

The Orlicz-Sobolev-Gauss Exponential Manifold

Giovanni Pistone
www.giannidiorestino.it



Liblice June 13 2016

My four parts

1. Amari's Information Geometry when the state space is not finite and the model is not parametric
2. An example: computing the Wasserstein's distance
3. Gauss-Orlicz-Sobolev model spaces
4. Second order geometry

*Cette conversation est dédiée à Michel Metivier,
mon maitre à Rennes (1973-75)*

Part I

Amari's Information Geometry when
the state space is not finite and the
model is not parametric

In IG the velocity is the score

- $\theta \mapsto p_\theta$ is a *curve*
- The **score** $\theta \mapsto \frac{d}{d\theta} \log p_\theta$ is an **estimating function** because $\mathbb{E}_\theta \left[\frac{d}{d\theta} \log p_\theta \right] = 0$
- Fisher-Rao computation:

$$\begin{aligned} \frac{d}{d\theta} \mathbb{E}_\theta [U] &= \int U(x) \frac{d}{d\theta} p(x; \theta) \mu(dx) \quad p(x; \theta) > 0 \\ &= \int U(x) \frac{d}{d\theta} \log p(x; \theta) p(x; \theta) \mu(dx) = \\ &= \left\langle U - \mathbb{E}_\theta [U], \frac{d}{d\theta} \log p_\theta \right\rangle_\theta \end{aligned}$$

- $U - \mathbb{E}_\theta [U]$ is the **statistical gradient** of $\theta \mapsto \mathbb{E}_\theta [U]$.
- Cf recent work by Ay, Jost, Lê, Schwachhöfer on *measure models*

IG is the geometry of the **statistical bundle**

- \mathcal{P} is a set of probabilities on a given sample space (Ω, \mathcal{F})
- For each $p \in \mathcal{P}$, $B_p \hookrightarrow L_0^1(p)$
- A *statistical bundle* is

$$T\mathcal{P} = \{(p, U) | p \in \mathcal{P}, U \in B_p\}$$

- We expect the fibers B_p to be *isomorphic* and express a tangent space at $p \in \mathcal{P}$
- A **chart at p** $\sigma_p : (q, V) \mapsto (s_p(q), \dot{s}_p(V)) \in B_p \times B_p$

- S.-i. Amari and M. Kumon. [Estimation in the presence of infinitely many nuisance parameters—geometry of estimating functions.](#)
Ann. Statist., 16(3):1044–1068, 1988
- P. Gibilisco and G. Pistone. [Connections on non-parametric statistical manifolds by Orlicz space geometry.](#)
IDAQP, 1(2):325–347, 1998
- Cf Otto, cf Lê

Fibers: $B_p = L^{(\cosh - 1)}(p)$

- The **exponential space** $L^{(\cosh - 1)}(p)$ and the **mixture space** $L^{(\cosh - 1)*}(p)$ are the Orlicz spaces respectively defined by the conjugate **Young functions**

$$(\cosh - 1)(x) = \cosh x - 1$$

- with

$$xy \leq (\cosh - 1)(x) + (\cosh - 1)_*(y)$$

- The **closed unit balls** of the exponential and mixture space are, respectively,

$$\left\{ f \mid \|f\|_{L^{(\cosh - 1)}(p)} \leq 1 \right\} = \left\{ f \mid \int (\cosh - 1)(f(x)) p(x) dx \leq 1 \right\}$$

$$\left\{ g \mid \|g\|_{L^{(\cosh - 1)*}(p)} \leq 1 \right\} = \left\{ g \mid \int (\cosh - 1)_*(g(x)) p(x) dx \leq 1 \right\}.$$

- $B_p = \{ U \in L^{(\cosh - 1)}(p) \mid \mathbb{E}_p[U] = 0 \}$ is the dual of ${}^*B_p = \{ V \in L^{(\cosh - 1)*}(p) \mid \mathbb{E}_p[V] = 0 \}$

B_p is the space of scores

B_p is exactly the space of scores of Gibbs model through p

Theorem

1. $U \in B_p$ iff $\mathbb{E}_p[U] = 0$ and $\mathbb{E}_p[(\cosh -1)(\rho U)] < \infty$ for some $\rho > 0$
2. $U \in B_p$ iff $\mathbb{E}_p[U] = 0$ and the moment generating function $\alpha \mapsto \mathbb{E}_p[e^{\theta U}]$ is finite in a neighbourhood of 0
3. The Gibbs model $\theta \mapsto \frac{e^{\theta U}}{\mathbb{E}_p[e^{\theta U}]}$ is defined in a neighborhood of 0 and $\frac{d}{d\theta} \mathbb{E}_p[e^{\theta U}]|_{\theta=0} = 0$.
4. The score of the Gibbs model at 0 is $\frac{d}{d\theta} \log p_\theta|_{\theta=0} = U$

This set-up applies to the set $\mathcal{P}_>$ of strictly positive densities.

