

TUTORIAL HANDOUT GENOVA JUNE 11

INFORMATION GEOMETRY AND ALGEBRAIC STATISTICS ON A FINITE STATE SPACE AND ON GAUSSIAN MODELS

LUIGI MALAGÒ AND GIOVANNI PISTONE

CONTENTS

1. Introduction	2
1.1. C. R. Rao	2
1.2. S.-H. Amari	2
2. Finite State Space: Full Simplex	4
2.1. Statistical bundle	4
2.2. Example: fractionalization	7
2.3. Example: entropy	8
2.4. Example: the polarization measure	10
2.5. Transports	10
2.6. Connections	11
2.7. Atlases	12
2.8. Using parameters	13
3. Finite State Space: Exponential Families	16
3.1. Statistical manifold	17
3.2. Gradient	18
3.3. Gradient flow in the mixture geometry	20
3.4. Gradient flow of the expected value function	20
4. Gaussian Models	21
4.1. Gaussian model in the Hermite basis	21
4.2. Optimisation	23
References	25

Date: June 11, 2015.

The authors wish to thank Henry Wynn for his comments on a draft of this tutorial. Some material reproduced here is part of work in progress by the authors and collaborators. G. Pistone is supported by the de Castro Statistics Initiative, Collegio Carlo Alberto, Moncalieri; and he is a member of GNAMPA-INDAM.

1. INTRODUCTION

It was shown by C. R. Rao in a paper published 1945 [23] that the set of positive probabilities $\Delta^\circ(\Omega)$ on a finite state space Ω is a *Riemannian manifold*, as it is defined in classical treatises such as [7] and [11], but in a way which is of interest for Statistics. It was later pointed out by Sun-Ichi Amari that it is actually possible to define two affine geometries of Hessian type [29] on top of the Rao's Riemannian geometry, but see also the original contribution by Steffen Lauritzen [12]. This development was somehow stimulated by two papers by Efron [8, 9]. The original work of Amari was published in the '80s, see Shun'ichi Amari [1], see monograph presentations in [10] and [3]. Amari gave to this new topic the name of *Information Geometry* and provided many applications, in particular in Machine Learning [2].

Information Geometry and Algebraic statistics are deeply connected because of the central place occupied by statistical exponential families [4] in both fields. There is possibly a simpler connection, which is the object of the first part of this presentation. The present tutorial is focused mainly on Differential Geometry.

The present tutorial treats only cases where the statistical model is *parametric*. However, there is an effort to use methods that scale well to the case where the statistical model is essentially infinite dimensional, e.g. [5, 6], *parry—dawid—lauritzen:2012*, [14], and, in general, all applications in Statistical Physics.

Where to start from? Here is my choice, but read the comments by C. R. Rao to Scholarpedia.

1.1. **C. R. Rao.** In [23] we find the following computation:

$$\begin{aligned}
 \frac{d}{dt} \mathbb{E}_{p(t)} [U] &= \frac{d}{dt} \sum_x U(x) p(x; t) \\
 &= \sum_x U(x) \frac{d}{dt} p(x; t) \\
 &= \sum_x U(x) \frac{d}{dt} \log(p(x; t)) p(x; t) \\
 &= \sum_x (U(x) - \mathbb{E}_{p(t)} [U]) \frac{d}{dt} \log(p(x; t)) p(x; t) \\
 &= \mathbb{E}_{p(t)} \left[(U - \mathbb{E}_{p(t)} [U]) \frac{d}{dt} \log(p(t)) \right] \\
 &= \left\langle (U - \mathbb{E}_{p(t)} [U]), \frac{d}{dt} \log(p(t)) \right\rangle_{p(t)}.
 \end{aligned}$$

Here the relevant point is fact the scalar product is computed at the running $p(t)$ and the Fisher's score $\frac{d}{dt} \log(p(t))$ appears as a measure of velocity.

1.2. **S.-H. Amari.** In [2] there are applications of computations of the following type. Give a function $f: \Delta^\circ(\Omega) \rightarrow \mathbb{R}$,

$$\begin{aligned}
\frac{d}{dt}f(p(t)) &= \sum_{x \in \Omega} \frac{\partial}{\partial p(x)} f(p(x; t): x \in \Omega) \frac{d}{dt}p(x; t) \\
&= \left\langle \text{grad } f(p(t)), \frac{d}{dt} \log(p(x; t)) \right\rangle_{p(t)} \\
&= \left\langle \text{grad } f(p(t)) - \mathbb{E}_{p(t)} [\text{grad } f(p(t))], \frac{d}{dt} \log(p(x; t)) \right\rangle_{p(t)},
\end{aligned}$$

where $\text{grad } f(p(t)) - \mathbb{E}_{p(t)} [\text{grad } f(p(t))]$ appears as the gradient of f in the scalar product $\langle \cdot, \cdot \rangle_{p(t)}$.

2. FINITE STATE SPACE: FULL SIMPLEX

Let Ω be a finite set with $n + 1 = \#\Omega$ points. We denote by $\Delta(\Omega)$ the simplex of the probability functions $p: \Omega \rightarrow \mathbb{R}_{\geq 0}$, $\sum_{x \in \Omega} p(x) = 1$. It is a n -simplex, i.e. an n -dimensional polytope which is the convex hull of its $n + 1$ vertices. It is a closed and convex subset of the affine space $A_1(\Omega) = \{q \in \mathbb{R}^\Omega \mid \sum_{x \in \Omega} q(x) = 1\}$ and it has non empty relative topological interior.

The interior of the simplex, Δ_n° , is the set of the strictly positive probability functions,

$$\Delta^\circ(\Omega) = \left\{ p \in \mathbb{R}^\Omega \mid \sum_{x \in \Omega} p(x) = 1, p(x) > 0 \right\}.$$

The border of the simplex is the union of all the faces of $\Delta(\Omega)$ as a convex set. We recall that a face of maximal dimension $n - 1$ is called *facet*. Each face is itself a simplex. An *edge* is a face of dimension 1.

Remark 1. The presentation below does not use explicitly any specific parameterization of the sets $\Delta^\circ(\Omega)$, $\Delta(\Omega)$, $A_1(\Omega)$. However, the actual extension of this theory to non finite sample space requires a careful handling as most of the topological features do not hold in such a case. One possibility is given by the so called exponential manifolds, see [21].

2.1. Statistical bundle. We first discuss the statistical geometry on the open simplex as deriving from a *vector bundle* with base $\Delta^\circ(\Omega)$. *The notion of vector bundle has been introduced in non-parametric Information Geometry by [13].* Later we will show that such a bundle can be identified with the tangent bundle of proper manifold structure. It is nevertheless interesting to observe that a number of geometrical properties do not require the actual definition of the statistical manifold, possibly opening the way to a generalization.

For each $p \in \Delta^\circ(\Omega)$ we consider the plane through the origin, orthogonal to the vector \vec{Op} . The set of positive probabilities each one associated to its plane forms a vector bundle which is the basic structure of our presentation of Information Geometry, see Fig. 1. Note that, because of our orientation to Statistics, we call each element of \mathbb{R}^Ω a *random variable*. A section mapping S from probabilities $p \in \Delta^\circ(\Omega)$ to the bundle, $\mathbb{E}_p[S(p)] = 0$ is an *estimating function* as the equation $F(\hat{p}, x) = 0$, $x \in \Omega$, provides an *estimator* that is a distinguished mapping from the sample space Ω to the simplex of probabilities $\Delta(\Omega)$.

We can give a formal definition as follows.

Definition 1.

- (1) For each $p \in \Delta^\circ(\Omega)$ let B_p be the vector space of random variables U that are p -centered,

$$B_p = \left\{ U: \Omega \rightarrow \mathbb{R} \mid \mathbb{E}_p[U] = \sum_{x \in \Omega} U(x)p(x) = 0 \right\}.$$

Each B_p is an Hilbert space for the scalar product $\langle U, V \rangle_p = \mathbb{E}_p[UV]$.

- (2) The statistical bundle of the open simplex $\Delta^\circ(\Omega)$ is the linear bundle on $\Delta^\circ(\Omega)$

$$T\Delta^\circ(\Omega) = \{(p, U) \mid p \in \Delta^\circ(\Omega), U \in B_p\}.$$

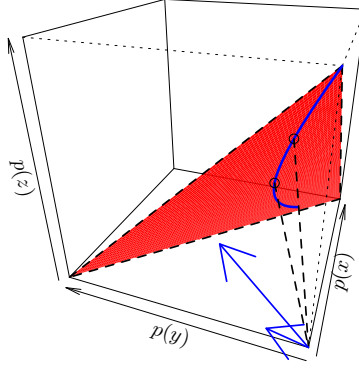


FIGURE 1. The red triangle is the simplex on the sample space with 3 points $\Omega = \{x, y, z\}$ viewed from below. The blue curve on the simplex is a one-dimensional statistical model. The probabilities p are represented by vectors from O to the point whose coordinates are $p = (p(x), p(y), p(z))$. The velocity vectors $Dp(t)$ of a curve $I \mapsto p(t)$ are represented by arrows; they are orthogonal to the vector from O to p .

It is an open subset of the variety of $\mathbb{R}^\Omega \times \mathbb{R}^\Omega$ defined by the polynomial equations

$$\begin{cases} \sum_{x \in \Omega} p(x) = 1, \\ \sum_{x \in \Omega} U(x)p(x) = 0. \end{cases}$$

(3) A vector field F of the statistical bundle is a section of the bundle i.e.,

$$F: \Delta^\circ(\Omega) \ni p \mapsto (p, F(p)) \in T\Delta^\circ(\Omega)$$

The term estimating function is also used in the statistical literature.

(4) If $I \ni t \mapsto p(t) \in \Delta^\circ(\Omega)$ is a C^1 curve, its score is defined by

$$Dp(t) = \frac{\dot{p}(t)}{p(t)} = \frac{d}{dt} \log p(t), \quad t \in I.$$

As the score $Dp(t)$ is a $p(t)$ -centered random variable which belongs to $B_{p(t)}$ for all $t \in I$, the mapping $I \ni t \mapsto (p(t), Dp(t))$ is a curve in the statistical bundle. Note that the notion of score extends to any curve $p(\cdot)$ in the affine space $A_1(\Omega)$ by its relation to the statistical gradient, i.e. $Dp(t)$ is implicitly defined by

$$\langle \text{grad } f(p(t)) - \mathbb{E}_{p(t)} [\text{grad } f(p(t))], \cdot \rangle_{p(t)} \quad f \in C^1(\mathbb{R}^\Omega).$$

(5) Given a function $f: \Delta^\circ(\Omega) \rightarrow \mathbb{R}$, its statistical gradient is a vector field $\nabla f: \Delta^\circ(\Omega) \ni p \mapsto (p, \nabla F(p)) \in T\Delta^\circ(\Omega)$ such that for each differentiable curve $t \mapsto p(t)$ it holds

$$\frac{d}{dt} f(p(t)) = \langle \nabla f(p(t)), Dp(t) \rangle_{p(t)}$$

Remark 2. The Information Geometry on the simplex does not coincide with the geometry of the embedding of the simplex $\Delta^\circ(\Omega) \rightarrow \mathbb{R}^\Omega$, in the sense the statistical bundle is not the tangent bundle of these embedding, see Fig. 1. It will become the tangent bundle of the proper geometric structure which is given by special atlases.

Remark 3. We could extend the statistical bundle by taking the linear fiberts on $\Delta(\Omega)$ or $A_1(\Omega)$. In such cases the bilinear form is not always a scalar product. In fact $\langle \cdot, \cdot \rangle_p$ is not faithful where at least a component of the probability vector is zero, while it is not positive definite outside the simplex $\Delta(\Omega)$.

Remark 4. The vector $Dp(t) \in B_{p(t)}$ is meant to represent the relative variation of the information in a one dimensional statistical model. The score is a representation of the velocity along a curve, because of the geometric interpretation of C. R. Rao's already mentioned classical computation of [23],

$$\begin{aligned} \frac{d}{dt} \mathbb{E}_{p(t)} [U] &= \frac{d}{dt} \sum_x U(x) p(x; t) = \sum_x U(x) \frac{d}{dt} p(x; t) = \\ & \sum_x U(x) \frac{d}{dt} \log(p(x; t)) p(x; t) = \sum_x (U(x) - \mathbb{E}_{p(t)} [U]) \frac{d}{dt} \log(p(x; t)) p(x; t) = \\ & E_t \left[(U - \mathbb{E}_{p(t)} [U]) \frac{d}{dt} \log(p(t)) \right] = \langle U - \mathbb{E}_{p(t)} [U], Dp(t) \rangle_{p(t)} \end{aligned}$$

We observe that the scalar product above is the scalar product on $B_{p(t)}$ because $U \mapsto U - \mathbb{E}_p [U]$ takes values in B_p .

