

GSI2015 2nd conference on Geometric Science of Information  
28-30 Oct 2015 Ecole Polytechnique Paris-Saclay

## Second-order Optimization over the Multivariate Gaussian Distribution



Luigi Malagò <sup>1</sup>    <sup>2</sup> Giovanni Pistone



<sup>1</sup>Shinshu University JP & INRIA Saclay FR

<sup>2</sup>de Castro Statistics, Collegio Carlo Alberto, Moncalieri IT

# Introduction

- This is the presentation by Giovanni of the paper with the same title in the Proceedings.
- Unfortunately, Giovanni is the least qualified of the two authors to present this specific application of Information Geometry, his specific field of expertise being non-parametric Information Geometry and its applications in Probability and Statistical Physics. Luigi is currently working in Japan and could not make it.
- Among the two of us, Luigi is the responsible for the idea of using gradient methods and later, Newton methods, in black box optimization. Our collaboration started with the preparation of the FOGA 2011 paper
- L. Malagò, M. Matteucci, and G. Pistone. [Towards the geometry of estimation of distribution algorithms based on the exponential family.](#)

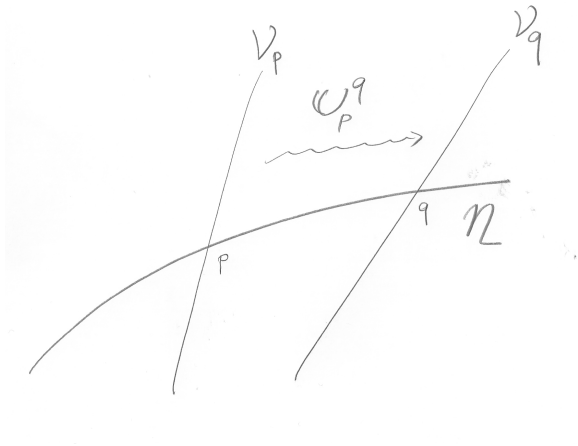
In *Proceedings of the 11th workshop on Foundations of genetic algorithms*, FOGA '11, pages 230–242, New York, NY, USA, 2011. ACM

# Summary

1. Geometry of the Exponential Family
2. Second-Order Optimization: The Newton Method
3. Applications to the Gaussian Distribution
4. Discussion and Future Work

- 
- A short introduction for Taylor formulæ on Gaussian exponential families is provided. The binary case has been previously discussed in
  - L. Malagò and G. Pistone. [Combinatorial optimization with information geometry: Newton method.](#) *Entropy*, 16:4260–4289, 2014.
  - Riemannian Newton methods are discussed in a Session of this Conference cf,
  - P.-A. Absil, R. Mahony, and R. Sepulchre. [Optimization algorithms on matrix manifolds.](#) Princeton University Press, Princeton, NJ, 2008.  
[With a foreword by Paul Van Dooren](#)
  - The focus of this short presentation is on a specific framework for Information Geometry we call *statistical bundle*.

# Hilbert vs Tangent vs Statistical Bundle



- S. Amari. [Dual connections on the Hilbert bundles of statistical models.](#)  
In *Geometrization of statistical theory* (Lancaster, 1987), pages 123–151, Lancaster, 1987. ULDM Publ
- R. E. Kass and P. W. Vos. [Geometrical foundations of asymptotic inference.](#)  
Wiley Series in Probability and Statistics: Probability and Statistics. John Wiley & Sons, Inc., New York, 1997.

## Statistical Bundle: Gaussian case

- $H_\alpha(\mathbf{x})$ ,  $\mathbf{x} \in \mathbb{R}^m$ , are Hermite polynomials of order 1 and 2.
- E.g,  $m = 3$ ,  $H_{010}(\mathbf{x}) = x_2$ ,  $H_{011}(\mathbf{x}) = x_2x_3$ ,  $H_{020}(\mathbf{x}) = x_2^2 - 1$ .
- The Gaussian model with sufficient statistics  $\mathcal{B} = \{X_1, \dots, X_n\} \subset \{H_\alpha \mid |\alpha| = 1, 2\}$ , is

$$\mathcal{N} = \left\{ p(\mathbf{x}; \boldsymbol{\theta}) = \exp \left( \sum_{j=1}^n \theta_j X_j - \psi(\boldsymbol{\theta}) \right) \right\}$$