- G. Pistone and C. Sempì. An infinite-dimensional geometric structure on the space of all the probability measures equivalent to a given one.

Ann. Statist., 23(5):1543–1561, October 1995

Isomorphism of the $L^{(\cosh^{-1})}(p)$ spaces

Theorem

$L^{(\cosh^{-1})}(p) = L^{(\cosh^{-1})}(q)$ as Banach spaces if $\theta \mapsto \int p^{1-\theta} q^\theta d\mu$ is finite on an open neighbourhood I of $[0, 1]$, i.e. It is an equivalence relation $p \sim q$ and we denote by $\mathcal{E}(p)$ the class containing p .

Proof.

Assume $U \in L^{(\cosh^{-1})}(p)$ and consider the restrictions to the axes of the convex function

$$(s, \theta) \mapsto \int e^{sU} p^{1-\theta} q^\theta d\mu = \int \exp\left(sU + \theta \log \frac{q}{p}\right) p d\mu$$

□

- G. Pistone and C. Sempì. *An infinite-dimensional geometric structure on the space of all the probability measures equivalent to a given one.*
Ann. Statist., 23(5):1543–1561, October 1995
- A. CENA. *Geometric structures on the non-parametric statistical manifold.*
PhD thesis, Dottorato in Matematica, Università di Milano, 2002

Portmanteau theorem

Theorem

The following statements are equivalent for $p, q \in \mathcal{P}_{>}$:

- $q \in \mathcal{E}(p)$;
- $p \smile q$;
- $\mathcal{E}(p) = \mathcal{E}(q)$;
- $L^{(\cosh - 1)}(p) = L^{(\cosh - 1)}(q)$;
- $\log\left(\frac{q}{p}\right) \in L^{\cosh - 1}(p) \cap L^{\cosh - 1}(q)$.
- $\frac{q}{p} \in L^{1+\epsilon}(p)$ and $\frac{p}{q} \in L^{1+\epsilon}(q)$ for some $\epsilon > 0$.

- A. Cena and G. Pistone. [Exponential statistical manifold](#).
Ann. Inst. Statist. Math., 59(1):27–56, 2007
- M. Santacroce, P. Siri, and B. Trivellato. [New results on mixture and exponential models by Orlicz spaces](#).
Bernoulli, 22(3):1431–1447, 2016

Maximal exponential family

- For each $p \in \mathcal{P}_>$, the *moment generating functional* is the positive lower-semi-continuous convex function $G_p: B_p \ni U \mapsto \mathbb{E}_p [e^U]$ and
- the *cumulant generating functional* is the non-negative lower semicontinuous convex function $K_p = \log G_p$.
- The interior of the proper domain

$$\mathcal{S}_p = \left\{ U \in L^{(\cosh^{-1})}(p) \mid G_p(U) < +\infty \right\}^\circ$$

is an open convex set containing the open unit ball of $L^{(\cosh^{-1})}(p)$.

- For each $p \in \mathcal{P}_>$, the *maximal exponential family* at p is

$$\mathcal{E}(p) = \left\{ e^{u - K_p(u)} \cdot p \mid u \in \mathcal{S}_p \right\}.$$

From now on the maximal exponential family of interest is the family of the Maxwell density on \mathbb{R}^n , $\mathcal{E}(M)$

e-chart at $p \in \mathcal{E}(M)$

- For each $p \in \mathcal{E}(M)$ we define a chart $s_p: \mathcal{E}(M) \rightarrow \mathcal{S}_p \subset B_p$.
- The chart is defined by

$$s_p(q) \mapsto \log\left(\frac{q}{p}\right) + D(p\|q) = \log\left(\frac{q}{p}\right) - \mathbb{E}_p\left[\log\left(\frac{q}{p}\right)\right]$$

- The inverse of the chart $e_p^{-1} = s_p: \mathcal{S}_p \rightarrow \mathcal{E}(M)$ is

$$e_p(U) = \exp(U - K_p(U)) \cdot p$$

- $\{s_p | p \in \mathcal{E}(M)\}$ is an affine atlas on $\mathcal{E}(M)$ that defines the *exponential manifold*
- The information closure of any $\mathcal{E}(M)$ is \mathcal{P}_{\geq} . The reverse information closure of any $\mathcal{E}(M)$ is $\mathcal{P}_{>}$.

- I. Csiszár and F. Matúš. [Information projections revisited](#).

IEEE Trans. Inform. Theory, 49(6):1474–1490, 2003

- D. Imparato and B. Trivellato. [Geometry of extended exponential models](#).