Remark 5. Consider the level surface of f at $p_0 \in \Delta^\circ(\Omega)$, that is $\{p \in \Delta^\circ(\Omega) | f(p) = f(p_0)\}$, and assume $\nabla f(p_0) \neq 0$. Then for each curve through p_0 , $I \mapsto p(t)$ with $p(0) = p_0$, such that $f(p(t)) = f(p_0)$, we have

$$0 = \frac{d}{dt} f(p(t)) \Big|_{t=0} = \langle \nabla f(p(t)), Dp(t) \rangle_{p(t)} \Big|_{t=0} = \langle \nabla f(p_0), Dp(t_0) \rangle_{p_0}$$

that is *all velocities $Dp(t_0)$ tangential to the level set are orthogonal to the statistical gradient*. Note that we have not jet defined a manifold such that the statistical bundle is equal to the tangent bundle.

Remark 6. If the function $f: \Delta^\circ(\Omega) \rightarrow \mathbb{R}$ extends to a C^1 function on an open subset of \mathbb{R}^Ω , then we can compute the statistical gradient via the ordinary gradient in the geometry of \mathbb{R}^Ω , namely $\text{grad } f(p) = \left(\frac{\partial}{\partial p(x)} f(p) : x \in \Omega \right)$. In fact,

$$\begin{aligned} \frac{d}{dt} f(p(t)) &= \sum_{x \in \Omega} \frac{\partial}{\partial p(x)} f(p(x; t)) \frac{d}{dt} p(x; t) = \\ & \langle \text{grad } f(p(t)), Dp(x; t) \rangle_{p(t)} = \langle \text{grad } f(p(t)) - \mathbb{E}_{p(t)} [\text{grad } f(p(t))], Dp(t) \rangle_{p(t)} = \\ & \langle \nabla f(p(t)), Dp(t) \rangle_{p(t)} , \end{aligned}$$

where $\text{grad } f(p(t)) - \mathbb{E}_{p(t)} [\text{grad } f(p(t))]$ is the projection of $\text{grad } f(p(t))$ onto B_p with respect to the scalar product $\langle \cdot, \cdot \rangle_{p(t)}$. Note that the statistical gradient is zero if, and only if, the ordinary gradient is constant.

Definition 2 (Flow).

- (1) *Given a vector field F , the trajectories along the vector field are the solution of the differential equation*

$$\frac{D}{dt} p(t) = F(p(t)) \quad \text{or} \quad \frac{d}{dt} p(t) = p(t) F(p(t)) .$$

- (2) *The flow is a mapping $S: \Delta^\circ(\Omega) \times \mathbb{R}_{>0} \ni (p, t) \mapsto S(p, t) \in \Delta^\circ(\Omega)$ such that $S(p, 0) = p$ and $t \mapsto S(p, t)$ is a trajectory along F .*
- (3) *Given $f: \Delta^\circ(\Omega) \rightarrow \mathbb{R}$ with statistical gradient $p \mapsto (p, \nabla f(p)) \in T\Delta^\circ(\Omega)$ a solution of the s -gradient flow equation, $s = \pm 1$, starting at $p_0 \in \Delta^\circ(\Omega)$ at time t_0 is a curve $I \mapsto (p(t), Dp(t)) \in T\Delta^\circ(\Omega)$ such that $Dp(t) = s \nabla f(p(t))$, $t \in I$, and such that $p(t_0) = p_0$.*

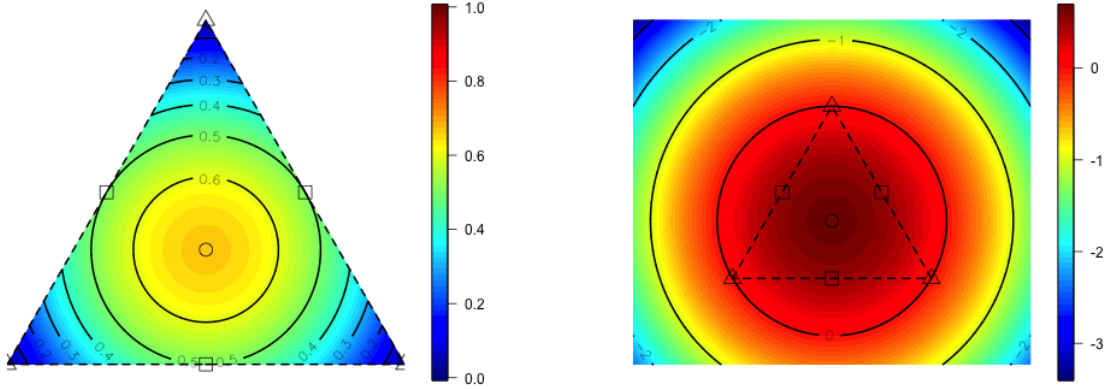


FIGURE 2. The FRAC function on the 3-point simplex (left) and extended to the affine space (right). Note difference between the flow of the ordinary gradient and the flow of the statistical gradient.

(4) *The s -gradient flow is a mapping*

$$\Delta^\circ(\Omega) \times I \ni (p, t) \mapsto S(p; t) \in \Delta^\circ(\Omega) ,$$

$0 \in I$, such that $t \mapsto S(p, t)$ is the solution of the gradient flow equation starting at p at time 0,

$$\frac{D}{dt} S(p, t) = \nabla f(S(p, t)), \quad t \in I \text{ and } S(p, 0) = p .$$

Remark 7. If we have a solution of the +1 gradient flow $Dp(t) = \nabla f(p(t))$, $t \in I$, then

$$\frac{d}{dt} f(p(t)) = \langle \nabla f(p(t)), Dp(t) \rangle_{p(t)} = \|\nabla f(p(t))\|_{p(t)}^2 = \|Dp(t)\|_{p(t)}^2 .$$

It follows that

$$f(p(t)) - f(p_0) = \int_0^t \|\nabla f(p(t))\|_{p(t)}^2 dt = \int_0^t \|Dp(t)\|_{p(t)}^2 dt .$$

Assume that $p \mapsto f(p)$ is bounded and $t \mapsto \|\nabla f(p(t))\|_{p(t)}^2 = \|Dp(t)\|_{p(t)}^2$ is uniformly continuous. Then $\lim_{t \rightarrow \infty} f(p(t)) - f(p_0) = \int_0^\infty \|\nabla f(p(t))\|_{p(t)}^2 dt = \int_0^\infty \|Dp(t)\|_{p(t)}^2 dt$ is finite, hence $\lim_{t \rightarrow \infty} \|\nabla f(p(t))\|_{p(t)} = \lim_{t \rightarrow \infty} \|Dp(t)\|_{p(t)} = 0$ (Barbalat lemma). If $p \mapsto \|\nabla f(p)\|_p^2$ has an isolated zero \bar{p} , then $\lim_{t \rightarrow \infty} p(t) = \bar{p}$

2.2. Example: fractionalization. Consider the function

$$\text{FRAC}(p) = \sum_{x \in \Omega} p(x)(1 - p(x)) = 1 - \sum_{x \in \Omega} p(x)^2, \quad p \in \Delta^\circ(\Omega) ,$$

see Fig. 2.

Its value is 0 if, and only if, p is a vertex of the simplex, otherwise it is positive. It is an index used in the Social Sciences to measure the fractionalization of the society in

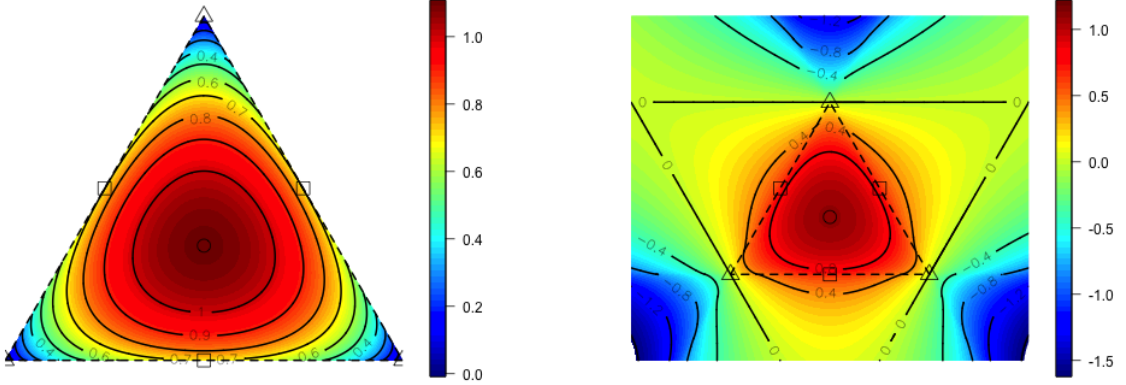


FIGURE 3. The entropy function on the 3-point simplex (left) and extended as $-\sum_{x \in \Omega} p(x) \log |p(x)|$ to the affine space (right). Note the large inverted triangle corresponding to the cases of the type $p(x) = 0$, $p(y) + p(z) = 0$.

various groups. In fact, the gradient is

$$\nabla \text{FRAC}(p) = \left(-2p(x) + 2 \sum_{y \in \Omega} p(y)^2 \middle| x \in \Omega \right)$$

which is 0 on the uniform distribution. The same Eq. extends to $A_1(\Omega)$, where the statistical gradient is 0 at the vertexes of the simplex. On the uniform distribution FRAC takes the maximal value $(n-1)/n$. The +1 gradient flow equation is

$$Dp(t) = -2p(x) + 2 \sum_{y \in \Omega} p(y)^2 \quad \text{or} \quad \dot{p}(t) = 2p(t) \left(\sum_{y \in \Omega} p(y)^2 - p(x) \right).$$

I do not know a closed form solution. The function $V(p) = \sum_{x \in \Omega} p(x)^2$ has statistical gradient $\nabla V(p) = \left(2(p(x) - \sum_{y \in \Omega} p(y)^2) \middle| x \in \Omega \right) = -\nabla \text{FRAC}(p)$, so that

$$\langle \nabla V(p), \nabla \text{FRAC}(p) \rangle_p < 0,$$

that is V is a Lyapunov function for FRAC.

2.3. Example: entropy. The combination of Dp as a definition of the derivative and of the metric $\langle \cdot, \cdot \rangle_{p(t)}$ on the vector fiber provides a computable rule for the gradient. We compute the gradient of the entropy $H(p) = -\mathbb{E}_p[\log p]$, $p \in \Delta^\circ(\Omega)$, see Fig. 3.

Let us compute the gradient:

$$\begin{aligned} \frac{d}{dt} H(p(t)) &= \\ - \frac{d}{dt} \sum_{x \in \Omega} p(x; t) \log(p(x; t)) &= - \sum_{x \in \Omega} (\log(p(x; t)) + 1) \dot{p}(x; t) = \\ & \langle -\log(p(t) - H(p(t))), Dp(t) \rangle_{p(t)}, \end{aligned}$$

hence $\nabla H(p) = -\log p - H(p)$. Note that the gradient is zero if, and only if, the distribution is uniform.

Let us solve the flow of the gradient equation

$$Dp(t) = \nabla H(p(t)), \quad p(0) = p_0 ,$$

that is

$$\frac{d}{dt} \log(p(t)) = -\log(p(t)) - H(p(t)) .$$

It follows

$$\frac{d}{dt} (e^t \log(p(t))) = -e^t H(p(t)) ,$$

hence

$$e^t \log(p(t)) = \log(p_0) - \int_0^t e^s H(p(s)) ds ,$$

so that

$$p(t) = \frac{p_0^{e^{-t}}}{\exp\left(\int_0^t e^{-(t-s)} H(p(s)) ds\right)}$$

and

$$\exp\left(\int_0^t e^{-(t-s)} H(p(s)) ds\right) = \sum_{x \in \Omega} p_0(x) e^{-t} .$$

In conclusion, we have proved that the solution of the gradient flow tends to the probability with maximal entropy, that is to the uniform distribution.

Let us consider the negative gradient flow,

$$Dp(t) = -\nabla H(p(t)), \quad p(0) = p_0 ,$$

that is

$$\frac{d}{dt} \log(p(t)) = \log(p(t)) + H(p(t)) .$$

It follows

$$\frac{d}{dt} (e^{-t} \log(p(t))) = e^{-t} H(p(t)) ,$$

hence

$$e^{-t} \log(p(t)) = \log(p_0) + \int_0^t e^{-s} H(p(s)) ds ,$$

so that

$$p(t) = \frac{p_0^{e^t}}{\exp\left(\int_0^t e^{(t-s)} H(p(s)) ds\right)}$$

and

$$\exp\left(\int_0^t e^{(t-s)} H(p(s)) ds\right) = \sum_{x \in \Omega} p_0(x) e^t .$$

In conclusion, it is easy to check that the solution of the negative gradient flow tend to the probability in $\Delta(\Omega)$ which is uniform on the set $\{x \in \Omega | p_0(x) = \max_y p_0(y)\}$.