- The *fibers* are  $\mathcal{V}_p = \text{Span}(X_j - \mathbb{E}_p[X_j] \mid j = 1, \dots, n)$
- The *statistical bundle* is

$$S\mathcal{N} = \{(p, U) \mid p \in \mathcal{N}, U \in \mathcal{V}_p\}$$

- Each  $U \in \mathcal{V}_p$ ,  $p \in \mathcal{N}$ , is a polynomial of degree up to 2 and  $t \mapsto \mathbb{E}_q[e^{tU}]$  is finite around 0,  $q \in \mathcal{N}$
- Every polynomial  $X$  belongs to  $\bigcap_{q \in \mathcal{N}} L^2(q)$

# Parallel transports

## Definition

- *e-transport*:

$${}^e\mathbb{U}_p^q: \mathcal{V}_p \ni U \mapsto U - \mathbb{E}_q[U] \in \mathcal{V}_q .$$

- *m-transport*: for each  $U \in \mathcal{V}_p$  and  $V \in \mathcal{V}_q$

$$\langle U, {}^m\mathbb{U}_q^p V \rangle_p = \langle {}^e\mathbb{U}_p^q U, V \rangle_q$$

## Properties

- ${}^e\mathbb{U}_q^r {}^e\mathbb{U}_p^q = {}^e\mathbb{U}_p^r$
- ${}^m\mathbb{U}_q^r {}^m\mathbb{U}_p^q = {}^m\mathbb{U}_p^r$
- $\langle {}^e\mathbb{U}_p^q U, {}^m\mathbb{U}_p^q V \rangle_q = \langle U, V \rangle_p$
- If  $\frac{q}{p}V \in L^2(p)$ , then  ${}^m\mathbb{U}_q^p V$  is its orthogonal projection onto  $\mathcal{V}_p$ .

## Parallel transports in coordinates I

We define on the statistical bundle  $S\mathcal{N}$  a system of *moving frames*.

1. The *exponential frame* of the fiber  $S_p\mathcal{N} = \mathcal{V}_p$  is the vector basis

$$\mathcal{B}_p = \{X_j - \mathbb{E}_p[X_j] \mid j = 1, \dots, n\}$$

2. Each element  $U \in \mathcal{V}_p$  is uniquely written as

$$U = \sum_{j=1}^n \alpha_j(U)(X_j - \mathbb{E}_p[X_j]) = \alpha(U)^T (\mathbf{X} - \mathbb{E}_p[\mathbf{X}])$$

3. The expression in the exponential frame of the scalar product is the Fisher information matrix:

$${}^e I_{ij}(p) = \langle X_i - \mathbb{E}_p[X_i], X_j - \mathbb{E}_p[X_j] \rangle_p = \text{Cov}_p(X_i, X_j) = \frac{\partial^2}{\partial \theta_i^2} \theta_j \psi(\boldsymbol{\theta})$$

- 4.

$$U \mapsto \alpha(U) = {}^e I(p)^{-1} \text{Cov}_p(\mathbf{X}, U)$$

## Parallel transports in coordinates II

5. The *mixture frame* of the fiber  $S_p\mathcal{N} = \mathcal{V}_p$  is

$${}^e l(p)^{-1} \mathcal{B}p = \left\{ \sum_{i=1}^n {}^e l^{ij}(p)(X_i - \mathbb{E}_p[X_i]) \mid j = 1, \dots, n \right\}$$

6. Each element  $V \in \mathcal{V}_p$  is uniquely written as

$$V = \sum_{j=1}^n \beta_j(V) \sum_{i=1}^n {}^e l^{ij}(p)(X_i - \mathbb{E}_p[X_i]) = \beta(V)^T {}^e l(p)^{-1}(\mathbf{X} - \mathbb{E}_p[\mathbf{X}])$$

7. The coordinates in the mixture basis are given in matrix form by

$$V \mapsto \beta(V) = \text{Cov}_p(\mathbf{X}, V) .$$

8. The matrix  ${}^m l(p) = {}^e l(p)^{-1}$  is the matrix expression of the metric in the mixture frame.

$$\alpha(U) = {}^m l(p)\beta(U), \quad \beta(U) = {}^e l(p)\alpha(U) .$$



## Parallel transports in the moving frames

- The e-transport acts on the exponential coordinates as the identity,

$$\alpha({}^e\mathbb{U}_p^q U) = \alpha(U)$$

- Equivalently,

$$= {}^e I(q)^{-1} \text{Cov}_q(\mathbf{X}, U) = {}^e I(p)^{-1} \text{Cov}_p(\mathbf{X}, U)$$

