial \theta_{k}} \right| \sqrt{I_{L}} \sqrt{I_{R}^{k}} + n^{2} I_{R}^{k}$$
$$= n^{2} \left(\left| \frac{\partial h_{M} \left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}} \right)}{\partial \theta_{k}} \right| \sqrt{I_{L}} - \sqrt{I_{R}^{k}} \right)^{2}$$
(29)

which implies:

$$\frac{1}{n^2} i_F \left(f_{\{\mathbf{X}\}_n}, f_{\{\mathbf{Y}\}_n \mid \boldsymbol{\theta}} \right)_k \ge \left(\left| \frac{\partial h_M \left(f_{\mathbf{X}}, f_{\mathbf{Y} \mid \boldsymbol{\theta}} \right)}{\partial \theta_k} \right| \sqrt{I_L} - \sqrt{I_R^k} \right)^2$$
(30)

Hence, using the negative solution of the square root:

$$\frac{1}{n}\sqrt{i_{F}\left(f_{\{\boldsymbol{X}\}_{n}},f_{\{\boldsymbol{Y}\}_{n}|\boldsymbol{\theta}}\right)_{k}} \geq -\left|\frac{\partial h_{M}\left(f_{\boldsymbol{X}},f_{\boldsymbol{Y}|\boldsymbol{\theta}}\right)}{\partial\theta_{k}}\right|\sqrt{I_{L}}+\sqrt{I_{R}^{k}}$$
(31)

Therefore:

$$\frac{\sqrt{I_{R}^{k}} - \frac{1}{n}\sqrt{i_{F}\left(f_{\{\mathbf{X}\}_{n}}, f_{\{\mathbf{Y}\}_{n}|\boldsymbol{\theta}}\right)_{k}}}{\sqrt{I_{L}}} \leq \left|\frac{\partial h_{M}\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}}\right)}{\partial \theta_{k}}\right|$$
(32)

Taking into account the positive solution of Equation (30):

$$\left|\frac{\partial h_M\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}}\right)}{\partial \theta_k}\right| \le \frac{\sqrt{I_R^k} + \frac{1}{n}\sqrt{i_F\left(f_{\{\mathbf{X}\}_n}, f_{\{\mathbf{Y}\}_n|\boldsymbol{\theta}}\right)_k}}{\sqrt{I_L}}$$
(33)

The previous result allows us to state the following remarks:

Remark 1 This theorem states algebraic bounds on the derivative of the mixed entropy that are valid for any value of *n*.

Remark 2 The expression:

$$I_{R}^{k} \equiv \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{X}\}_{n}} \left| 2\varepsilon \frac{\partial \delta}{\partial \theta_{k}} \right|^{2} d\boldsymbol{u} = \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{X}\}_{n}} \left| -\frac{\partial}{\partial \theta_{k}} \left(-\frac{1}{n} \ln f_{\{\mathbf{Y}\}_{n}|\boldsymbol{\theta}} - h_{M} \left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}} \right) \right) \right|^{2} d\boldsymbol{u}$$
(34)

refers to the accumulation of derivatives of very small differences. This differences become smaller and smaller as the weak law of large numbers is enacted.

8. Unknown Source Density Function, f_x

Now, we study what happens when the source density function, f_x , is not known.

It can be proven that $D_{KL}(f_X || f_{Y|\theta}) \ge 0$ ([2][p. 232]). Thus, from the corresponding definitions given at the start of this work:

$$h_M\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}}\right) \ge h_S(f_{\mathbf{X}}) \tag{35}$$

Therefore, the density estimation problem can be framed as the following optimization program:

$$\boldsymbol{\theta}_{\circ} \equiv \boldsymbol{\theta}_{\circ}(f_{\boldsymbol{X}}, f_{\boldsymbol{Y}|\boldsymbol{\theta}}) = \arg\min_{\boldsymbol{\theta}} h_{M}\left(f_{\boldsymbol{X}}, f_{\boldsymbol{Y}|\boldsymbol{\theta}}\right)$$
(36)

Remark 3 The solutions of this optimization program pose the following options:

- $h_M(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}_0}) = h_S(f_{\mathbf{X}})$, hence $f_{\mathbf{Y}|\boldsymbol{\theta}_0} = f_{\mathbf{X}}$.
- $h_M(f_{\mathbf{x}}, f_{\mathbf{y}|\theta_o}) > h_s(f_{\mathbf{x}})$, thus $f_{\mathbf{y}|\theta_o} \neq f_{\mathbf{x}}$. This happens when the estimator cannot implement the desired density function, which is beyond its reach. One way of checking whether the global minimum has been reached is to compare the mixed entropy value to that of the Shannon differential entropy. If they differ, the global minimum has not been reached yet.