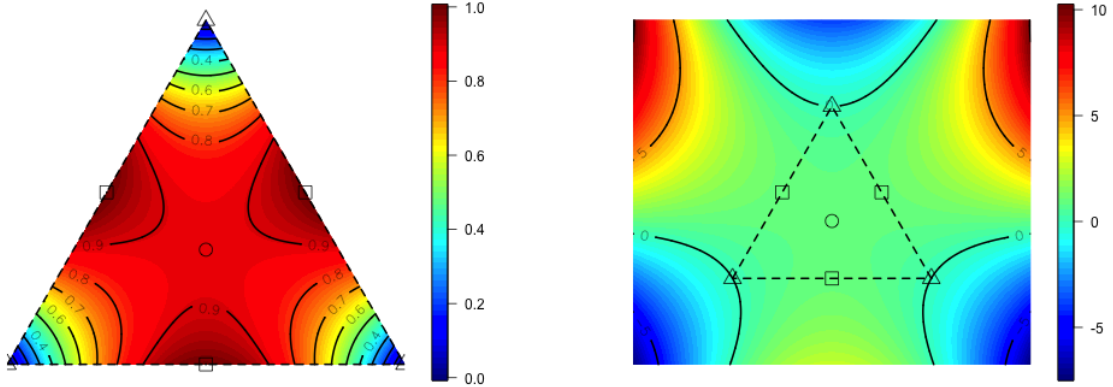


FIGURE 4. Polarization

2.4. **Example: the polarization measure.** M. Roynal-Quenol has defined the index of polarization in her PhD thesis, discussed at LSE on 2001, see [25] to be

$$\text{POL}: \Delta_n \ni p \mapsto 1 - 4 \sum_{x=0}^n \left(\frac{1}{2} - p(x) \right)^2 p(x) = 4 \sum_{x=0}^n p(x)^2 (1 - p(x))$$

From the first form it follows that the value is 1 if, and only if, $p(x) = 1/2$ or $p(x) = 0$ for all x , that is p is the uniform distribution on a couple of points, $\binom{n+1}{2} = \frac{n(n+1)}{2}$ cases. Otherwise it is < 1 . From the second form, we see it is ≥ 0 and 0 on the vertexes of the simplex only.

Note that we can extend the function POL on the affine space

$$A_1(\Omega) = \left\{ q \in \mathbb{R}^\Omega \mid \sum_{x \in \Omega} q(x) = 1 \right\} .$$

See in Fig. 4 the two cases. On such extension it is still true that the function is zero on the vertexes of the simplex $\Delta(\Omega)$, but it is not anymore true that the maximum is reached on the middle points of the edges.

The statistical gradient is the vector

$$\nabla \text{POL}(p) = \left(8p(x) - 12p(x)^2 - 8 \sum_{y \in \Omega} p(y)^2 + 12 \sum_{y \in \Omega} p(y)^3 \mid x \in \Omega \right)$$

This example shall be discussed again in the context of parameterization.

2.5. **Transports.** For each random variable $U \in B_p$, note that $\mathbb{E}_q[U - \mathbb{E}_q[U]] = 0$ and $\mathbb{E}_q \left[\frac{p}{q} U \right] = 0$, so that we can define the following transport mappings.

Definition 3.

(1) *The e-transport is the family of linear mappings*

$${}^+ \mathbb{U}_p^q: B_p \ni U \mapsto U - \mathbb{E}_q[U] \in B_q .$$

(2) The m-transport is the family of linear mappings

$${}^{-}\mathbb{U}_p^q: B_p \ni U \mapsto \frac{q}{p}U \in B_q .$$

(3) The h-transport is the family of linear mappings

$${}^0\mathbb{U}_p^q: B_p \ni U \mapsto \sqrt{\frac{p}{q}}u - \left(1 + \mathbb{E}_q \left[\sqrt{\frac{p}{q}} \right] \right)^{-1} \left(1 + \sqrt{\frac{p}{q}}\right) \mathbb{E}_q \left[\sqrt{\frac{p}{q}}u \right] \in B_q$$

Proposition 4.

- (1) ${}^{+}\mathbb{U}_q^r + {}^{+}\mathbb{U}_p^q = {}^{+}\mathbb{U}_p^r$, ${}^{-}\mathbb{U}_q^r - {}^{-}\mathbb{U}_p^q = {}^{-}\mathbb{U}_p^r$, ${}^0\mathbb{U}_q^p {}^0\mathbb{U}_p^q u = u$.
- (2) $\langle {}^{+}\mathbb{U}_p^q U, V \rangle_q = \langle U, {}^{-}\mathbb{U}_q^p V \rangle_p$.
- (3) $\langle {}^{+}\mathbb{U}_p^q U, {}^{-}\mathbb{U}_p^q V \rangle_q = \langle U, V \rangle_p$.
- (4) $\langle {}^0\mathbb{U}_p^q U, {}^0\mathbb{U}_p^q V \rangle_q = \langle U, V \rangle_p$.

2.6. Connections. Second order geometry is usually based on a notion of covariant derivative or a notion of connection. It leads to interesting applications in optimization, such as the use of a Newton method based on the computation of the proper Hessian. Here we restrict to a definition of acceleration based on the use of transports.

Let us compute the acceleration of a curve $I \mapsto p(t)$. The velocity is $t \mapsto (p(t), Dp(t)) = (p(t), \frac{d}{dt} \log(p(t))) \in T\Delta^\circ(\Omega)$. The vector $Dp(t)$ has to be checked against an element of $B_{p(t)}$, say ${}^{-}\mathbb{U}_p^{p(t)}v$. We can compute an acceleration as

$$\begin{aligned} \frac{d}{dt} \langle Dp(t), {}^{-}\mathbb{U}_p^{p(t)}v \rangle_{p(t)} &= \frac{d}{dt} \langle {}^{+}\mathbb{U}_{p(t)}^p Dp(t), v \rangle_p \\ &= \left\langle \frac{d}{dt} {}^{+}\mathbb{U}_{p(t)}^p Dp(t), v \right\rangle_p \\ &= \left\langle {}^{+}\mathbb{U}_p^{p(t)} \frac{d}{dt} {}^{+}\mathbb{U}_{p(t)}^p Dp(t), {}^{-}\mathbb{U}_p^{p(t)}v \right\rangle_{p(t)} \end{aligned}$$

The *exponential acceleration* is

$$\begin{aligned} {}^{+}\mathbb{U}_p^{p(t)} \frac{d}{dt} {}^{+}\mathbb{U}_{p(t)}^p Dp(t) &= {}^{+}\mathbb{U}_p^{p(t)} \frac{d}{dt} \left(\frac{\dot{p}(t)}{p(t)} - \mathbb{E}_p \left[\frac{\dot{p}(t)}{p(t)} \right] \right) \\ &= {}^{+}\mathbb{U}_p^{p(t)} \left(\frac{\ddot{p}(t)p(t) - \dot{p}(t)^2}{p(t)^2} - \mathbb{E}_p \left[\frac{\ddot{p}(t)p(t) - \dot{p}(t)^2}{p(t)^2} \right] \right) \\ &= \frac{\ddot{p}(t)p(t) - \dot{p}(t)^2}{p(t)^2} - \mathbb{E}_{p(t)} \left[\frac{\ddot{p}(t)p(t) - \dot{p}(t)^2}{p(t)^2} \right] \\ &= \frac{\ddot{p}(t)}{p(t)} - (Dp(t))^2 + \mathbb{E}_{p(t)} [(Dp(t))^2] . \end{aligned}$$

Exponential families have null exponential acceleration: for $p(t) = \exp(tu - \psi(t))p$, we have

$$\begin{aligned} Dp(t) &= u - \dot{\psi}(t) \\ \frac{d}{dt} Dp(t) &= -\ddot{\psi}(t) \\ \mathbb{E}_{p(t)} [(Dp(t))^2] &= \ddot{\psi}(t) \end{aligned}$$

We could have computed the acceleration as

$$\begin{aligned} \frac{d}{dt} \langle Dp(t), {}^+\mathbb{U}_p^{p(t)}v \rangle_{p(t)} &= \frac{d}{dt} \left\langle -\mathbb{U}_{p(t)}^p Dp(t), v \right\rangle_p \\ &= \left\langle \frac{d}{dt} -\mathbb{U}_{p(t)}^p Dp(t), v \right\rangle_p \\ &= \left\langle -\mathbb{U}_p^{p(t)} \frac{d}{dt} -\mathbb{U}_{p(t)}^p Dp(t), {}^+\mathbb{U}_p^{p(t)}v \right\rangle_{p(t)} \end{aligned}$$

hence we can compute the *mixture acceleration* as follows

$$\begin{aligned} -\mathbb{U}_p^{p(t)} \frac{d}{dt} -\mathbb{U}_{p(t)}^p Dp(t) &= \frac{p}{p(t)} \frac{d}{dt} \left(\frac{p(t)}{p} \frac{\dot{p}(t)}{p(t)} \right) \\ &= \frac{\ddot{p}(t)}{p(t)}. \end{aligned}$$

It follows that *mixture models have null mixture acceleration*.

We could compute the *Riemannian acceleration* as ${}^0\mathbb{U}_p^{p(t)} \frac{d}{dt} {}^0\mathbb{U}_{p(t)}^p Dp(t)$. See some developments in [20, 21, 14]

2.7. Atlases. We now turn to the explicit introduction of atlases of charts. Each of these chart correspond to a different geometry, in particular to different notion of geodesics.

Definition 5 (Exponential atlas). *For each $p \in \Delta^\circ(\Omega)$ we define*

$$s_p: T\Delta^\circ(\Omega) \ni (q, w) \mapsto \left(\log \frac{q}{p} - \mathbb{E}_p \left[\log \frac{q}{p} \right], {}^+\mathbb{U}_q^p w \right) \in B_p \times B_p$$

Proposition 6 (Exponential atlas).

- (1) *If $u = s_p(q)$, then $q = e^{u - K_p(u)} \cdot p$ with $K_p(u) = \log \mathbb{E}_p [e^u]$.*
- (2) *The patches are*

$$s_p^{-1}: (u, v) \mapsto (e^{u - K_p(u)} \cdot p, v - dK_p(u)[v])$$

- (3) *The transitions are*

$$s_{p_1} \circ s_{p_2}^{-1}: (u, v) \mapsto \left({}^+\mathbb{U}_{p_2}^{p_1} u + \log \frac{p_2}{p_1} - \mathbb{E}_{p_1} \left[\log \frac{p_2}{p_1} \right] \right) \in B_{p_1} \times B_{p_1}$$

- (4) *The tangent bundle identifies with the statistical bundle,*

$$\frac{d}{dt} s_p(p(t)) = {}^+\mathbb{U}_{p(t)}^p Dp(t).$$

Definition 7 (Mixture atlas). *For each $p \in \Delta^\circ(\Omega)$ we define*

$$\eta_p: T\Delta^\circ(\Omega) \ni (q, w) \mapsto \left(\frac{q}{p} - 1, -\mathbb{U}_q^p w \right) \in B_p \times B_p$$

Proposition 8 (Mixture atlas).

- (1) *If $u = \eta_p(q)$, then $q = (1 + u)p$.*
- (2) *The patches are*

$$\eta_p^{-1}: (u, v) \mapsto ((1 + u)p, (1 + u)w)$$

- (3) *The transitions are*

$$\eta_{p_1} \circ \eta_{p_2}^{-1}: (u, v) \mapsto \left((1 + u) \frac{p_2}{p_1} - 1, -\mathbb{U}_{p_1}^{p_2} v \right)$$

(4) *The tangent bundle identifies with the statistical bundle,*

$$\frac{d}{dt}\eta_p(p(t)) = {}^{-\mathbb{U}}_{p(t)} Dp(t) .$$

It is possible to define the Riemannian atlas, see [21].

2.8. Using parameters. We still study the geometry of the full simplex, but now we introduce parameters. This section is taken from [22].