- The m-transport acts on the mixture coordinates as the identity,

$$\beta({}^m\mathbb{U}_p^q V) = \beta(V)$$

### REMARK

A *section* or *vector field* of the statistical bundle is a mapping  $F: \mathcal{N} \ni p \mapsto F(p) \in \mathcal{V}_p$ . As there are *two* distinguished charts on the model (exponential  $p \mapsto \boldsymbol{\theta}(p)$  and mixture  $p \mapsto \boldsymbol{\eta}(p) = \nabla\psi(\boldsymbol{\theta}(p))$ ) and *two* distinguished frames on each fiber, there are in general *four* distinguished expression of each section.

# Score and statistical gradient

## Definition

$t \mapsto p(t)$  is a curve in the model  $\mathcal{N}$  and  $f: \mathcal{N} \rightarrow \mathbb{R}$ .

- The *score* of the curve  $t \mapsto p(t)$  is a curve in the statistical bundle  $t \mapsto (p(t), Dp(t)) \in S\mathcal{N}$  such that for all  $X \in \text{Span}(1, X_1, \dots, X_n)$  it holds

$$\frac{d}{dt} \mathbb{E}_{p(t)} [X] = \langle X - \mathbb{E}_{p(t)} [X], Dp(t) \rangle_{p(t)}$$

- Usually,

$$Dp(t) = \frac{\dot{p}(t)}{p(t)} = \frac{d}{dt} \log p(t)$$

- The *statistical gradient* of  $f: \mathcal{N} \rightarrow \mathbb{R}$  is a section of the statistical bundle,  $p \mapsto (p, \text{grad } f(p)) \in S\mathcal{N}$  such that for each regular curve  $t \mapsto p(t)$ , it holds

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

## Score and statistical gradient in coordinates

- Let the regular curve  $t \mapsto p(t)$  be expressed in the exponential coordinates by  $t \mapsto \theta(t)$ . The score  $t \mapsto Dp(t)$  is expressed in the exponential frame by  $t \mapsto \dot{\theta}(t)$  that is,

$$Dp(t) = \sum_{j=1}^n \dot{\theta}_j(t) (X_j - \frac{\partial}{\partial \theta_j} \psi(\theta(t)))$$

- Let the regular curve  $t \mapsto p(t)$  be expressed in the mixture coordinates by  $t \mapsto \eta(t) = \nabla \psi(\theta(t))$ . The score is expressed in the mixture frame as  $t \mapsto \dot{\eta}(t)$ .
- Let  $X$  be a random variable which belongs to all  $L^2(p)$ ,  $p \in \mathcal{N}$  and  $f(p) = \mathbb{E}_p [f]$ . Then  $p \mapsto \text{grad } f(p)$  exists and equals the orthogonal projection of  $X$  onto  $\mathcal{V}_p$ , namely

$$\text{grad}(p \mapsto \mathbb{E}_p [X]) = \\ {}^e I(p)^{-1} \text{Cov}_p(\mathbf{X}, X) (\mathbf{X} - \mathbb{E}_p [\mathbf{X}]), \quad \mathbf{X} = (X_1, \dots, X_n).$$

- The expressions of  $\text{grad } f$  are of interest in optimization.

## Taylor formula in the Statistical Bundle

- For a curve  $t \mapsto p(t) \in \mathcal{N}$  connecting  $p = p(0)$  to  $q = p(1)$  and a function  $f: \mathcal{N} \rightarrow \mathbb{R}$  the Taylor formula is

$$f(q) = f(p) + \left. \frac{d}{dt} f(p(t)) \right|_{t=0} + \frac{1}{2} \left. \frac{d^2}{dt^2} f(p(t)) \right|_{t=0} + R_2(f, p, q)$$

- The first derivative is computed with the statistical gradient and the score

$$f(q) = f(p) + \langle \text{grad } f(p(0)), Dp(0) \rangle_p + \frac{1}{2} \left. \frac{d}{dt} \langle \text{grad } f(p(t)), Dp(t) \rangle_{p(t)} \right|_{t=0} + R_2(f, p, q)$$