In *Algebraic and geometric methods in statistics*, pages 307–326. Cambridge Univ. Press, Cambridge, 2010

e-chart at $(p, U) \in T\mathcal{E}(M)$

- A curve $t \mapsto p(t)$, $p(0) = p$ in the exponential manifold $\mathcal{E}(M)$ is expressed in the chart s_p as $p(t) = e^{U(t) - K_p(U(t))} \cdot p$.
- The expression of the velocity at $t = 0$ is $\dot{U}(0) = \left. \frac{d}{dt} \log p(t) \right|_{t=0}$
- It follows that the **exponential bundle**

$$T\mathcal{E}(M) = \{(p, U) | p \in \mathcal{E}(M), U \in B_p\}$$

is the expression of the tangent bundle of the exponential manifold

- The transition map $s_{p_2} \circ e_{p_1} : \mathcal{S}_{p_1} \rightarrow \mathcal{S}_{p_2}$ is affine with derivative ${}^e\mathbb{U}_{p_1}^{p_2} : B_{p_1} \rightarrow B_{p_2}$ given by ${}^e\mathbb{U}_{p_1}^{p_2} U = U - \mathbb{E}_{p_2}[U]$
- We define an atlas of charts on $T\mathcal{E}(M)$ by

$$\sigma_p(q, V) = (s_p(q), {}^e\mathbb{U}_q^p V)$$

- G. Pistone. [Nonparametric information geometry](#).

In F. Nielsen and F. Barbaresco, editors, *Geometric science of information*, volume 8085 of *Lecture Notes in Comput. Sci.*, pages 5–36. Springer, Heidelberg, 2013.

First International Conference, GSI 2013 Paris, France, August 28-30, 2013 Proceedings

Cumulant functional

- The r-divergence $q \mapsto \mathbf{D}(p \| q)$ is represented in the chart centered at p by $\mathbf{D}(p \| e_p(U)) = K_p(U) = \log \mathbb{E}_p [e^U]$.
- $K_p : B_p \rightarrow \mathbb{R}_{\geq} \cup \{+\infty\}$ is convex and its proper domain contains the open unit ball of B_p . It is infinitely Gâteaux-differentiable on the interior \mathcal{S}_p of its proper domain and analytic on the unit ball of B_p .
- For all $V, V_1, V_2, V_3 \in B_p$ the first derivatives are:

$$\begin{aligned}d K_p(U)[V] &= \mathbb{E}_q [V] \\d^2 K_p(U)[V_1, V_2] &= \text{Cov}_q (V_1, V_2) \\d^3 K_p(U)[V_1, V_2, V_3] &= \text{Cov}_q (V_1, V_2, V_3)\end{aligned}$$

- G. Pistone and M. Rogantin. [The exponential statistical manifold: mean parameters, orthogonality and space transformations.](#)
Bernoulli, 5(4):721–760, August 1999
- A. Cena and G. Pistone. [Exponential statistical manifold.](#)
Ann. Inst. Statist. Math., 59(1):27–56, 2007

Pre-dual statistical bundle

- Recall $L^{(\cosh - 1)*}(M)$ is the pre-dual of $L^{(\cosh - 1)}(M)$
- Define the **pre-dual statistical bundle** with fibers ${}^*B_p = \{V \in L^{(\cosh - 1)*}(M) \mid \mathbb{E}_p[V] = 0\}$.
- Compute the adjoint of the transport ${}^e\mathbb{U}_p^q$. For $U \in B_p$ and $V \in {}^*B_q$,

$$\begin{aligned}\langle {}^e\mathbb{U}_p^q U, V \rangle_q &= \langle U - \mathbb{E}_q[U], V \rangle_q = \mathbb{E}_q[UV] \\ &= \mathbb{E}_p \left[\frac{q}{p} UV \right] = \left\langle U, \frac{q}{p} V \right\rangle_p = \langle U, {}^m\mathbb{U}_q^p V \rangle_p\end{aligned}$$

- Define the charts on ${}^*T\mathcal{E}(M)$ by

$$\sigma_p^*(q, W) = (s_p(q), {}^m\mathbb{U}_q^p W)$$

Statistical gradient

- The **score** of the curve $t \mapsto p(t)$ is a curve in the statistical bundle $t \mapsto (p(t), Dp(t)) \in T\mathcal{E}(M)$ such that for all $X \in L^{(\cosh-1)*}(M)$ it holds

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

- $Dp(t)$ is the expression in the exponential atlas of the velocity

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

- The **statistical gradient** of $F: \mathcal{E}(M) \rightarrow \mathbb{R}$ is a section of the pre-dual statistical bundle ${}^*T\mathcal{E}(M)$, $p \mapsto (p, \text{grad } F(p)) \in {}^*T\mathcal{E}(p)$ such that for each regular curve

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

- 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

- G. Pistone. [Examples of the application of nonparametric information geometry to statistical physics.](#) *Entropy*, 15(10):4042–4065, 2013

Part II

An example: computing the Wasserstein distance

- This is an example of the use of the formalism. There is considerable literature, e.g. F. Otto
- Unpublished talk at the workshop *Computational information geometry for image and signal processing* Sep 21st-25th, 2015 at ICMS, Edinburgh. Finite state space.
- Unpublished work in progress with Luigi Malagò.