Computations are usually performed in a *parametrization*

$$\pi: \Theta \ni \boldsymbol{\theta} \mapsto \pi(\boldsymbol{\theta}) \in \Delta^\circ(\Omega),$$

Θ being an open set in \mathbb{R}^n . The j -th coordinate curve is obtained by fixing the other $n - 1$ components and moving θ_j only. The scores of the j -th coordinate curves are the random variables

$$D_j\pi(\boldsymbol{\theta}) = \frac{\partial}{\partial\theta_j} \log \pi(\boldsymbol{\theta}), \quad j = 1, \dots, n.$$

The sequence $(D_j\pi(\boldsymbol{\theta}): j = 1, \dots, n)$ is a vector basis of the tangent space $B_{\pi(\boldsymbol{\theta})}$. The representation of the scalar product in such a basis is

$$\left\langle \sum_{i=1}^n \alpha_i D_i\pi(\boldsymbol{\theta}), \sum_{j=1}^n \beta_j D_j\pi(\boldsymbol{\theta}) \right\rangle_{\pi(\boldsymbol{\theta})} = \sum_{i,j=1}^n \alpha_i \beta_j I_{ij}(\boldsymbol{\theta}) .$$

Definition 9 (Fisher Information). *The matrix $I(\boldsymbol{\theta}) = \left[\langle D_i\pi(\boldsymbol{\theta}), D_j\pi(\boldsymbol{\theta}) \rangle_{\pi(\boldsymbol{\theta})} \right]_{i,j=1}^n$ is the Fisher information matrix.*

If $\boldsymbol{\theta} \mapsto \tilde{\phi}(\boldsymbol{\theta})$ is the expression in the parameters of a function $\phi: \Delta^\circ(\Omega) \rightarrow \mathbb{R}$, that is $\tilde{\phi}(\boldsymbol{\theta}) = \phi(\pi(\boldsymbol{\theta}))$, and $t \mapsto \boldsymbol{\theta}(t)$ is the expression in the parameters of a generic curve $p: I \rightarrow \Delta^\circ(\Omega)$, then the components of the gradient in (11) are expressed in terms of the ordinary gradient by observing that

$$\frac{d}{dt}\phi(p(t)) = \frac{d}{dt}\tilde{\phi}(\boldsymbol{\theta}(t)) = \sum_{j=1}^n \frac{\partial}{\partial\theta_j} \tilde{\phi}(\boldsymbol{\theta}(t)) \dot{\theta}_j(t).$$

As $Dp(t) = \sum_j D_j\pi(\boldsymbol{\theta}(t)) \dot{\theta}_j(t)$, we obtain from (11)

$$(1) \quad \frac{d}{dt}\phi(p(t)) = \langle \nabla\phi(p(t)), Dp(t) \rangle_{p(t)} = \sum_j \left\langle \nabla\phi(p(t)), D_j\pi(\boldsymbol{\theta}(t)) \dot{\theta}_j(t) \right\rangle_{p(t)} .$$

Definition 10 ([2]). *The natural gradient is a vector $\tilde{\nabla}\tilde{\phi}(\boldsymbol{\theta})$ whose components are the coordinates of the gradient $\nabla\tilde{\phi}(\pi(\boldsymbol{\theta})) \in B_{\pi(\boldsymbol{\theta})}$ in its π -basis, that is*

$$\nabla\phi(\pi(\boldsymbol{\theta})) = \sum_{j=1}^n (\tilde{\nabla}\tilde{\phi}(\boldsymbol{\theta}))_j D_j\pi(\boldsymbol{\theta}).$$

By substitution of the expression in (1) we obtain

$$(2) \quad \tilde{\nabla}\tilde{\phi}(\boldsymbol{\theta}) = \nabla\tilde{\phi}(\boldsymbol{\theta}) I^{-1}(\boldsymbol{\theta}).$$

The *common parametrization* of the (flat) simplex $\Delta^\circ(\Omega)$ is the projection on the solid simplex $\Gamma_n = \left\{ \boldsymbol{\eta} \in \mathbb{R}^n \mid 0 < \eta_j, \sum_{j=1}^n \eta_j < 1 \right\}$, that is

$$\pi: \Gamma_n \ni \boldsymbol{\eta} \mapsto \left(1 - \sum_{j=1}^n \eta_j, \eta_1, \dots, \eta_n \right) \in \Delta^\circ(\Omega),$$

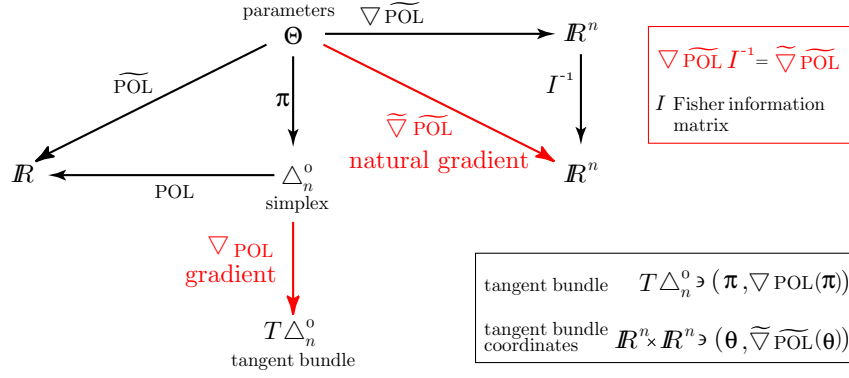


FIGURE 5. Diagram of the action of the natural gradient in a given parametrization $\pi: \Theta \rightarrow \Delta^\circ(\Omega)$.

in which case $\partial_j \pi(\boldsymbol{\eta})$, $j = 1, \dots, n$, is the random variable with values -1 at $x = 0$, 1 at $x = j$, 0 otherwise, hence $\partial_j \pi(\boldsymbol{\eta}) = ((X = j) - (X = 0))$ and

$$D_j \pi(\boldsymbol{\eta}) = ((X = j) - (X = 0)) / \pi(\boldsymbol{\eta}).$$

The element (j, h) of the Fisher information matrix is

$$I_{jh}(\boldsymbol{\eta}) = \mathbb{E}_{\pi(\boldsymbol{\eta})} \left[\frac{(X = j) - (X = 0)}{\pi(X; \boldsymbol{\eta})} \frac{(X = h) - (X = 0)}{\pi(X; \boldsymbol{\eta})} \right] = \sum_x \pi(x, \boldsymbol{\eta})^{-1} ((x = j)(j = h) + (x = 0)) = \eta_j^{-1} (j = h) + \left(1 - \sum_k \eta_k \right)^{-1}$$

hence

$$I(\boldsymbol{\eta}) = \text{diag}(\boldsymbol{\eta})^{-1} + \left(1 - \sum_{j=1}^n \eta_j \right)^{-1} [\mathbb{1}_{i,j=1}^n].$$

As an example we consider $n = 3$. The Fisher information matrix, its inverse and the determinant of the inverse are, respectively,

$$I(\eta_1, \eta_2, \eta_3) = (1 - \eta_1 - \eta_2 - \eta_3)^{-1} \begin{bmatrix} \eta_1^{-1}(1 - \eta_2 - \eta_3) & 1 & 1 \\ 1 & \eta_2^{-1}(1 - \eta_1 - \eta_3) & 1 \\ 1 & 1 & \eta_3^{-1}(1 - \eta_1 - \eta_2) \end{bmatrix},$$

$$I(\eta_1, \eta_2, \eta_3)^{-1} = \begin{bmatrix} (1 - \eta_1)\eta_1 & -\eta_1\eta_2 & -\eta_1\eta_3 \\ -\eta_1\eta_2 & (1 - \eta_2)\eta_2 & -\eta_2\eta_3 \\ -\eta_1\eta_3 & -\eta_2\eta_3 & (1 - \eta_3)\eta_3 \end{bmatrix},$$

$$\det(I(\eta_1, \eta_2, \eta_3)^{-1}) = (1 - \eta_1 - \eta_2 - \eta_3)\eta_1\eta_2\eta_3.$$

Note that the computation of the inverse of $I(\boldsymbol{\eta})$ is an application of the Sherman-Morrison formula and the computation of the determinant of $I(\boldsymbol{\eta})^{-1}$ is an application of the matrix determinant lemma.

For general n , we have the following Proposition, whose interest stems from the definition of natural gradient, see Eq. (2).

Proposition 11.

- (1) *The inverse of the Fisher information matrix is*

$$I(\boldsymbol{\eta})^{-1} = \text{diag}(\boldsymbol{\eta}) - \boldsymbol{\eta}\boldsymbol{\eta}^t$$

- (2) In particular, $I(\boldsymbol{\eta})^{-1}$ is zero on the vertexes of the simplex, only.
(3) The determinant of the inverse Fisher information matrix is

$$\det(I(\boldsymbol{\eta})^{-1}) = \left(1 - \sum_{i=1}^n \eta_i\right) \prod_{i=1}^n \eta_i.$$

- (4) The determinant of $I(\boldsymbol{\eta})^{-1}$ is zero on the borders of the simplex, only.
(5) On the interior of each facet, the rank of $I(\boldsymbol{\eta})^{-1}$ is $n - 1$ and the $n - 1$ linear independent column vectors generate the subspace parallel to the facet itself.

Proof.

- (1) By direct computation, $I(\boldsymbol{\eta})I(\boldsymbol{\eta})^{-1}$ is the identity matrix.
(2) The diagonal elements of $I(\boldsymbol{\eta})^{-1}$ are zero if $\eta_j = 1$ or $\eta_j = 0$, for $j = 1, \dots, n$. If, for a given j , $\eta_j = 1$, then the elements of $I(\boldsymbol{\eta})^{-1}$ are zero if $\eta_h = 0$, $h \neq j$. The remaining case corresponds to $\eta_j = 0$ for all j . Then $I(\boldsymbol{\eta})^{-1} = 0$ on all the vertexes of the simplex.
(3) It follows from Matrix Determinant Lemma.
(4) The determinant factors in terms corresponding to the equations of the facets.
(5) Given i , the conditions $\eta_i = 0$ and $\eta_j \neq 0, 1$ for all $j \neq i$, define the interior of the facet orthogonal to standard base vector e_i . In this case the i -th row and the i -th column of $I(\boldsymbol{\eta})^{-1}$ are zero and the complement matrix corresponds to the inverse of a Fisher information matrix in dimension $n - 1$ with non zero determinant. It follows that the subspace generated by the columns has dimension $n - 1$ and coincides with the space orthogonal to η_i . Consider the facet defined by $(1 - \sum_{i=1}^n \eta_i) = 0$, $\eta_i \neq 0, 1$ for all i . For a given j , the matrix without the j -th row and the j -th column has determinant $\left(1 - \sum_{i=1, i \neq j}^n \eta_i\right) \prod_{i=1, i \neq j}^n \eta_i$. On the considered facet this determinant is different to zero and $I(\boldsymbol{\eta})^{-1}$ has rank $n - 1$ and their columns are orthogonal to the constant vector. \square

An other parametrization is the *exponential parametrization* based on the exponential family with sufficient statistics $X_j = (X = j)$, $j = 1, \dots, n$,

$$\pi: \mathbb{R}^n \ni \boldsymbol{\theta} \mapsto \exp\left(\sum_{j=1}^n \theta_j X_j - \psi(\boldsymbol{\theta})\right) \frac{1}{n+1}$$

where

$$\psi(\boldsymbol{\theta}) = \log\left(1 + \sum_j e^{\theta_j}\right) - \log(n+1).$$

See [22] for an illustration of the gradient flow in case of the function POL.

3. FINITE STATE SPACE: EXPONENTIAL FAMILIES

This section is taken from our paper “Natural Gradient Flow in the Mixture Geometry of a Discrete Exponential Family” to appear in *Entropy*.

Let Ω be a finite set of points $\mathbf{x} = (x_1, \dots, x_n)$ and μ the counting measure on Ω . In this case a density $p \in \mathcal{P}_{\geq}$ is a probability function i.e., $p: \Omega \rightarrow \mathbb{R}_+$ such that $\sum_{\mathbf{x} \in \Omega} p(\mathbf{x}) = 1$.

Given a set $\mathcal{B} = \{T_1, \dots, T_d\}$ of random variables such that, if $\sum_{j=1}^d c_j T_j$ is constant, then $c_1 = \dots = c_d = 0$. E.g., if $\sum_{\mathbf{x} \in \Omega} T_j(\mathbf{x}) = 0$, $j = 0, \dots, d$, and the \mathcal{B} is a linear basis. We say that \mathcal{B} is a set of *affinely independent* random variables. If \mathcal{B} is a linear basis, it is affinely independent if and only if $\{1, T_1, \dots, T_d\}$ is a linear basis.

We consider the statistical model \mathcal{E} whose elements are uniquely identified by the natural parameters $\boldsymbol{\theta}$ in the exponential family with sufficient statistics \mathcal{B} , namely

$$p_{\boldsymbol{\theta}} \in \mathcal{E} \quad \Leftrightarrow \quad \log p_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{i=1}^d \theta_i T_i(\mathbf{x}) - \psi(\boldsymbol{\theta}), \quad \boldsymbol{\theta} \in \mathbb{R}^d,$$

see [4].

The proper convex function $\psi: \mathbb{R}^d$,

$$\boldsymbol{\theta} \mapsto \psi(\boldsymbol{\theta}) = \log \sum_{\mathbf{x} \in \Omega} e^{\boldsymbol{\theta} \cdot \mathbf{T}(\mathbf{x})} = \boldsymbol{\theta} \cdot \mathbb{E}_{p_{\boldsymbol{\theta}}}[\mathbf{T}] - \mathbb{E}_{p_{\boldsymbol{\theta}}}[\log(p_{\boldsymbol{\theta}})]$$

is the *cumulant generating function* of the sufficient statistics \mathbf{T} , in particular,

$$\nabla \psi(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}}[\mathbf{T}], \quad \text{Hess } \psi(\boldsymbol{\theta}) = \text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}, \mathbf{T}).$$

Moreover, the entropy of $p_{\boldsymbol{\theta}}$ is

$$H(p_{\boldsymbol{\theta}}) = -\mathbb{E}_{p_{\boldsymbol{\theta}}}[\log(p_{\boldsymbol{\theta}})] = \psi(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \nabla \psi(\boldsymbol{\theta}).$$

The mapping $\nabla \psi$ is 1-to-1 onto the interior M° of the *marginal polytope*, that is the convex span of the values of the sufficient statistics $M = \{\mathbf{T}(\mathbf{x}) | \mathbf{x} \in \Omega\}$. Note that no extra condition is required because on a finite state space all random variables are bounded. Nonetheless, even in this case the proof is not trivial, see [4].