## Acceleration and Hessian

$$\begin{aligned} \frac{d}{dt} \langle \text{grad } f(\rho(t)), D\rho(t) \rangle_{\rho(t)} \Big|_{t=0} &= \\ \frac{d}{dt} \left\langle e^{\mathbb{U}_{\rho(t)}^{\rho(0)}} \text{grad } f(\rho(t)), m^{\mathbb{U}_{\rho(t)}^{\rho(0)}} D\rho(t) \right\rangle_{\rho(0)} \Big|_{t=0} \end{aligned}$$

$$\begin{aligned} \frac{d}{dt} \langle \text{grad } f(\rho(t)), D\rho(t) \rangle_{\rho(t)} \Big|_{t=0} &= \\ \frac{d}{dt} \left\langle m^{\mathbb{U}_{\rho(t)}^{\rho(0)}} \text{grad } f(\rho(t)), m^{\mathbb{U}_{\rho(t)}^{\rho(0)}} D\rho(t) \right\rangle_{\rho(0)} \Big|_{t=0} \end{aligned}$$

$$\begin{aligned} \frac{d}{dt} \langle \text{grad } f(\rho(t)), D\rho(t) \rangle_{\rho(t)} \Big|_{t=0} &= \\ \frac{d}{dt} \mathbb{E}_{\rho(0)} \left[ \frac{\rho(t)}{\rho(0)} \text{grad } f(\rho(t)) D\rho(t) \right] \Big|_{t=0} \end{aligned}$$

# Accelerations

- Let us define the *acceleration* at  $t$  of a curve  $t \mapsto p(t) \in \mathcal{N}$ . The velocity is defined to be

$$t \mapsto (p(t), Dp(t)) = (p(t), \frac{d}{dt} \log(p(t))) \in S\mathcal{N}$$

- The *exponential acceleration* is

$${}^e D^2 p(t) = \frac{d}{ds} {}^e \mathbb{U}_{p(s)}^{p(t)} Dp(s) \Big|_{s=t}$$

- The *mixture acceleration* is

$${}^m D^2 p(t) = \frac{d}{ds} {}^m \mathbb{U}_{p(s)}^{p(t)} Dp(s) \Big|_{s=t}$$

- The *Riemannian acceleration* is

$${}^0 D^2 p(t) = \frac{1}{2} ({}^e D^2 p(t) + {}^m D^2 p(t))$$

## Covariant derivatives I

- $p \mapsto (p, F(p))$ ,  $p \mapsto (p, G(p))$ , are sections of  $SN$ , with expressions in the moving frames

$$F(p) = \sum_{j=1}^n \alpha^j(p)(X_j - \mathbb{E}_p[X_j]) ,$$

$$F(p) = \sum_{j=1}^n \beta^j(p) \sum_{i=1}^n {}^e l^{ij}(p)(X_i - \mathbb{E}_p[X_i]) ,$$

$$G(p) = \sum_{j=1}^n \gamma^j(p)(X_j - \mathbb{E}_p[X_j]) ,$$

$$G(p) = \sum_{j=1}^n \delta^j(p) \sum_{i=1}^n {}^e l^{ij}(p)(X_i - \mathbb{E}_p[X_i]) .$$

## Covariant derivatives II

- The *exponential covariant derivative* is the vector field  $p \mapsto (p, {}^e D_G F(p))$ , where

$$\begin{aligned} {}^e D_G F(p) &= \sum_{j=1}^n \langle \text{grad } \alpha^j(p), G(p) \rangle_p (X_j - \mathbb{E}_p [X_j]) \\ &= \sum_{j=1}^n \sum_{i=1}^n \gamma^i(p) (\partial_i \text{grad } \alpha^j(p)) (X_j - \mathbb{E}_p [X_j]) \end{aligned}$$

- The *mixture covariant derivative* is the vector field  $p \mapsto (p, {}^m D_G F(p))$ , where

$$\begin{aligned} {}^m D_G F(p) &= \sum_{j=1}^n \langle \text{grad } \beta^j(p), G(p) \rangle_p \sum_{i=1}^n {}^e I^{ij}(p) (X_i - \mathbb{E}_p [X_i]) = \\ &= \sum_{j=1}^n \sum_{k=1}^n \gamma^k(p) \langle \text{grad } \beta^j(p), X_k - \mathbb{E}_p [X_k] \rangle_p \sum_{i=1}^n {}^e I^{ij}(p) (X_i - \mathbb{E}_p [X_i]) \end{aligned}$$