According to this result, it is straightforward to search for parameters that make true the following expression:

$$\nabla_{\boldsymbol{\theta}} h_M \left(f_{\boldsymbol{X}}, f_{\boldsymbol{Y}|\boldsymbol{\theta}} \right) \Big|_{\boldsymbol{\theta} = \boldsymbol{\theta}_{\circ}} = \boldsymbol{0}$$
(37)

However, this is out of the limits in our density estimation problem, given that we do not have access to f_x . Thank to the weak law of large numbers:

$$\lim_{n \to \infty} -\frac{1}{n} \ln f_{\{\mathbf{Y}\}_n \mid \theta} = h_M \left(f_{\mathbf{X}}, f_{\mathbf{Y} \mid \theta} \right)$$
(38)

in probability. Thus, for sufficiently large values of n:

$$-\frac{1}{n}\ln f_{\{\mathbf{Y}\}_n|\boldsymbol{\theta}} \approx h_M\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}}\right)$$
(39)

Hence, the optimization program posed in Equation (36) can be approximated by:

$$\boldsymbol{\theta}_{\circ} \equiv \boldsymbol{\theta}_{\circ}(n, \{\boldsymbol{x}\}_{n}, f_{\boldsymbol{Y}|\boldsymbol{\theta}}) = \arg\min_{\boldsymbol{\theta}} \left(-\frac{1}{n} \ln f_{\{\boldsymbol{Y}\}_{n}|\boldsymbol{\theta}} \right)$$
(40)

or:

$$\boldsymbol{\theta}_{\circ} \equiv \boldsymbol{\theta}_{\circ}(n, \{\boldsymbol{x}\}_{n}, f_{\boldsymbol{Y}|\boldsymbol{\theta}}) = \arg \max_{\boldsymbol{\theta}} \left(\frac{1}{n} \ln f_{\{\boldsymbol{Y}\}_{n}|\boldsymbol{\theta}}\right)$$
(41)

which is the *maximum log likelihood* framework ([3][p. 261], [4][260], [5][p. 65]).

Remark 4 As in the mixed entropy case, the maximum log likelihood framework can exhibit the following cases:

• The optimization program finds a parameter combination that makes the estimator equal to the unknown source density function, thus making true the following expression for a sufficiently large *n*:

$$-\frac{1}{n}\ln f_{\{\mathbf{Y}\}_n|\boldsymbol{\theta}_{\mathbf{o}}} \approx h_M\left(f_{\mathbf{X}}, f_{\mathbf{Y}|\boldsymbol{\theta}_{\mathbf{o}}}\right) \approx h_S(f_{\mathbf{X}})$$
(42)

• The estimator does not find the desired target function. This happens when the family of functions that can be implemented by the estimator does not include the target function. It also happens when the maximum log likelihood program is riddled with multiple local minima, and the optimization process gets trapped in one of them. Only under very specific conditions, it is possible to obtain the global maximum of the maximum log likelihood problem [10]. In this case, we cannot use a comparison between the value of the log likelihood term and that of the Shannon differential entropy to determine the goodness of our solution, as we could do in the case of the mixed entropy, because calculating the latter is analogous to the very problem we are trying to solve.

Hence, it is natural to work with this approximated program and look for parameters that make the following expression true:

$$\nabla_{\boldsymbol{\theta}} \left(\frac{1}{n} \ln f_{\{\boldsymbol{Y}\}_n | \boldsymbol{\theta}} \right) \bigg|_{\boldsymbol{\theta} = \boldsymbol{\theta}_{\circ}} = \mathbf{0}$$
(43)

9. An Example: Mismatched Density Functions

9.1. Known Density Functions Case

Let us define a source density function:

$$f_X(x) = \lambda e^{-\lambda x} H(x) \tag{44}$$

with $x \in \mathbb{R}$ and H(x), the Heaviside step function. This density function is completely defined by the parameter, λ . Its differential entropy is given by:

$$h_s(f_x) = 1 - \ln \lambda = 1 + \ln \frac{1}{\lambda}$$
(45)

We also have a data set, $\{x\}_n$, generated by the source density function.

The density function that depends on parameters is:

$$f_{Y|\mu,\sigma} = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(\mu-x)^2}{2\sigma^2}\right)$$
(46)

with $x \in \mathbb{R}$, μ its mean and σ its standard deviation. Its differential entropy is given by:

$$h_{S}(f_{Y|\mu,\sigma}) = \ln\left(\sigma\sqrt{2\pi e}\right) = 1 + \ln\sigma + \ln\sqrt{\frac{2\pi}{e}}$$
(47)

which does not depend on the mean, μ .