Convex conjugation applies, see [27, §25], with the definition

$$\psi_*(\boldsymbol{\eta}) = \sup \{ \boldsymbol{\theta} \in \mathbb{R}^d | \boldsymbol{\theta} \cdot \boldsymbol{\eta} - \psi(\boldsymbol{\theta}) \}, \quad \boldsymbol{\eta} \in \mathbb{R}^d.$$

The concave function $\boldsymbol{\theta} \mapsto \boldsymbol{\eta} \cdot \boldsymbol{\theta} - \psi(\boldsymbol{\theta})$ has divergence mapping $\boldsymbol{\theta} \mapsto \boldsymbol{\eta} - \nabla \psi(\boldsymbol{\theta})$ and the equation $\boldsymbol{\eta} = \nabla \psi(\boldsymbol{\theta})$ has solution if and only if $\boldsymbol{\eta}$ belongs to the interior M° of the marginal polytope. The restriction $\phi = \psi_*|_{M^\circ}$ is the *Legendre conjugate* of ψ , and it is computed by

$$\phi: M^\circ \ni \boldsymbol{\eta} \mapsto \boldsymbol{\eta} \cdot (\nabla \psi)^{-1}(\boldsymbol{\eta}) - \psi \circ (\nabla \psi)^{-1}(\boldsymbol{\eta}) \in \mathbb{R}.$$

The Legendre conjugate ϕ is such that $\nabla \phi = (\nabla \psi)^{-1}$ and it provides an alternative parameterization of \mathcal{E} with the so called *expectation* or *mixture* parameter $\boldsymbol{\eta} = \nabla \psi(\boldsymbol{\theta})$,

$$(3) \quad p_{\boldsymbol{\eta}} = \exp((\mathbf{T} - \boldsymbol{\eta}) \cdot \nabla \phi(\boldsymbol{\eta}) + \phi(\boldsymbol{\eta})).$$

While in the $\boldsymbol{\theta}$ -parameters the entropy is $H(p_{\boldsymbol{\theta}}) = \psi(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \nabla \psi(\boldsymbol{\theta})$, in the $\boldsymbol{\eta}$ -parameters the ϕ function gives the the negative of the entropy: $-H(p_{\boldsymbol{\eta}}) = \mathbb{E}_{p_{\boldsymbol{\eta}}}[\log p_{\boldsymbol{\eta}}] = \phi(\boldsymbol{\eta})$.

Proposition 12.

$$(1) \quad \text{Hess } \phi(\boldsymbol{\eta}) = (\text{Hess } \psi(\boldsymbol{\theta}))^{-1} \text{ when } \boldsymbol{\eta} = \nabla \psi(\boldsymbol{\theta}).$$

- (2) The Fisher information matrix of the statistical model given by the exponential family in the $\boldsymbol{\theta}$ -parameters is $I_e(\boldsymbol{\theta}) = \text{Cov}_{p_{\boldsymbol{\theta}}}(\nabla \log p_{\boldsymbol{\theta}}, \nabla \log p_{\boldsymbol{\theta}}) = \text{Hess } \psi(\boldsymbol{\theta})$.
- (3) The Fisher information matrix of the statistical model given by the exponential family in the $\boldsymbol{\eta}$ -parameters is $I_m(\boldsymbol{\theta}) = \text{Cov}_{p_{\boldsymbol{\eta}}}(\nabla \log p_{\boldsymbol{\eta}}, \nabla \log p_{\boldsymbol{\eta}}) = \text{Hess } \phi(\boldsymbol{\eta})$.

Proof. Derivation of the equality $\nabla \phi = (\nabla \psi)^{-1}$ gives the first item. The second item is a property of the cumulant generating function ψ . The third item follows from Eq. (3). \square

3.1. Statistical manifold. The exponential family \mathcal{E} is an elementary manifold in either the $\boldsymbol{\theta}$ - or the $\boldsymbol{\eta}$ -parameterization, named respectively *exponential* or *mixture* parameterization. We discuss now the proper definition of the tangent bundle $T\mathcal{E}$.

Definition 13 (Velocity). *If $I \ni t \mapsto p_t$, I open interval, is a differentiable curve in \mathcal{E} , then its velocity vector is identified with its Fisher score:*

$$\frac{D}{dt}p(t) = \frac{d}{dt} \log(p_t) .$$

The capital- D notation is taken from differential geometry, see the classical monograph [7].

Definition 14 (Tangent space). *In the expression of the curve by the exponential parameters, the velocity is*

$$(4) \quad \frac{D}{dt}p(t) = \frac{d}{dt} \log(p_t) = \frac{d}{dt} (\boldsymbol{\theta}(t) \cdot \mathbf{T} - \psi(\boldsymbol{\theta}(t))) = \dot{\boldsymbol{\theta}}(t) \cdot (\mathbf{T} - \mathbb{E}_{\boldsymbol{\theta}(t)}[\mathbf{T}]) ,$$

that is it equals the statistics whose coordinates are $\dot{\boldsymbol{\theta}}(t)$ in the basis of the sufficient statistics centered at p_t . As a consequence, we identify the tangent space at each $p \in \mathcal{E}$ with the vector space of centered sufficient statistics, that is

$$T_p\mathcal{E} = \text{Span} (T_j - \mathbb{E}_p[T_j] | j = 1, \dots, d) .$$

In the mixture parameterization of Eq. (3) the computation of the velocity is

$$(5) \quad \frac{D}{dt}p(t) = \frac{d}{dt} \log(p_t) = \frac{d}{dt} (\nabla \phi(\boldsymbol{\eta}(t)) \cdot (\mathbf{T} - \boldsymbol{\eta}(t)) + \phi(\boldsymbol{\eta}(t))) = \\ (\text{Hess } \phi(\boldsymbol{\eta}(t)) \dot{\boldsymbol{\eta}}(t)) \cdot (\mathbf{T} - \boldsymbol{\eta}(t)) = \dot{\boldsymbol{\eta}}(t) \cdot [\text{Hess } \phi(\boldsymbol{\eta}(t)) (\mathbf{T} - \boldsymbol{\eta}(t))] .$$

The last equality provides the interpretation of $\dot{\boldsymbol{\eta}}(t)$ as the coordinate of the velocity in the *conjugate* vector basis $\text{Hess } \phi(\boldsymbol{\eta}(t)) (\mathbf{T} - \boldsymbol{\eta}(t))$, that is the basis of velocities along the $\boldsymbol{\eta}$ coordinates.

In conclusion, the first order geometry is characterized as follows.

Definition 15 (Tangent bundle $T\mathcal{E}$). *The tangent space at each $p \in \mathcal{E}$ is a vector space of random variables $T_p\mathcal{E} = \text{Span} (T_j - \mathbb{E}_p[T_j] | j = 1, \dots, d)$ and the tangent bundle $T\mathcal{E} = \{(p, V) | p \in \mathcal{E}, V \in T_p\mathcal{E}\}$, as a manifold, is defined by the chart*

$$(6) \quad T\mathcal{E} \ni (e^{\boldsymbol{\theta} \cdot \mathbf{T} - \psi(\boldsymbol{\theta})}, \mathbf{v} \cdot (\mathbf{T} - \mathbb{E}_{\boldsymbol{\theta}}[\mathbf{T}])) \mapsto (\boldsymbol{\theta}, \mathbf{v}) .$$

Proposition 16.

(1) If $V = \mathbf{v} \cdot (\mathbf{T} - \boldsymbol{\eta}) \in T_{p_\eta} \mathcal{E}$, then V is represented in the conjugate basis as

$$(7) \quad V = \mathbf{v} \cdot (\mathbf{T} - \boldsymbol{\eta}) = \mathbf{v} \cdot (\text{Hess } \phi(\boldsymbol{\eta}))^{-1} \text{Hess } \phi(\boldsymbol{\eta}) (\mathbf{T} - \boldsymbol{\eta}) = \\ ((\text{Hess } \phi(\boldsymbol{\eta}))^{-1} \mathbf{v}) \cdot \text{Hess } \phi(\boldsymbol{\eta}) (\mathbf{T} - \boldsymbol{\eta}).$$

(2) The mapping $(\text{Hess } \phi(\boldsymbol{\eta}))^{-1}$ maps the coordinates \mathbf{v} of a tangent vector $V \in T_{p_\eta} \mathcal{E}$ with respect to the basis of centered sufficient statistics to the coordinates \mathbf{v}^* with respect to the conjugate basis.

(3) In the $\boldsymbol{\theta}$ -parameters the transformation is $\mathbf{v} \mapsto \mathbf{v}^* = \text{Hess } \psi(\boldsymbol{\theta}) \mathbf{v}$.

Remark 8. In the finite state space case it is not necessary to go on to the formal construction of a dual tangent bundle because all finite dimensional vector spaces are isomorphic. However, this step is compulsory in the infinite state space case, as it was done in [21]. Moreover, the explicit construction of natural connections and natural parallel transports of the tangent and dual tangent bundle is unavoidable when considering the second order calculus as it was done in [21, 17] in order to compute Hessians and implement Newton methods of optimization. However, the scope of the present paper is restricted to a basic study of gradient flows, hence from now on we focus on the Riemannian structure and disregard all second order topics.

Proposition 17 (Riemannian metric). *The tangent bundle has a Riemannian structure with the natural scalar product of each $T_p \mathcal{E}$, $\langle V, W \rangle_p = \mathbb{E}_p [VW]$. In the basis of sufficient statistics the metric is expressed by the Fisher information matrix $I(p) = \text{Cov}_p(\mathbf{T}, \mathbf{T})$, while in the conjugate basis it is expressed by the inverse Fisher matrix $I^{-1}(p)$.*

Proof. In the basis of the sufficient statistics, $V = \mathbf{v} \cdot (\mathbf{T} - \mathbb{E}_p[\mathbf{T}])$, $W = \mathbf{w} \cdot (\mathbf{T} - \mathbb{E}_p[\mathbf{T}])$, so that

$$(8) \quad \langle V, W \rangle_p = \mathbf{v}' \mathbb{E}_p [(\mathbf{T} - \mathbb{E}_p[\mathbf{T}]) (\mathbf{T} - \mathbb{E}_p[\mathbf{T}])'] \mathbf{w} = \mathbf{v}' \text{Cov}_p(\mathbf{T}, \mathbf{T}) \mathbf{w} = \mathbf{v}' I(p) \mathbf{w},$$

where $I(p) = \text{Cov}_p(\mathbf{T}, \mathbf{T})$ is the Fisher information matrix.

If $p = p_\theta = p_\eta$, the conjugate basis at p is

$$(9) \quad \text{Hess } \phi(\boldsymbol{\eta})(\mathbf{T} - \boldsymbol{\eta}) = \text{Hess } \psi(\boldsymbol{\theta})^{-1} (\mathbf{T} - \nabla \phi(\boldsymbol{\theta})) = I^{-1}(p) (\mathbf{T} - \mathbb{E}_p[\mathbf{T}]),$$

so that for elements of the tangent space expressed in the conjugate basis we have $V = \mathbf{v}^* \cdot I^{-1}(p) (\mathbf{T} - \mathbb{E}_p[\mathbf{T}])$, $W = \mathbf{w}^* \cdot I^{-1}(p) (\mathbf{T} - \mathbb{E}_p[\mathbf{T}])$, thus

$$(10) \quad \langle V, W \rangle_p = \mathbf{v}^{*'} \mathbb{E}_p [I^{-1}(p) \cdot (\mathbf{T} - \mathbb{E}_p[\mathbf{T}]) (\mathbf{T} - \mathbb{E}_p[\mathbf{T}])' I^{-1}(p)] \mathbf{w}^* = \mathbf{v}^{*'} I^{-1}(p) \mathbf{w}^*.$$

□

3.2. Gradient. For each C^1 real function $F: \mathcal{E} \rightarrow \mathbb{R}$, its gradient is defined by taking the derivative along a C^1 curve $I \mapsto p(t)$, $p = p(0)$, and writing it in the Riemannian metrics,

$$(11) \quad \left. \frac{d}{dt} \hat{F}(\boldsymbol{\theta}(t)) \right|_{t=0} = \left\langle \nabla F(p), \left. \frac{D}{dt} p(t) \right|_{t=0} \right\rangle_p, \quad \nabla F(p) \in T_p \mathcal{E}.$$

If $\boldsymbol{\theta} \mapsto \hat{F}(\boldsymbol{\theta})$ is the expression of F in the parameter $\boldsymbol{\theta}$, and $t \mapsto \boldsymbol{\theta}(t)$ is the expression of the curve, then $\left. \frac{d}{dt} \hat{F}(\boldsymbol{\theta}(t)) \right|_{t=0} = \nabla \hat{F}(\boldsymbol{\theta}(t)) \cdot \dot{\boldsymbol{\theta}}(t)$ so that at $p = p_{\boldsymbol{\theta}(0)}$, with velocity $V =$

$\frac{D}{dt}p(t)|_{t=0} = \dot{\boldsymbol{\theta}}(0) \cdot (\mathbf{T} - \nabla\psi(\boldsymbol{\theta}(0)))$, so that we obtain the celebrated Amari's *natural gradient* of [2]:

$$(12) \quad \langle \nabla F(p), V \rangle_p = \left(\text{Hess } \psi(\boldsymbol{\theta}(0))^{-1} \nabla \hat{F}(\boldsymbol{\theta}(0)) \right)' \text{Hess } \psi(\boldsymbol{\theta}(0)) \dot{\boldsymbol{\theta}}(0) .$$

If $\boldsymbol{\eta} \mapsto \check{F}(\boldsymbol{\eta})$ is the expression of F in the parameter $\boldsymbol{\eta}$, and $t \mapsto \boldsymbol{\eta}(t)$ is the expression of the curve, then $\frac{d}{dt}\check{F}(\boldsymbol{\eta}(t)) = \nabla \check{F}(\boldsymbol{\eta}(t)) \cdot \dot{\boldsymbol{\eta}}(t)$ so that at $p = p_{\boldsymbol{\eta}(0)}$, with velocity $V = \frac{d}{dt} \log(p(t))|_{t=0} = \dot{\boldsymbol{\eta}}(0) \cdot \text{Hess } \phi(\boldsymbol{\eta}(0))(\mathbf{T} - \boldsymbol{\eta}(0))$,

$$(13) \quad \langle \nabla F(p), V \rangle_p = (\text{Hess } \phi(\boldsymbol{\eta}(0))^{-1} \nabla \hat{F}(\boldsymbol{\eta}(0)))' \text{Hess } \phi(\boldsymbol{\eta}(0)) \dot{\boldsymbol{\eta}}(0).$$

We summarize all notions of gradient in the following definition.