Transport plan

$(\mathbb{R}^{2n}, M_{2n}) = (\mathbb{R}^n, M_n) \otimes (\mathbb{R}^n, M_n)$ with projection X and Y

- The marginalization mapping

$$\mathcal{E}(M_{2n}) \ni \gamma \mapsto (\gamma_1, \gamma_2) \in \mathcal{E}(M_n) \times \mathcal{E}(M)_n$$

has fibers

$$\Gamma(\mu_1, \mu_2) = \{\gamma \in \mathcal{E}(M_{2n}) \mid X_{\#}\gamma = \mu_1, Y_{\#}\gamma = \mu_2\}$$

which are convex subsets

- If $t \mapsto \gamma(t) \in \Gamma(\mu_1, \mu_2)$, then

$$\begin{aligned} 0 &= \frac{d}{dt} \mathbb{E}_{\gamma(t)} [f \circ X] = \langle f \circ X - \mathbb{E}_{\mu_1} [f], D\gamma(t) \rangle_{\gamma(t)} = \\ &\quad \langle f \circ X - \mathbb{E}_{\mu_1} [f], \mathbb{E}_{\gamma(t)} (D\gamma(t) | X) \rangle_{\gamma(t)} \end{aligned}$$

-

$$\mathbb{E}_{\gamma(t)} (D\gamma(t) | X) = 0, \quad \mathbb{E}_{\gamma(t)} (D\gamma(t) | Y) = 0, \quad D\gamma(t) \in B_{\gamma(t)}$$

Splitting of $T\Gamma(\mu_1, \mu_2)$

- Consider subspaces of the ANOVA

$$B_{\gamma,1} = \{f \circ X | f \in B_{\gamma_1}\},$$

$$B_{\gamma,2} = \{f \circ Y | f \in B_{\gamma_2}\},$$

$$^*B_{\gamma,12} = (B_{\gamma,0} + B_{\gamma,1} + B_{\gamma,2})^\perp$$

- For $U \in B_\gamma$, let $U = U_1 + U_2 + U_{12}$ be a splitting. Then

$$\mathbb{E}_\gamma(U - (U_1 + U_2) | X) = 0$$

$$\mathbb{E}_\gamma(U - (U_1 + U_2) | Y) = 0$$

- or in terms of transport

$$\mathbb{E}_{M_{2n}}({}^m\mathbb{U}_\gamma^{M_{2n}} U | X) - \mathbb{E}_{M_{2n}}\left(\frac{\gamma}{M_{2n}} \middle| X\right) U_1 + \mathbb{E}_{M_{2n}}({}^m\mathbb{U}_\gamma^{M_{2n}} U_2 | X) = 0$$

$$\mathbb{E}_{M_{2n}}({}^m\mathbb{U}_\gamma^{M_{2n}} U | X) - \mathbb{E}_{M_{2n}}({}^m\mathbb{U}_\gamma^{M_{2n}} U_1 | X) + \mathbb{E}_{M_{2n}}\left(\frac{\gamma}{M_{2n}} \middle| X\right) U_2 = 0$$

Gradient flow

- Given a cost function $w: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$, the *expected cost function* is $W: \mathcal{E}(M_{2n}) \ni \gamma \mapsto \mathbb{E}_\gamma[w]$. Then the function W restricted to the open plan $\Gamma(\mu_1, \mu_2) \cap \mathcal{E}(M_{2n})$ has statistical gradient obtained by the projection of the unconstrained gradient $w - \mathbb{E}_\gamma[w]$ onto the the space of interactions $B_{\gamma,12}$

$$\text{grad } W(\gamma): \gamma \mapsto w - \mathbb{E}_\gamma[w] - w_{1,\gamma} - w_{2,\gamma}$$

- The gradient flow equation is

$$D\gamma(t) = - (w - \mathbb{E}_\gamma[w] - w_{1,\gamma} - w_{2,\gamma}) .$$

- Any solution $t \mapsto \gamma(t)$ of the gradient flow converges to a measure $\gamma^* = \lim_{t \rightarrow \infty} \gamma(t)$, in $\Gamma(\mu_1, \mu_2)$ (but not in $\mathcal{E}(M_{2n})$) such that

$$\mathbb{E}_{\gamma^*}[w] = \min \{ \mathbb{E}_\gamma[w] \mid \gamma \in \Gamma(\mu_1, \mu_2) \}$$