9.2. Compliance with Theorem 2

We first analyze the bounds associated with parameter μ . For these density functions:

$$\frac{1}{n}\ln f_{\{Y\}_n|\mu,\sigma} = \frac{1}{n}\sum_{k=1}^n \ln f_{Y|\mu,\sigma}(x_k) = \ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{1}{2n\sigma^2}\sum_{k=1}^n (\mu - x_k)^2$$
(48)

Hence:

$$\frac{1}{n}\frac{\partial \ln f_{\{Y\}_n|\mu,\sigma}}{\partial \mu} = -\frac{1}{\sigma^2} \left(\mu - \frac{1}{n}\sum_{k=1}^n x_k\right)$$
(49)

We calculate the mixed information for μ :

$$i_{F} \left(f_{\{X\}_{n}}, f_{\{Y\}_{n}|\mu,\sigma} \right)_{\mu} = \int_{\mathbb{R}^{n}_{+}} f_{\{X\}_{n}} \left(\frac{\partial \ln f_{\{Y\}_{n}|\mu,\sigma}}{\partial \mu} \right)^{2} d\boldsymbol{u}$$
$$= \int_{\mathbb{R}^{n}_{+}} f_{\{X\}_{n}} \left(\frac{n}{\sigma^{2}} \left(\mu - \frac{1}{n} \sum_{k=1}^{n} u_{k} \right) \right)^{2} d\boldsymbol{u}$$
$$= \frac{n^{2}}{\sigma^{4}} \int_{\mathbb{R}^{n}_{+}} f_{\{X\}_{n}} \left(\mu - \frac{1}{n} \sum_{k=1}^{n} u_{k} \right)^{2} d\boldsymbol{u}$$
(50)

Hence:

$$\frac{1}{n}\sqrt{i_F\left(f_{\{X\}_n}, f_{\{Y\}_n \mid \mu, \sigma}\right)_{\mu}} = \frac{1}{\sigma^2}\sqrt{\int_{\mathbb{R}^n_+} f_{\{X\}_n}\left(\mu - \frac{1}{n}\sum_{k=1}^n u_k\right)^2} d\boldsymbol{u}$$
(51)

Г

Also:

$$h_{M}(f_{X}, f_{Y|\mu,\sigma}) = -\int_{-\infty}^{\infty} f_{X}(u) \ln f_{Y|\mu,\sigma}(u) du$$

$$= -\int_{-\infty}^{\infty} \lambda e^{-\lambda u} H(u) \ln f_{Y|\mu,\sigma}(u) du$$

$$= -\int_{0}^{\infty} \lambda e^{-\lambda u} \left(\ln \left(\frac{1}{\sigma\sqrt{2\pi}} \right) - \frac{(\mu - u)^{2}}{2\sigma^{2}} \right) du$$

$$= \ln \left(\sigma\sqrt{2\pi} \right) \int_{0}^{\infty} \lambda e^{-\lambda u} du + \frac{\lambda}{2\sigma^{2}} \int_{0}^{\infty} (\mu - u)^{2} e^{-\lambda u} du$$

$$= \ln \left(\sigma\sqrt{2\pi} \right) + \frac{1}{\sigma^{2}} \left(\frac{\mu^{2}}{2} - \frac{\mu}{\lambda} + \frac{1}{\lambda^{2}} \right)$$
(52)

and:

$$\frac{\partial h_M\left(f_X, f_{Y|\mu,\sigma}\right)}{\partial \mu} = \frac{1}{\sigma^2} \left(\mu - \frac{1}{\lambda}\right) \tag{53}$$

Finally, we calculate the micro-structure term using Equation (17):

$$2\varepsilon \frac{\partial \delta}{\partial \mu} = \frac{1}{\sigma^2} \left(\frac{1}{n} \sum_{k=1}^n x_k - \frac{1}{\lambda} \right)$$
(54)

This allows us to calculate:

$$I_{R}^{\mu} = \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{x}\}_{n}} \left| 2\varepsilon \frac{\partial \delta}{\partial \mu} \right|^{2} d\boldsymbol{u}$$

$$= \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{x}\}_{n}} \left| \frac{1}{\sigma^{2}} \left(\frac{1}{n} \sum_{k=1}^{n} u_{k} - \frac{1}{\lambda} \right) \right|^{2} d\boldsymbol{u}$$

$$= \frac{1}{\sigma^{4}} \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{x}\}_{n}} \left(\frac{1}{n} \sum_{k=1}^{n} u_{k} - \frac{1}{\lambda} \right)^{2} d\boldsymbol{u}$$
(55)