Definition 18 (Gradients).

- (1) *The random variable $\nabla F(p)$ uniquely defined by Eq. (11) is called statistical gradient of F at p . The mapping $\nabla F: \mathcal{E} \ni p \mapsto \nabla F(p)$ is a vector field of $T\mathcal{E}$.*
- (2) *The vector $\tilde{\nabla} \hat{F}(\boldsymbol{\theta}) = \text{Hess } \psi(\boldsymbol{\theta})^{-1} \nabla \hat{F}(\boldsymbol{\theta})$ of (12) is the expression of the statistical gradient in the $\boldsymbol{\theta}$ in the basis of sufficient statistics, and it is called Amari's natural gradient, while $\nabla \hat{F}(\boldsymbol{\theta})$, which is the expression in the conjugate basis of the sufficient statistics, is called Amari's vanilla gradient.*
- (3) *The vector $\tilde{\nabla} \check{F}(\boldsymbol{\eta}) = \text{Hess } \phi(\boldsymbol{\eta})^{-1} \nabla \check{F}(\boldsymbol{\eta})$ of (12) is the expression of the statistical gradient in the $\boldsymbol{\eta}$ parameter and in the conjugate basis of sufficient statistics, and it is called Amari's natural gradient, while $\nabla \check{F}(\boldsymbol{\eta})$, which is the expression in the basis of sufficient statistics, is called Amari's vanilla gradient.*

Given a vector field of \mathcal{E} i.e., a mapping G defined on \mathcal{E} such that $G(p) \in T_p\mathcal{E}$ —which is called a *section* of the tangent bundle in the standard differential geometric language—an integral curve from p is a curve $I \ni t \mapsto p(t)$ such that $p(0) = p$ and $\frac{D}{dt}p(t) = G(p(t))$. In the $\boldsymbol{\theta}$ parameters $G(p_{\boldsymbol{\theta}}) = \hat{\mathbf{G}}(\boldsymbol{\theta}) \cdot (\mathbf{T} - \nabla\psi(\boldsymbol{\theta}))$, so that the differential equation is expressed by $\dot{\boldsymbol{\theta}}(t) = \hat{\mathbf{G}}(\boldsymbol{\theta}(t))$. In the $\boldsymbol{\eta}$ parameters, $G(p_{\boldsymbol{\eta}}) = \check{\mathbf{G}}(\boldsymbol{\eta}) \cdot \text{Hess } \phi(\boldsymbol{\eta})(\mathbf{T} - \boldsymbol{\eta})$ and the differential equation is $\dot{\boldsymbol{\eta}}(t) = \check{\mathbf{G}}(\boldsymbol{\eta}(t))$.

Definition 19 (Gradient flow).

- *The gradient flow of the real function $F: \mathcal{E}$ is the flow of the differential equation $\frac{D}{dt}p(t) = \nabla F(p(t))$ i.e. $\frac{d}{dt}p(t) = p(t) \nabla F(p(t))$.*
- *The expression in the $\boldsymbol{\theta}$ parameters is $\dot{\boldsymbol{\theta}}(t) = \tilde{\nabla} \hat{F}(\boldsymbol{\theta}(t))$.*
- *The expression in the $\boldsymbol{\eta}$ parameters is $\dot{\boldsymbol{\eta}}(t) = \tilde{\nabla} \check{F}(\boldsymbol{\eta}(t))$.*

The cases of gradient computation we have discussed above are just a special case of a generic argument. Let us briefly study the gradient flow in a general chart $f: \boldsymbol{\zeta} \mapsto p_{\boldsymbol{\zeta}}$. Consider the change of parametrization from $\boldsymbol{\zeta}$ to $\boldsymbol{\theta}$,

$$\boldsymbol{\zeta} \mapsto p_{\boldsymbol{\zeta}} \mapsto \boldsymbol{\theta}(p_{\boldsymbol{\zeta}}) = I(p_{\boldsymbol{\zeta}})^{-1} \text{Cov}_{p_{\boldsymbol{\zeta}}}(\mathbf{T}, \log p_{\boldsymbol{\zeta}}) ,$$

and denote the Jacobian matrix of the parameters' change by $J(\boldsymbol{\zeta})$. We have

$$\begin{aligned} \log p_{\boldsymbol{\zeta}} &= \mathbf{T} \cdot \boldsymbol{\theta}(\boldsymbol{\zeta}) - \psi(\boldsymbol{\theta}(\boldsymbol{\zeta})) \\ &= \mathbf{T} \cdot I(p_{\boldsymbol{\zeta}})^{-1} \text{Cov}_{p_{\boldsymbol{\zeta}}}(\mathbf{T}, \log p_{\boldsymbol{\zeta}}) - \psi(I(p_{\boldsymbol{\zeta}})^{-1} \text{Cov}_{p_{\boldsymbol{\zeta}}}(\mathbf{T}, \log p_{\boldsymbol{\zeta}})) , \end{aligned}$$

and the $\boldsymbol{\zeta}$ -coordinate basis of the tangent space $T_{p_{\boldsymbol{\zeta}}}\mathcal{E}$ consists of the components of the gradient with respect to $\boldsymbol{\zeta}$,

$$\nabla(\zeta \mapsto \log p_\zeta) = J^{-1}(\zeta) (\mathbf{T} - \mathbb{E}_{p_\zeta} [\mathbf{T}])$$

It should be noted that in this case the expression Fisher information matrix does not have the form of an Hessian of a potential function. In fact, the case of the exponential and the mixture parameters point to a special structure which is called *Hessian manifold*, see [29].

3.3. Gradient flow in the mixture geometry. From now on, we are going to focus on the expression of the gradient flow in the $\boldsymbol{\eta}$ parameters. From Def. 18 we have

$$\tilde{\nabla} \check{F}(\boldsymbol{\eta}) = \text{Hess } \phi(\boldsymbol{\eta})^{-1} \nabla \check{F}(\boldsymbol{\eta}) = \text{Hess } \psi(\nabla \phi(\boldsymbol{\eta})) \nabla \check{F}(\boldsymbol{\eta}) = I(p_\boldsymbol{\eta}) \nabla \check{F}(\boldsymbol{\eta}) ,$$

where $I(p) = \text{Cov}_p(\mathbf{T}, \mathbf{T})$. As $p \mapsto \text{Cov}_p(\mathbf{T}, \mathbf{T})$ is the restriction to the simplex of a quadratic function, while $p \mapsto \boldsymbol{\eta}$ is the restriction to the exponential family \mathcal{E} of a linear function, in some cases we can naturally consider the extension of the gradient flow equation outside M° . One notable case is when the function F is a relaxation of a non constant state space function $f: \Omega \rightarrow \mathbb{R}$, as it is defined in e.g. [15].

Proposition 20. *Let $f: \Omega \rightarrow \mathbb{R}$ and let $F(p) = \mathbb{E}_p[f]$ be its relaxation on $p \in \mathcal{E}$, It follows:*

- (1) $\nabla F(p)$ is the least square projection of f onto $T_p \mathcal{E}$, that is

$$\nabla F(p) = I(p)^{-1} \text{Cov}_p(f, \mathbf{T}) \cdot (\mathbf{T} - \mathbb{E}_p[\mathbf{T}]) .$$

- (2) *The expressions in the exponential parameters $\boldsymbol{\theta}$ are $\tilde{\nabla} \hat{F}(\boldsymbol{\theta}) = (\text{Hess } \psi(\boldsymbol{\theta}))^{-1} \text{Cov}_{\boldsymbol{\theta}}(f, \mathbf{T})$, $\nabla \hat{F}(\boldsymbol{\theta}) = \text{Cov}_{\boldsymbol{\theta}}(f, \mathbf{T})$, respectively.*
- (3) *The expressions in the mixture parameters $\boldsymbol{\eta}$ are $\tilde{\nabla} \check{F}(\boldsymbol{\eta}) = \text{Cov}_{\boldsymbol{\eta}}(f, \mathbf{T})$ and $\nabla \check{F}(\boldsymbol{\eta}) = \text{Hess } \phi(\boldsymbol{\eta}) \text{Cov}_{\boldsymbol{\eta}}(f, \mathbf{T})$, respectively.*

Proof. On a generic curve thought p with velocity V , we have

$$\left. \frac{d}{dt} \mathbb{E}_{p(t)}[f] \right|_{t=0} = \text{Cov}_p(f, V) = \langle f, V \rangle_p .$$

If $V \in T_p \mathcal{E}$ we can orthogonally project f to get $\langle \nabla F, V \rangle_p = \langle (I^{-1}(p) \text{Cov}_p(f, \mathbf{T})) \cdot (\mathbf{T} - \mathbb{E}_p[\mathbf{T}]), V \rangle_p$. □

3.4. Gradient flow of the expected value function. Let us briefly recall the behaviour of the gradient flow in the relaxation case. Let $\boldsymbol{\theta}_n, n = 1, 2, \dots$, be a minimizing sequence for \hat{F} and let \bar{p} be a limit point of the sequence $(p_{\boldsymbol{\theta}_n})_n$. It follows that \bar{p} has a defective support, in particular $\bar{p} \notin \mathcal{E}$, see [26] and [24]. For a proof along lines coherent with the present paper, see [16, Th. 1]. It is found that the support $\underline{F} \subset \Omega$ is *exposed*, that is $\mathbf{T}(\underline{F})$ is a face of the marginal polytope $M = \text{con} \{\mathbf{T}(\mathbf{x}) | \mathbf{x} \in \Omega\}$. In particular, $\mathbb{E}_{\bar{p}}[\mathbf{T}] = \bar{\boldsymbol{\eta}}$ belongs to a face of the marginal polytope M . If \mathbf{a} is the (interior) orthogonal of the face, that is $\mathbf{a} \cdot \mathbf{T}(\mathbf{x}) + b \geq 0$ for all $\mathbf{x} \in \Omega$ and $\mathbf{a} \cdot \mathbf{T}(\mathbf{x}) + b = 0$ on the exposed set, then $\mathbf{a} \cdot (\mathbf{T}(\mathbf{x}) - \bar{\boldsymbol{\eta}}) = 0$ on the face, so that $\mathbf{a} \cdot \text{Cov}_{\bar{p}}(f, \mathbf{T}) = 0$. If we extend the mapping $\boldsymbol{\eta} \mapsto \text{Cov}_{\boldsymbol{\eta}}(f, \mathbf{T})$ on the closed marginal polytope M to be the limit of the vector field of the gradient on the faces of the marginal polytope, we expect to see that such a vector field is tangent to the faces.

4. GAUSSIAN MODELS

This material is taken from a conference poster presented by G.P. at the 2014 SIS conference.

Gaussian models form an exponential family with special analytical features. The following notes are intended to underline some of these special features namely the connection between differentiation with respect to parameters and differentiation with respect to the space variables.

The geometry of the multivariate Gaussian model $N(\vec{\mu}, \Sigma)$ has been studied in detail by Skovgaard in [30], where normal densities are parameterised by the mean parameter $\vec{\mu}$ and the covariance matrix Σ , and the relevant Riemannian geometry is based on an explicit form of the Fisher information.