Part III

Gauss-Orlicz-Sobolev model spaces

- M. R. Grasselli. [Dual connections in nonparametric classical information geometry.](#)
Ann. Inst. Statist. Math., 62(5):873–896, 2010
- B. Lods and G. Pistone. [Information geometry formalism for the spatially homogeneous Boltzmann equation.](#)
Entropy, 17(6):4323–4363, 2015
- D. Brigo and G. Pistone. [Projection based dimensionality reduction for measure valued evolution equations in statistical manifolds.](#)
arXiv:1601.04189, 2016
- D. Brigo and G. Pistone. [Eigenfunctions based maximum likelihood estimation of the fokker planck equation and hellinger projection.](#)
submitted, 2016
- Luigi Montrucchio and GP. Unpublished working paper (2016) based on N. Newton deformed logarithm

Inclusions

1. If $1 < a < \infty$,

$$L^\infty(M) \subset L^{(\cosh-1)}(M) \subset L^a(M) \subset L^{(\cosh-1)*}(M) \subset L^1(M)$$

2. Local inclusions hold, if $1 \leq a < \infty$, $\Omega_R = \{x \in \mathbb{R}^n \mid |x| < R\}$,

$$L^{(\cosh-1)}(M) \hookrightarrow L^a(\Omega_R), \quad L^{(\cosh-1)*}(M) \hookrightarrow L^1(\Omega_R)$$

3. The Orlicz space $L^{(\cosh-1)}(M)$ contains all functions $f \in C^2(\mathbb{R}^n; \mathbb{R})$ whose Hessian is uniformly bounded in operator's norm. In particular, it contains all polynomials with degree up to 2 and, moreover, all functions which are bounded by such a polynomial.
4. The Orlicz space $L^{(\cosh-1)*}(M)$ contains all random variables $f: \mathbb{R}^d \rightarrow \mathbb{R}$ which are bounded by a polynomial, in particular, all polynomials.

Pointwise density

The space $L^{(\cosh - 1)}(M)$ is *not* separable nor reflexive. However, we have the following monotone class theorem

- For each $f \in L^{(\cosh - 1)}(M)$ there exist a nonnegative function $h \in L^{(\cosh - 1)}(M)$ and a sequence $f_n \in C_0(\mathbb{R}^n)$ with $|f_n| \leq h$, $n = 1, 2, \dots$, such that $\lim_{n \rightarrow \infty} f_n = f$ a.s..
- The space $C_0(\mathbb{R}^n)$ is strongly dense in $L^{(\cosh - 1)*}(M)$ and it is weakly*-dense in $L^{(\cosh - 1)}(M)$.
- For each $f \in L^{(\cosh - 1)}(M)$ there exist a nonnegative function $h \in L^{(\cosh - 1)}(M)$ and a sequence $\phi_n \in C_0^\infty(\mathbb{R}^n)$ with $|\phi_n| \leq h$, $n = 1, 2, \dots$, such that $\lim_{n \rightarrow \infty} \phi_n = f$ a.s..
- The space $C_0^\infty(\mathbb{R}^n)$ is strongly dense in $L^{(\cosh - 1)*}(M)$ and it is weakly*-dense in $L^{(\cosh - 1)}(M)$.

Orlicz class

Definition

We define the *exponential (Orlicz) class*, $C_0^{(\cosh^{-1})}(M)$, to be the closure of $C_0(\mathbb{R}^n)$ in the exponential (Orlicz) space $L^{(\cosh^{-1})}(M)$.

Theorem

Assume $f \in L^{(\cosh^{-1})}(M)$ and write $f_R(x) = f(x)(|x| > R)$. The following conditions are equivalent:

1. The real function $\rho \mapsto \int (\cosh^{-1})(\rho f(x)) M(x) dx$ is finite for all $\rho > 0$.
2. $f \in C^{\cosh^{-1}}(M)$.
3. $\lim_{R \rightarrow \infty} \|f_R\|_{L^{(\cosh^{-1})}(M)} = 0$.