Now, for σ , we have:

$$\frac{1}{n}\frac{\partial \ln f_{\{Y\}n\mid\mu,\sigma}}{\partial\sigma} = \frac{1}{\sigma^3} \left(-\sigma^2 + \frac{1}{n}\sum_{k=1}^n \left(\mu - u_k\right)^2 \right)$$
(56)

The mixed information for σ is defined by:

$$i_{F} \left(f_{\{X\}_{n}}, f_{\{Y\}_{n}|\mu,\sigma} \right)_{\mu} = \int_{\mathbb{R}^{n}_{+}} f_{\{X\}_{n}} \left(\frac{\partial \ln f_{\{Y\}_{n}|\mu,\sigma}}{\partial \sigma} \right)^{2} d\boldsymbol{u} \\ = \int_{\mathbb{R}^{n}_{+}} f_{\{X\}_{n}} \left(\frac{n}{\sigma^{3}} \left(-\sigma^{2} + \frac{1}{n} \sum_{k=1}^{n} (\mu - u_{k})^{2} \right) \right)^{2} d\boldsymbol{u} \\ = \frac{n^{2}}{\sigma^{6}} \int_{\mathbb{R}^{n}_{+}} f_{\{X\}_{n}} \left(-\sigma^{2} + \frac{1}{n} \sum_{k=1}^{n} (\mu - u_{k})^{2} \right)^{2} d\boldsymbol{u}$$
(57)

Hence:

$$\frac{1}{n}\sqrt{i_F\left(f_{\{X\}_n}, f_{\{Y\}_n \mid \mu, \sigma}\right)_{\mu}} = \frac{1}{\sigma^3}\sqrt{\int_{\mathbb{R}^n_+} f_{\{X\}_n}\left(-\sigma^2 + \frac{1}{n}\sum_{k=1}^n \left(\mu - u_k\right)^2\right)^2} d\boldsymbol{u}$$
(58)

Again, the derivative of the mixed entropy with respect to σ :

$$\frac{\partial h_M\left(f_X, f_{Y|\mu,\sigma}\right)}{\partial \sigma} = \frac{1}{\sigma^3} \left(\left(\sigma^2 - \frac{1}{\lambda^2}\right) - \left(\mu - \frac{1}{\lambda}\right)^2 \right)$$
(59)

Hence, the associated micro-differences expression is:

$$2\varepsilon \frac{\partial \delta}{\partial \sigma} = \frac{1}{\sigma^3} \left(\left(\frac{1}{n} \sum_{k=1}^n \left(\mu - u_k \right)^2 - \frac{1}{\lambda^2} \right) - \left(\mu - \frac{1}{\lambda} \right)^2 \right)$$
(60)

which allows us to calculate:

$$I_{R}^{\sigma} = \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{x}\}_{n}} \left| 2\varepsilon \frac{\partial \delta}{\partial \sigma} \right|^{2} d\boldsymbol{u}$$

$$= \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{x}\}_{n}} \left| \frac{1}{\sigma^{3}} \left(\left(\frac{1}{n} \sum_{k=1}^{n} \left(\mu - u_{k} \right)^{2} - \frac{1}{\lambda^{2}} \right) - \left(\mu - \frac{1}{\lambda} \right)^{2} \right) \right|^{2} d\boldsymbol{u}$$

$$= \frac{1}{\sigma^{6}} \int_{\mathcal{M}_{\varepsilon}^{n}} f_{\{\mathbf{x}\}_{n}} \left(\left(\frac{1}{n} \sum_{k=1}^{n} \left(\mu - u_{k} \right)^{2} - \frac{1}{\lambda^{2}} \right) - \left(\mu - \frac{1}{\lambda} \right)^{2} \right)^{2} d\boldsymbol{u}$$
(61)

9.3. Unknown Source Density Function, f_X

If f_x is unknown, then it is not possible to calculate the previous bounds. In this case, we resort to the optimization program specified in Equation (41). Hence, we use Equations (49) and (56), in order to obtain:

$$\mu_n = \frac{1}{n} \sum_{k=1}^n x_k \tag{62}$$

$$\sigma_n^2 = \frac{1}{n} \sum_{k=1}^n (\mu_n - x_k)^2$$
(63)

Using this choice of values immediately makes both mixed information expressions described by Equations (51) and (58) equal to zero, without the need for large values of n.