We review Hermite polynomials from [18, ch V]. If $Z \sim \nu_1 = N(0, 1)$ and f, g real smooth functions, then $\mathbb{E}[df(Z)g(Z)] = \mathbb{E}[f(Z)\delta g(Z)]$ where $df(x) = f'(x)$ and $\delta g = xg(x) - g'(x)$ is the Stein operator. The Hermite polynomial of order $n = 0, 1, \dots$ is $H_n(x) = \delta^n 1$ e.g., $H_0(x) = 1$, $H_1(x) = x$, $H_2(x) = x^2 - 1$, ... It follows that each H_n is a monic polynomial of degree n , $dH_n = nH_{n-1}$, $\mathbb{E}[H_n(Z)H_m(Z)] = 0$ for $n \neq m$, $\mathbb{E}[H_n(Z)^2] = n!$. The sequence $H_n/\sqrt{n!}$, $n = 0, 1, \dots$ is a complete orthonormal basis of $L^2(\nu_1)$. In dimension d , for each multi-index α , we define $H_\alpha(\vec{x}) = \prod_{i=1}^d H_{\alpha_i}(x_i)$ to get an orthogonal basis of $L^2(\nu_d)$, $\nu_d = N_d(\vec{0}, I_d)$. If we define $d^\alpha = \prod_{i=1}^d d_{x_i}^{\alpha_i}$, $\delta^\alpha = \prod_{i=1}^d \delta_{x_i}^{\alpha_i}$, we have for functions $f, g: \mathbb{R}^d \rightarrow \mathbb{R}$ and $\vec{Z} \sim \nu_d$ that $\mathbb{E}[d^\alpha f(\vec{Z})g(\vec{Z})] = \mathbb{E}[f(\vec{Z})\delta^\alpha g(\vec{Z})]$. If f is infinitely differentiable, then its Fourier-Hermite expansion is $\sum_\alpha \mathbb{E}[d^\alpha f(\vec{Z})] H_\alpha(\vec{Z})/\alpha!$. Sometimes it is convenient to use $\tilde{H}_\alpha = H_\alpha/\alpha!$.

4.1. Gaussian model in the Hermite basis. Given a vector of means $\boldsymbol{\mu} \in \mathbb{R}^m$ and a full-rank covariance matrix $\Sigma \in S_m^+$, with $\Sigma = [\sigma_{ij}]$ and $\Sigma^{-1} = [\sigma^{ij}]$, the exponent $-(1/2)(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$ in the Gaussian density $N(\boldsymbol{\mu}, \Sigma)$ can be written

$$-\frac{1}{2} \left(\boldsymbol{\mu}^t \Sigma^{-1} \boldsymbol{\mu} + \text{Tr}(\Sigma^{-1}) + 2 \sum_i (\boldsymbol{\mu}^t \Sigma^{-1})_i H_i(\mathbf{x}) + 2 \sum_{i < j} \sigma^{ij} H_{ij}(\mathbf{x}) + \sum_i \sigma^{ii} H_{ii}(\mathbf{x}) \right),$$

where $H_i(\mathbf{x}) = H_1(x_i) = x_i$ and $H_{ii}(\mathbf{x}) = H_2(x_i) = x_i^2 - 1$ for $i = 1, \dots, m$, and $H_{ij} = H_1(x_i)H_1(x_j) = x_i x_j$ for $1 \leq i < j \leq m$. The likelihood of $N(\boldsymbol{\mu}, \Sigma)$ with respect to the standard Gaussian with density $w(\mathbf{x}) = (2\pi)^{-1/2} \exp(-\mathbf{x}^t \mathbf{x}/2)$ has exponent

$$-\frac{1}{2} \boldsymbol{\mu}^t \Sigma^{-1} \boldsymbol{\mu} - \frac{1}{2} \text{Tr}(\Sigma^{-1}) - \frac{m}{2} + \sum_i (\boldsymbol{\mu}^t \Sigma^{-1})_i H_i(\mathbf{x}) - \sum_{i < j} \sigma^{ij} H_{ij}(\mathbf{x}) - \sum_i (\sigma^{ii} - 1) \frac{H_{ii}(\mathbf{x})}{2}$$

Vice-versa, given $I - \Theta \in S_m^+$ and $\boldsymbol{\theta} \in \mathbb{R}^n$, then

$$(14) \quad p(\mathbf{x}; \theta_i, \theta_{ij}: i \leq j) = \exp \left(\sum_i \theta_i H_i(\mathbf{x}) + \sum_{i < j} \theta_{ij} H_{ij}(\mathbf{x}) + \sum_i \theta_{ii} \frac{H_{ii}(\mathbf{x})}{2} - \psi(\theta_i, \theta_{ij}: i \leq j) \right) w(\mathbf{x})$$

is the multivariate Gaussian density with $\Sigma^{-1}\boldsymbol{\mu} = \boldsymbol{\theta} = (\theta_i: i = 1, \dots, n)$, $I - \Sigma^{-1} = \Theta$ with upper entries $(\theta_{ij}: i < j)$, and cumulant generating function

$$(15) \quad \psi(\theta_i, \theta_{ij}: i \leq j) = \frac{1}{2}\boldsymbol{\theta}^t(I - \Theta)^{-1}\boldsymbol{\theta} - \frac{1}{2}\text{Tr}(\Theta) - \frac{1}{2}\log \det(I - \Theta).$$

In Eq. (14) the Gaussian model is presented as an exponential family with natural parameters $(\theta_i: i = 1, \dots, m; \theta_{ij}: 1 \leq i \leq j \leq m)$ in the open convex set $\mathbb{R}^n \times (I + S_m^-)$ and w -orthogonal sufficient statistics. From $(\partial/\partial\theta_i)\boldsymbol{\theta} = \mathbf{e}_i$, $(\partial/\partial\theta_{ij})\Theta = E^{ij}$ and Eq. (15) we can compute the first derivatives of the cumulant generating function ψ , that is the expected values of the sufficient statistics,

$$(16) \quad \frac{\partial}{\partial\theta_i}\psi = \boldsymbol{\theta}^t(I - \Theta)^{-1}\mathbf{e}_i = \mu_i,$$

$$(17) \quad \begin{aligned} \frac{\partial}{\partial\theta_{ij}}\psi &= \frac{1}{2}\text{Tr}((I - \Theta)^{-1}E^{ij}) + \frac{1}{2}\boldsymbol{\theta}^t(I - \Theta)^{-1}E^{ij}(I - \Theta)^{-1}\boldsymbol{\theta} \\ &= \sigma_{ij} + \mu_i\mu_j, \quad i < j, \end{aligned}$$

$$(18) \quad \begin{aligned} \frac{\partial}{\partial\theta_{ii}}\psi &= \frac{1}{2}\text{Tr}((I - \Theta)^{-1}E^{ii}) + \frac{1}{2}\boldsymbol{\theta}^t(I - \Theta)^{-1}E^{ii}(I - \Theta)^{-1}\boldsymbol{\theta} - \frac{1}{2} \\ &= \frac{1}{2}(\sigma_{ii} + \mu_i^2 - 1). \end{aligned}$$

The second derivatives, that is the covariances of the the sufficient statistics, are

$$\frac{\partial^2}{\partial\theta_i\partial\theta_j}\psi = \mathbf{e}_j^t(I - \Theta)^{-1}\mathbf{e}_i, \quad \frac{\partial^2}{\partial\theta_i\partial\theta_{jh}}\psi = \boldsymbol{\theta}^t(I - \Theta)^{-1}E^{jk}(I - \Theta)^{-1}\mathbf{e}_i,$$

and

$$\begin{aligned} \frac{\partial^2}{\partial\theta_{ij}\partial\theta_{hk}}\psi &= \frac{1}{2}\text{Tr}((I - \Theta)^{-1}E^{hk}(I - \Theta)^{-1}E^{ij}) \\ &\quad + \frac{1}{2}\boldsymbol{\theta}^t(I - \Theta)^{-1}E^{hk}(I - \Theta)^{-1}E^{ij}(I - \Theta)^{-1}\boldsymbol{\theta} \\ &\quad + \frac{1}{2}\boldsymbol{\theta}^t(I - \Theta)^{-1}E^{ij}(I - \Theta)^{-1}E^{hk}(I - \Theta)^{-1}\boldsymbol{\theta}. \end{aligned}$$

This formulæ are to be compared with the expression of the Riemannian metric in [30]. Yo Sheena [28] has a different parameterisation in which the Fisher matrix is diagonal.

We have used up to now standard change-of-parameter computations. We turn now to exploit specific properties of the Hermite system. Let us write $U(\mathbf{x}; \boldsymbol{\theta}, \Theta) = \sum_i \theta_i H_i(\mathbf{x}) + \sum_{i < j} \theta_{ij} H_{ij}(\mathbf{x}) + \sum_i \theta_{ii} \frac{H_{ii}(\mathbf{x})}{2}$. The vector space generated by the sufficient statistics $\text{Span}(U_{\boldsymbol{\theta}, \Theta} | \boldsymbol{\theta}, \Theta)$ is the space of polynomials up to degree 2 in the variables X_1, \dots, X_n that are centered with respect to w . In the geometrical picture, it is the tangent space at w of the Gaussian model, while the tangent space at $p_{\boldsymbol{\theta}, \Theta}$ is generated by the Fisher's scores, i.e. the partial derivatives of the log-density, see the discussion in [21]. We have

$$\begin{aligned} \partial U(\mathbf{x}; \boldsymbol{\theta}, \Theta) / \partial x_i &= \\ \frac{\partial}{\partial x_i} &\left(\theta_i H_1(x_i) + H_1(x_i) \sum_{j < i} \theta_{ji} H_1(x_j) + \frac{1}{2} \theta_{ii} H_2(x_i) + H_1(x_i) \sum_{i < j} \theta_{ij} H_1(x_j) \right) \\ &= \theta_i + \sum_{j < i} \theta_{ji} H_1(x_j) + \theta_{ii} H_1(x_i) + \sum_{i < j} \theta_{ij} H_1(x_j) \end{aligned}$$

and $\partial^2 U(\mathbf{x}; \boldsymbol{\theta}, \Theta) / \partial x_i \partial x_j = \theta_{ij}$. In matrix form, the basic relation between parameters of the Gaussian model and Hermite polynomials is

$$(19) \quad \nabla_{\mathbf{x}} U(\mathbf{x}; \boldsymbol{\theta}, \Theta) = \boldsymbol{\theta} + \Theta \mathbf{x}, \quad \text{Hess}_{\mathbf{x}} U(\mathbf{x}; \boldsymbol{\theta}, \Theta) = \Theta.$$

Let us write the expectation parameters as $\eta_i = \mathbb{E}_{\boldsymbol{\theta}, \Theta} [H_i]$, $\eta_{ij} = \mathbb{E}_{\boldsymbol{\theta}, \Theta} [H_{ij}]$, $i < j$, $\eta_{ii} = \mathbb{E}_{\boldsymbol{\theta}, \Theta} [H_{ii}] / 2$, and $\mathbb{E}_{\mathbf{0}, I} = \mathbb{E}$. We can compute the η 's as moments, instead of derivatives of the cumulant generating function. From $H_i = \delta_i 1$,

$$\begin{aligned} \eta_i &= \mathbb{E} [H_i e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)}] = \mathbb{E} [\partial_i e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)}] \\ &= \mathbb{E} \left[\left(\theta_i + \sum_j \theta_{ij} H_j \right) e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)} \right] = \theta_i + \sum_j \theta_{ij} \eta_j, \end{aligned}$$

or $\boldsymbol{\eta} = \boldsymbol{\theta} + \Theta \boldsymbol{\eta}$, $\boldsymbol{\eta} = (I - \Theta)^{-1} \boldsymbol{\theta}$, cfr. Eq. (16). For η_{ij} we need

$$\begin{aligned} \partial_i \partial_j e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)} &= \left(\theta_{ij} + \left(\theta_i + \sum_h \theta_{ih} H_h \right) \left(\theta_j + \sum_k \theta_{jk} H_k \right) \right) e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)} \\ &= \left(\theta_{ij} + \theta_i \theta_j + \sum_h (\theta_i \theta_{jh} + \theta_j \theta_{ih}) H_h + \sum_{h,k} \theta_{ik} \theta_{jh} H_h H_k \right) e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)}. \end{aligned}$$

From $H_{ij} = \delta^i \delta_j 1$, $\mathbb{E}_{\boldsymbol{\theta}, \Theta} [H_h H_k] = \eta_{hk}$ if $h \neq k$, $\mathbb{E}_{\boldsymbol{\theta}, \Theta} [H_h^2] = 2\eta_{hh} + 1$, we obtain

$$\eta_{ij} = \theta_{ij} + \theta_i \theta_j + \sum_h (\theta_i \theta_{jh} + \theta_j \theta_{ih}) \eta_h + \sum_{h \neq k} \theta_{ik} \theta_{jh} \eta_{hk} + \sum_h \theta_{ih} \theta_{jh} (2\eta_{hh} + 1),$$

to be compared with Eqs. (17) and (18).