For example $(x \mapsto \|x\|^2) \in L^{(\cosh^{-1})}(M) \setminus C_0^{(\cosh^{-1})}(M)$

Translation models

- We look for statistical models induced by the geometry of the state space. E.g. the n -dimensional model defined by the translation of $p = e^{U - K_M(U)} \cdot M \in \mathcal{E}(M)$

$$p(x; h) = p(x - h) = e^{U(x-h) - K_M(U)} e^{h \cdot x - \frac{|h|^2}{2}} \cdot M$$

-

$$\begin{aligned} \mathbb{E}_M [U(X - h) + h \cdot X] &= \int U(x - h) M(x) dx = \\ &= \int U(x) e^{-h \cdot x - \frac{|h|^2}{2}} M(x) dx = \mathbb{E}_M \left[U e^{-h \cdot X - \frac{|h|^2}{2}} \right] \end{aligned}$$

-

$$p(x; h) = p(x - h) = \exp(U_h - K_M(U_h)) \cdot M$$

with $U_h = \tau_h U + h \cdot X - \mathbb{E}_M[\tau_h U] \in B_M$ and

$$K_M(\tau_h U) = K_M(U) - \frac{1}{2} |h|^2$$

Translations in $L^{(\cosh - 1)}(M)$

- For each $h \in \mathbb{R}^n$, the mapping $f \mapsto \tau_h f$ is linear and bounded from $L^{(\cosh - 1)}(M)$ to itself and $\|\tau_h f\|_{L^{(\cosh - 1)}(M)} \leq 2 \|f\|_{L^{(\cosh - 1)}(M)}$ if $|h| \leq \sqrt{\log 2}$.
- For each $f \in L^{(\cosh - 1)}(M)$ and $h \in \mathbb{R}^n$ we have $\tau_h f \in L^{(\cosh - 1)}(M)$. For all $g \in L^{(\cosh - 1)*}(M)$ we have

$$\langle \tau_h f, g \rangle_M = \langle f, \tau_h^* g \rangle_M, \quad \tau_h^* g(x) = e^{-h \cdot x - \frac{|h|^2}{2}} \tau_{-h} g(x),$$

and $|h| \leq \sqrt{2}$ implies $\|\tau_h^* g\|_{L^{(\cosh - 1)*}(M)} \leq 4 \|g\|_{L^{(\cosh - 1)*}(M)}$.

Moreover, $h \mapsto \tau_h^* g$ is continuous in $L^{(\cosh - 1)*}(M)$.

- If $f \in C_0^{(\cosh - 1)}(M)$ then $\tau_h f \in C_0^{(\cosh - 1)}(M)$, $h \in \mathbb{R}^n$ and the mapping $\mathbb{R}^n: h \mapsto \tau_h f$ is continuous in $L^{(\cosh - 1)}(M)$.

Translation by a probability

- Let

$$\tau_\mu f(x) = \int f(x - y) \mu(dy) = f * \mu(x)$$

for $\mu \in \mathcal{P}_e$, namely $\mathbb{E}_M \left[e^{\frac{1}{2}|h|^2} \right] < \infty$.

- The mapping $f \mapsto \tau_\mu f$ is linear and bounded from $L^{(\cosh - 1)}(M)$ to itself. If, moreover, $\int e^{|h|^2/2} \mu(dh) \leq \sqrt{2}$, then its norm is bounded by 2.
- If $f \in C_0^{(\cosh - 1)}(M)$ then $\tau_\mu f \in C_0^{(\cosh - 1)}(M)$. The mapping $\mathcal{P}: \mu \mapsto \tau_\mu f$ is continuous at δ_0 from the weak convergence to the $L^{(\cosh - 1)}(M)$ norm.

Mollifiers

- Let be given a family of mollifiers $\omega_\lambda \in C_0^\infty(\mathbb{R}^n)$. Let $f \in C_0^{(\cosh^{-1})}(M)$. For each $\lambda > 0$ the function

$$\tau_{\omega_\lambda} f(x) = \int f(x-y) \lambda^{-n} \omega(\lambda^{-1}y) dy = f * \omega_\lambda(x)$$

belongs to $C^\infty(\mathbb{R}^n)$ and $\lim_{\lambda \rightarrow 0} f * \omega_\lambda = f$ in $L^{(\cosh^{-1})}(M)$

- For each $f \in L^{(\cosh^{-1})}(M)$ there exists a sequence f_n , in $C_0^\infty(\mathbb{R}^n)$, $n \in \mathbb{N}$, and a bound $h \in L^{(\cosh^{-1})}(M)$ such that $|f_n(x)| \leq h(x)$ and $\lim_{n \rightarrow \infty} f_n(x) = f(x)$ for all x .
- For each $f \in C_0^{(\cosh^{-1})}(M)$ there exists a sequence f_n , in $C_0^\infty(\mathbb{R}^n)$, $n \in \mathbb{N}$, with $\lim_{n \rightarrow \infty} \|f_n - f\|_{L^{(\cosh^{-1})}(M)} = 0$.