Furthermore, if we use:

$$\mu_{\circ} = \lim_{n \to \infty} \mu_n = \frac{1}{\lambda} \tag{64}$$

$$\sigma_{\circ} = \lim_{n \to \infty} \sigma_n = \frac{1}{\lambda}$$
(65)

then the micro-differences integrals also become equal to zero.

This example shows several things:

- The use of the maximum log likelihood solutions and large values of n guarantees that the derivative of the mixed entropy equals zero; hence, the derivative of the relative entropy with respect to its parameters is zero, too.
- However, this is not enough to guarantee that the density functions will match each other (a Gaussian cannot estimate an exponential). The maximum log likelihood solution only guarantees that the estimator will do the best it can do.
- If not enough data is available, the derivative of the mixed entropy will be different from zero, because the micro-differences terms are not zero.

10. Discussion

If one wanted to estimate a density function out of a data set, one could minimize the relative entropy: once we reach its minimum, we can guarantee that the density function implemented by the estimator will be equal to that which originated the data. This is true only if the family of functions that the estimator can effectively implement includes that of the source density function. Moreover, realizing that it is not possible to measure the relative entropy, one would quickly settle to maximize the log likelihood and obtain the desired density function. Keeping this in mind, the results presented in this work seem just another way of saying the same. However, this is not so. The difference is subtle. Whereas there are many density function estimation principles and it could be discussed whether minimizing the relative entropy is the most convenient one, the bounds presented in this work show that minimizing the relative entropy automatically activates in the presence of large data sets and that one does not need to decide

whether that estimation framework is the most convenient or not. Something similar happens with the weak law of large numbers: there are many ways of determining the expected value of the density function that generated the data, but everybody knows that the weak law of large numbers guarantees that in the presence of large data sets, this value can be effectively approximated by the average of the values. Our bounds, also based on the weak law of large numbers, state something similar: in the presence of large data sets, the source density function is perceived as the solution of the maximum log likelihood solution.

11. Conclusions

Theorem 2, which is the main result of this work, consists in a set of bounds on the partial derivative of the mixed entropy, which is a proxy to the relative entropy, with respect to the parameters of the estimator. These bounds relate the mixed information value, the capacity of the estimator to produce micro-differences and the number of random variables present in the analysis in such a way that it is possible to determine that these bounds activate automatically when the data set is large. If the micro-differences integral is zero, the mixed information derivative is zero for large data sets, independently of any other considerations. In other words, minimizing the mixed entropy is the preferred density estimation framework when large data sets are available.

Given the impossibility of estimating the mixed entropy, it is shown that the correct numerical approximation to this problem is finding the solution to the maximum log likelihood problem. Again, when the estimator is associated to a micro-differences integral equal to zero, this framework becomes the optimal one in the presence of large data sets.

Acknowledgments

The authors thank Jaime Cisternas, Jorge Silva and Jose Principe for their helpful insights about this work.

Conflict of Interest

The authors declare no conflict of interest.

References

- 1. Shannon, C. A mathematical theory of communication. Bell Syst. Tech. J. 1948, 27, 379-423.
- 2. Cover, T.; Thomas, J. *Elements of Information Theory*; John Wiley and Sons, Inc.: Hoboken, NJ, USA, 1991.
- Hogg, R.V.; Craig, A.T. Introduction to Mathematical Statistics; Prentice Hall: Upper Saddle River, NJ, USA, 1995.
- 4. Papoulis, A. *Probability, Random Variables, and Stochastic Processes*; McGraw-Hill: New York, NY, USA, 1991.
- 5. Van Trees, H.L. *Detection, Estimation, and Modulation Theory: Part 1*; John Wiley and Sons, Inc.: Hoboken, NJ, USA, 2001.

- 7. Romera, E.; Dehesa, J. The Fisher-Shannon information plane, an electron correlation tool. *J. Chem. Phys.* **2004**, *120*, 8906–8912.
- 8. Dimitrov, V. On Shannon-Jaynes Entropy and Fisher Information. In Proceedings of the 27th International Workshop on Bayesian Inference and Maximum Entropy Methods in Science and Engineering, Saratoga Springs, New York, NY, USA, 8–13 July 2007.
- 9. Taubman, D.; Marcellin, M. *JPEG2000: Image Compression Fundamentals, Standards, and Practice*; Kluwer Academic Publishers: Dordrecht, The Netherlands, 2002.
- 10. Boyd, S.; Vandenberghe, L. *Convex Optimization*; Cambridge University Press: Cambridge, UK, 2004.

© 2013 by the authors; licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution license (http://creativecommons.org/licenses/by/3.0/).