4.2. Optimisation. Let $f: \mathbb{R}^m \rightarrow R$ be a continuous bounded function, with maximum at a point $\mathbf{m} \in \mathbb{R}^m$. We define the *relaxed function* $F(\boldsymbol{\theta}, \Theta) = \mathbb{E}_{\boldsymbol{\theta}, \Theta} [f] = \mathbb{E} [f e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)}]$. Then $F(\boldsymbol{\theta}, \Theta) \leq f(\mathbf{m})$ and for each sequence $(\boldsymbol{\theta}_n, \Theta_n)$, $n = 1, 2, \dots$, such that $\lim_{n \rightarrow \infty} (I - \Theta_n)^{-1} = \lim_{n \rightarrow \infty} \Sigma_n = 0$ and $\lim_{n \rightarrow \infty} (I - \Theta_n)^{-1} \boldsymbol{\theta}_n = \lim_{n \rightarrow \infty} \boldsymbol{\mu}_n = \mathbf{m}$, we have $\lim_{n \rightarrow \infty} F(\boldsymbol{\theta}_n, \Theta_n) = f(\mathbf{m})$. This remark has been used in Optimisation when the function f is a *black box* that is when no analytic expression is known, but the function can be computed at each point \mathbf{x} , see for example [19]. In fact, the gradient of the relaxed function has components

$$\begin{aligned} \frac{\partial}{\partial \theta_i} F &= \text{Cov}_{\boldsymbol{\theta}, \Theta} (f, H_i), \\ \frac{\partial}{\partial \theta_{ij}} F &= \text{Cov}_{\boldsymbol{\theta}, \Theta} (f, H_{ij}), \quad i < j, \\ \frac{\partial}{\partial \theta_{ii}} F &= \frac{1}{2} \text{Cov}_{\boldsymbol{\theta}, \Theta} (f, H_{ii}), \end{aligned}$$

so that the direction of steepest ascent at $(\boldsymbol{\theta}, \Theta)$ can be learned from a sample of $e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)} v$ for example from sample covariances. This method does not require any smoothness in the original function and it is expected to have a better robustness vs local maxima than the ordinary gradient search because a mean values of the function f are used. A reduction of dimensionality is obtained by considering sub-models, for example Θ diagonal.

We note that the gradient of the relaxed function is related with the f , ∇f , Hess f as follows. We have $\text{Cov}_{\boldsymbol{\theta}, \Theta}(f, H_i) = \mathbb{E}_{\boldsymbol{\theta}, \Theta}[f H_i] - \mathbb{E}_{\boldsymbol{\theta}, \Theta}[f] \mathbb{E}_{\boldsymbol{\theta}, \Theta}[H_i]$ and

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\theta}, \Theta}[f H_i] &= \mathbb{E}[H_i f e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)}] = \mathbb{E}[\partial_i (f e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)})] \\ &= \mathbb{E}[(\partial_i f + f \partial_i U_{\boldsymbol{\theta}, \Theta}) e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)}] = \mathbb{E}_{\boldsymbol{\theta}, \Theta}[\partial_i f] + \theta_i \mathbb{E}_{\boldsymbol{\theta}, \Theta}[f] + \sum_j \theta_{ij} \mathbb{E}_{\boldsymbol{\theta}, \Theta}[f H_j]. \end{aligned}$$

If \mathbf{H}_1 is the vector with components H_1, \dots, H_m ,

$$\mathbb{E}_{\boldsymbol{\theta}, \Theta}[f \mathbf{H}_1] = \mathbb{E}_{\boldsymbol{\theta}, \Theta}[(I - \Theta)^{-1} \nabla f] + \mathbb{E}_{\boldsymbol{\theta}, \Theta}[f] (I - \Theta)^{-1} \boldsymbol{\theta},$$

so that $\nabla_{\boldsymbol{\theta}} F = \mathbb{E}_{\boldsymbol{\theta}, \Theta}[(I - \Theta)^{-1} \nabla f] = \mathbb{E}_{\boldsymbol{\mu}, \Sigma}[\Sigma \nabla f]$.

In a similar way, $\text{Cov}_{\boldsymbol{\theta}, \Theta}(f, H_{ij}) = \mathbb{E}_{\boldsymbol{\theta}, \Theta}[f H_{ij}] - \mathbb{E}_{\boldsymbol{\theta}, \Theta}[f] \mathbb{E}_{\boldsymbol{\theta}, \Theta}[H_{ij}]$ and

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\theta}, \Theta}[f H_{ij}] &= \mathbb{E}[H_{ij} f e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)}] = \mathbb{E}[\partial_i \partial_j (f e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)})] \\ &= \mathbb{E}[\partial_i [(\partial_j f + f \partial_j U_{\boldsymbol{\theta}, \Theta}) e^{U_{\boldsymbol{\theta}, \Theta} - \psi(\boldsymbol{\theta}, \Theta)}]] \\ &= \mathbb{E}_{\boldsymbol{\theta}, \Theta}[\partial_i \partial_j f + \partial_i f \partial_j U_{\boldsymbol{\theta}, \Theta} + \partial_j f \partial_i U_{\boldsymbol{\theta}, \Theta} + f \partial_i \partial_j U_{\boldsymbol{\theta}, \Theta} + \partial_i U_{\boldsymbol{\theta}, \Theta} \partial_j U_{\boldsymbol{\theta}, \Theta}]. \end{aligned}$$

Now we can substitute in the equation above $\partial_i U_{\boldsymbol{\theta}, \Theta} = \theta_i + \sum_h \theta_{ih} H_h$, $\partial_j U_{\boldsymbol{\theta}, \Theta} = \theta_j + \sum_h \theta_{jh} H_h$, $\partial_i \partial_j U_{\boldsymbol{\theta}, \Theta} = \theta_{ij}$ and

$$\begin{aligned} \partial_i U_{\boldsymbol{\theta}, \Theta} \partial_j U_{\boldsymbol{\theta}, \Theta} &= (\theta_i + \sum_h \theta_{ih} H_h)(\theta_j + \sum_k \theta_{jk} H_k) \\ &= \theta_i \theta_j + \sum_h (\theta_i \theta_{jh} + \theta_j \theta_{ih}) H_h + \sum_{h, k} \theta_{ih} \theta_{jk} H_h H_k \\ &= \theta_i \theta_j + \sum_h (\theta_i \theta_{jh} + \theta_j \theta_{ih}) H_h + 2 \sum_{h < k} \theta_{ih} \theta_{jk} H_{hk} + \sum_h \theta_{ih} \theta_{jh} (H_{hh} + 1), \end{aligned}$$

to obtain the required relation. We leave the rest of the computation to the reader.

REFERENCES

1. Shun-ichi Amari, *Differential-geometrical methods in statistics*, Lecture Notes in Statistics, vol. 28, Springer-Verlag, New York, 1985. MR 86m:62053
2. Shun-ichi Amari, *Natural gradient works efficiently in learning*, *Neural Computation* **10** (1998), no. 2, 251–276.
3. Shun-ichi Amari and Hiroshi Nagaoka, *Methods of information geometry*, American Mathematical Society, Providence, RI, 2000, Translated from the 1993 Japanese original by Daishi Harada. MR 1 800 071
4. Lawrence D. Brown, *Fundamentals of statistical exponential families with applications in statistical decision theory*, IMS Lecture Notes. Monograph Series, no. 9, Institute of Mathematical Statistics, 1986. MR MR882001 (88h:62018)
5. A. P. Dawid, *Discussion of a paper by Bradley Efron*, *Ann. Statist.* **3** (1975), no. 6, 1231–1234.
6. ———, *Further comments on: “Some comments on a paper by Bradley Efron”* (*Ann. Statist.* **3** (1975), 1189–1242), *Ann. Statist.* **5** (1977), no. 6, 1249. MR MR0471125 (57 #10863)
7. Manfredo Perdigão do Carmo, *Riemannian geometry*, Mathematics: Theory & Applications, Birkhäuser Boston Inc., Boston, MA, 1992, Translated from the second Portuguese edition by Francis Flaherty. MR 1138207 (92i:53001)
8. Bradley Efron, *Defining the curvature of a statistical problem (with applications to second order efficiency)*, *Ann. Statist.* **3** (1975), no. 6, 1189–1242, With a discussion by C. R. Rao, Don A. Pierce, D. R. Cox, D. V. Lindley, Lucien LeCam, J. K. Ghosh, J. Pfanzagl, Niels Keiding, A. P. Dawid, Jim Reeds and with a reply by the author. MR MR0428531 (55 #1552)
9. ———, *The geometry of exponential families*, *Ann. Statist.* **6** (1978), no. 2, 362–376. MR 57 #10890
10. Robert E. Kass and Paul W. Vos, *Geometrical foundations of asymptotic inference*, Wiley Series in Probability and Statistics: Probability and Statistics, John Wiley & Sons, Inc., New York, 1997, A Wiley-Interscience Publication. MR 1461540 (99b:62032)
11. Serge Lang, *Differential and Riemannian manifolds*, third ed., Graduate Texts in Mathematics, vol. 160, Springer-Verlag, New York, 1995. MR 96d:53001
12. Steffen L. Lauritzen, *Statistical manifolds*, Differential geometry in statistical inference, Institute of Mathematical Statistics Lecture Notes—Monograph Series, 10, Institute of Mathematical Statistics, Hayward, CA, 1987.
13. Hông Vân Lê, *The uniqueness of the Fisher metric as information metric*, ArXiv.1306.1465, 2014.
14. Bertrand Lods and Giovanni Pistone, *Information geometry formalism for the spatially homogeneous Boltzmann equation*, arXiv:1502.06774, 2015.
15. Luigi Malagò, Matteo Matteucci, and Giovanni Pistone, *Towards the geometry of estimation of distribution algorithms based on the exponential family*, Proceedings of the 11th workshop on Foundations of genetic algorithms (New York, NY, USA), FOGA '11, ACM, 2011, pp. 230–242.
16. Luigi Malagò and Giovanni Pistone, *A note on the border of an exponential family*, arXiv:1012.0637v1, 2010.
17. ———, *Combinatorial optimization with information geometry: Newton method*, *Entropy* **16** (2014), 4260–4289.
18. Paul Malliavin, *Integration and probability*, Graduate Texts in Mathematics, vol. 157, Springer-Verlag, New York, 1995, With the collaboration of Hélène Airault, Leslie Kay and Gérard Letac, Edited and translated from the French by Kay, With a foreword by Mark Pinsky. MR MR1335234 (97f:28001a)
19. Y. Ollivier, L. Arnold, A. Auger, and N. Hansen, *Information-Geometric Optimization Algorithms: A Unifying Picture via Invariance Principles*, arXiv:1106.3708, 2011v1; 2013v2.
20. Giovanni Pistone, *Examples of the application of nonparametric information geometry to statistical physics*, *Entropy* **15** (2013), no. 10, 4042–4065. MR 3130268
21. ———, *Nonparametric information geometry*, Geometric science of information (Frank Nielsen and Frédéric Barbaresco, eds.), Lecture Notes in Comput. Sci., vol. 8085, Springer, Heidelberg, 2013, First International Conference, GSI 2013 Paris, France, August 28–30, 2013 Proceedings, pp. 5–36. MR 3126029
22. Giovanni Pistone and Maria Piera Rogantin, *The gradient flow of the polarization measure. with an appendix*, arXiv:1502.06718, 2015.
23. C. Radhakrishna Rao, *Information and the accuracy attainable in the estimation of statistical parameters*, *Bull. Calcutta Math. Soc.* **37** (1945), 81–91. MR MR0015748 (7,464a)

24. Johannes Rauh, Thomas Kahle, and Nihat Ay, *Support sets in exponential families and oriented matroid theory*, International Journal of Approximate Reasoning **52** (2011), no. 2, 613–626, Special Issue Workshop on Uncertainty Processing WUPES09.
25. Marta Reynal-Querol, *Ethnicity, political systems and civil war*, Journal of Conflict Resolution **46** (2002), no. 1, 29–54.
26. Alessandro Rinaldo, Stephen E. Fienberg, and Yi Zhou, *On the geometry of discrete exponential families with application to exponential random graph models*, Electron. J. Stat. **3** (2009), 446–484. MR 2507456 (2010f:62077)
27. R. Tyrrell Rockafellar, *Convex analysis*, Princeton Mathematical Series, No. 28, Princeton University Press, Princeton, N.J., 1970. MR MR0274683 (43 #445)
28. Yo Sheena, *Inference on the eigenvalues of the covariance matrix of a multivariate normal distribution—geometrical view—*, ArXiv1211.5733, November 2012.
29. Hirohiko Shima, *The geometry of Hessian structures*, World Scientific Publishing Co. Pte. Ltd., Hackensack, NJ, 2007. MR 2293045 (2008f:53011)
30. Lene Theil Skovgaard, *A Riemannian geometry of the multivariate normal model*, Scand. J. Statist. **11** (1984), no. 4, 211–223. MR 793171 (86m:62102)

L. MALAGÒ: DEPARTMENT OF ELECTRICAL AND ELECTRONIC ENGINEERING, SHINSHU UNIVERSITY, 4-17-1 WAKASATO, NAGANO 380-8553 JAPAN

E-mail address: malago@shinshu-u.ac.jp

URL: <http://homes.di.unimi.it/~malago/>

G. PISTONE: DE CASTRO STATISTICS, COLLEGIO CARLO ALBERTO, VIA REAL COLLEGIO 30, 10024 MONCALIERI, ITALY

E-mail address: giovanni.pistone@carloalberto.org

URL: www.giannidiorestino.it