Differentiable densities

Definition

The Orlicz-Sobolev (O-S) spaces with weight M are

$$W^{1,(\cosh^{-1})}(M) = \left\{ f \in L^{(\cosh^{-1})}(M) \mid \partial_j f \in L^{(\cosh^{-1})}(M), j = 1, \dots, n \right\}$$

$$W^{1,(\cosh^{-1})^*}(M) = \left\{ f \in L^{(\cosh^{-1})^*}(M) \mid \partial_j f \in L^{(\cosh^{-1})^*}(M), j = 1, \dots, n \right\}$$

where ∂_j is the derivative in the sense of distributions. The spaces $W^{1,(\cosh^{-1})}(M)$ and $W^{1,(\cosh^{-1})^*}(M)$ are both Banach spaces for the graph norms.

-

$$\langle \partial_j f, \phi M \rangle = - \langle f, \partial_j \phi M - X_j \phi M \rangle = \langle f, (X_j - \partial_j) \phi M \rangle$$

$$\langle \partial_j f, \phi \rangle_M = \langle f, \delta_j \phi \rangle_M \quad \delta_j \phi = (X_j - \partial_j) \phi$$

- Cf Malliavin Calculus

Inclusions

Theorem

Let $R > 0$ and let Ω_R denote the open sphere of radius R .

1. We have the embeddings

$$W^{1,(\cosh-1)}(\mathbb{R}^n) \subset W^{1,(\cosh-1)}(M) \subset W^{1,(\cosh-1)}(\Omega_R) \subset W^{1,p}(\Omega_R), \quad p \geq 1.$$

2. We have the embeddings

$$W^{1,p}(\mathbb{R}^n) \subset W^{1,(\cosh-1)^*}(\mathbb{R}^n) \subset W^{1,(\cosh-1)^*}(M) \subset W^{1,(\cosh-1)^*}(\Omega_R) \subset W^{1,1}(\Omega_R), \quad p > 1.$$

3. Each $u \in W^{1,(\cosh-1)}(M)$ is a.s. continuous and Hölder of all orders on each $\overline{\Omega}_R$.

Directional derivative

Theorem

- For each $f \in W^{1,(\cosh-1)}(M)$, each unit vector $h \in S^n$, and all $t \in \mathbb{R}$, it holds

$$f(x + th) - f(x) = t \int_0^1 \sum_{j=1}^n \partial_j f(x + sth) h_j ds .$$

Moreover, $|t| \leq \sqrt{2}$ implies

$$\|f(x + th) - f(x)\|_{L^{(\cosh-1)}(M)} \leq 2t \|\nabla f\|_{L^{(\cosh-1)}(M)} ,$$

especially, $\lim_{t \rightarrow 0} \|f(x + th) - f(x)\|_{L^{(\cosh-1)}(M)} = 0$ uniformly in h .

- For each $f \in W^{1,(\cosh-1)}(M)$ and each $g \in L^{(\cosh-1)*}(M)$, the mapping $h \mapsto \langle \tau_h f, g \rangle_M$ is differentiable. Viceversa, if $f \in L^{(\cosh-1)}(M)$ and $h \mapsto \tau_h f$ is weakly differentiable, then $f \in W^{1,(\cosh-1)}(M)$
- If $\partial_j f \in C_0^{(\cosh-1)}(M)$, $j = 1, \dots, n$, then strong differentiability in $L^{(\cosh-1)}(M)$ holds.

Orlicz-Sobolev class

Definition

The *Orlicz-Sobolev-Gauss exponential class* is

$$C_0^{1,(\cosh-1)}(M) = \left\{ f \in W^{1,(\cosh-1)}(M) \mid f, \partial_j f \in C_0^{(\cosh-1)}(M), j = 1, \dots, n \right\}$$

- The translation model is qualified as

$$p_h = \tau_h p, \quad p = e^{U - K_M(U)} \cdot M, \quad U \in \mathcal{S}_M \cap C_0^{1,(\cosh-1)}(M)$$

- The score in the direction j is $(x_j - te_j) - \partial_j U(x - te_j)$:

$$\begin{aligned} \frac{\frac{d}{dt} p(x - te_j)}{p(x - te_j)} &= \frac{\frac{d}{dt} e^{U(x-te_j) - K_M(U)} M(x - te_j)}{p(x - te_j)} = \\ &= \frac{-\partial_j U(x - te_j) e^{U(x-te_j) - K_M(U)} M(x - te_j) + (x_j - te_j) e^{U(x-te_j) - K_M(U)} M(x - te_j)}{p(x - te_j)} = \\ &= (x_j - te_j) - \partial_j U(x - te_j) \end{aligned}$$