## Covariant derivatives III

- The *Riemannian covariant derivative* is the vector field  $p \mapsto (p, {}^0D_G F(p))$  with

$${}^0D_G F = \frac{1}{2} ({}^eD_G F + {}^mD_G F) .$$

# Hessians

- Let  $f: \mathcal{N} \rightarrow \mathbb{R}$  be a mapping with gradient  $p \mapsto (p, \text{grad } f(p))$ . Let  $p \mapsto (p, G(p))$  be a vector field (section) of  $S\mathcal{N}$ .
- The *exponential Hessian* of  $f$  is the vector field  $p \mapsto (p, {}^e\text{Hess}_G f(p))$ , with

$${}^e\text{Hess}_G f(p) = {}^eD_G \text{grad } f(p) .$$

- The *mixture Hessian* of  $f$  is the vector field  $p \mapsto (p, {}^m\text{Hess}_G f(p))$ , with

$${}^m\text{Hess}_G f(p) = {}^mD_G \text{grad } f(p) .$$

- The *Riemannian Hessian* of  $F$  is the vector field  $p \mapsto (p, {}^0\text{Hess}_G F(p))$ , with

$${}^0\text{Hess}_G f(p) = {}^0D_G \text{grad } f(p) .$$

## Taylor's formulæ I

1.  $t \mapsto p(t)$  is the *mixture geodesic* connecting  $p = p(0)$  to  $q = p(1)$ .

$$f(q) = f(p) + \langle \text{grad } f(p), Dp(0) \rangle_p + \frac{1}{2} \langle {}^e\text{Hess}_{Dp(0)} f(p), Dp(0) \rangle_p + R_2^+(p, q)$$

$$R_2^+(p, q) = \int_0^1 dt \left( (1-t) \langle {}^e\text{Hess}_{Dp(t)} f(p(t)), Dp(t) \rangle_{p(t)} \right) - \frac{1}{2} \langle {}^e\text{Hess}_{Dp(0)} f(p), Dp(0) \rangle_p$$

## Taylor's formulæ II

2.  $t \mapsto p(t)$  is the *exponential geodesic* connecting  $p = p(0)$  to  $q = p(1)$ .

$$f(q) = f(p) + \langle \text{grad } f(p), Dp(0) \rangle_p + \frac{1}{2} \langle {}^m\text{Hess}_{Dp(0)} f(p), Dp(0) \rangle_p + R_2^-(p, q)$$

$$R_2^-(p, q) = \int_0^1 dt \left( (1-t) \langle {}^m\text{Hess}_{Dp(t)} f(p(t)), Dp(t) \rangle_{p(t)} \right) - \frac{1}{2} \langle {}^m\text{Hess}_{Dp(0)} f(p), Dp(0) \rangle_p$$

## Taylor's formulæ III

3.  $t \mapsto p(t)$  is the *Riemannian geodesic* connecting  $p = p(0)$  to  $q = p(1)$ .

$$f(q) = f(p) + \langle \text{grad } f(p), Dp(0) \rangle_p + \frac{1}{2} \langle {}^0\text{Hess}_{Dp(0)} f(p), Dp(0) \rangle_p + R_2^0(p, q)$$

where

$$R_2^0(p, q) = \int_0^1 dt (1-t) \langle {}^0\text{Hess}_{Dp(t)} f(p(t)), Dp(t) \rangle_{p(t)} - \frac{1}{2} \langle {}^0\text{Hess}_{Dp(0)} f(p), Dp(0) \rangle_p$$

## Newton step

- Let  $t \mapsto p(t)$  be the *exponential geodesic* starting at  $p = p(0)$  with  $Dp(0) = U$ .
- Assume  $U$  is a critical point of

$$\mathcal{V}_{p(0)} \ni U \mapsto f(p) + \langle \text{grad } f(p(0)), U \rangle_{p(0)} + \frac{1}{2} \langle {}^m\text{Hess}_U f(p), U \rangle_{p(0)}$$

that is

$$\text{grad } f(p(0)) + {}^m\text{Hess}_U f(p) = 0$$

- If  $q = p(1)$ , then

$$f(q) = f(p) - \frac{1}{2} \langle {}^0\text{Hess}_U f(p), U \rangle_p + R_2^0(p, q)$$

## Conclusion and work in progress

- Comparisons between the Riemannian Newton method e.g., Absil et al., and the statistical bundle setup are being performed.
- In particular, the use of alternative Hessians is of special interest.