Calculus in $C_0^{1,(\cosh-1)}(M)$

Theorem

- For each $f \in C_0^{1,(\cosh-1)}(M)$ the sequence $f * \omega_n$, $n \in \mathbb{N}$, belongs to $C^\infty(\mathbb{R}^n) \cap W^{1,(\cosh-1)}(M)$. Precisely, for each n and $j = 1, \dots, n$, we have the equality $\partial_j(f * \omega_n) = (\partial_j f) * \omega_n$; the sequences $f * \omega_n$, respectively $\partial_j f * \omega_n$, $j = 1, \dots, n$, converge to f , respectively $\partial_j f$, $j = 1, \dots, n$, strongly in $W^{1,(\cosh-1)*}(M)$.
- Same statement is true if $f \in W^{1,(\cosh-1)*}(M)$.
- Let be given $f \in C_0^{(\cosh-1)}(M)$ and $g \in W^{1,(\cosh-1)*}(M)$. Then $fg \in W^{1,1}(M)$ and $\partial_j(fg) = \partial_j fg + f \partial_j g$.
- Let be given $F \in C^1(\mathbb{R})$ with $\|F'\|_\infty < \infty$. For each $U \in C_0^{(\cosh-1)}(M)$, we have $F \circ U, F' \circ U \partial_j U \in C_0^{(\cosh-1)}(M)$ and $\partial_j F \circ U = F' \circ U \partial_j U$, in particular $F(U) \in C_0^{1,(\cosh-1)}(M)$.

Exponential family modeled on $C_0^{1,(\cosh-1)}(M)$

- Restrict the exponential family $\mathcal{E}(M)$ to $C_0^{1,(\cosh-1)}(M)$,

$$\mathcal{E}_1(M) = \left\{ e^{U - K_M(U)} \cdot M \mid U \in C_0^{1,(\cosh-1)}(M) \cap \mathcal{S}_M \right\}$$

- Because of $C_0^{1,(\cosh-1)}(M) \hookrightarrow L^{\cosh-1}(M)$ the domain $C_0^{1,(\cosh-1)}(M) \cap \mathcal{S}_M$ is open and the cumulant functional $K_M : C_0^{1,(\cosh-1)}(M) \cap \mathcal{S}_M \rightarrow \mathbb{R}$ remains convex and differentiable.
- Every feature of the exponential manifold carries over to this case. Define $B_1(p) = B_p \cap C_0^{1,(\cosh-1)}(M)$ to be models for the tangent spaces of $\mathcal{E}_1(M)$. The e-transport acts on these spaces

$${}^e\mathbb{U}_f^g : B_1(f) \ni U \mapsto U - \mathbb{E}_g[U] \in B_1(g) ,$$

so that we can define the statistical bundle to be

$$T\mathcal{E}_1(M) = \{(g, V) \mid g \in \mathcal{E}_1(M), V \in B_1(g)\}$$

and take as charts the restrictions of the charts defined on $T\mathcal{E}(M)$.

Application: Hyvärinen divergence

- For each $f, g \in \mathcal{E}_1(M)$ the Hyvärinen divergence is

$$\text{DH}(g|f) = \mathbb{E}_g \left[|\nabla \log f - \nabla \log g|^2 \right].$$

- The expression in the chart centered at M is

$$\text{DH}_M(v|u) = \text{DH}(e_M(v)|e_M(u)) = \mathbb{E}_M \left[|\nabla u - \nabla v|^2 e^{v-K_M(v)} \right],$$

where $f = e_M(u)$, $g = e_M(v)$.

- $\text{grad}(f \mapsto \text{DH}(g|f)) = -2\nabla \log g \cdot \nabla \log \frac{f}{g} - 2\Delta \log \frac{f}{g}$
- $\text{grad}(g \mapsto \text{DH}(f|g)) = 2\nabla \log g \cdot \nabla \log \frac{f}{g} + 2\Delta \log \frac{f}{g} + \text{DH}(f|g)$

Example: Elliptic operator

- Elliptic operator as section of the tangent bundle is

$$\mathcal{A}p(x) = p(x)^{-1} \sum_{i,j=1}^d \frac{\partial}{\partial x_i} \left(a_{ij}(x) \frac{\partial}{\partial x_j} p(x) \right), \quad x \in \mathbb{R}^d.$$

- The expression in the statistical bundle is

$$\begin{aligned} U \mapsto \widehat{\mathcal{A}}_M(U) &= e^{U-K_M(U)} \mathcal{A}(e^{U-K_M(U)} \cdot M) = \\ &= \frac{e^{U-K_M(U)}}{e^{U-K_M(U)} \cdot M} \mathcal{A}(e^{U-K_M(U)} \cdot M) = M^{-1} \mathcal{L}^*(e^{U-K_M(U)} \cdot M) \end{aligned}$$

- Computation gives

$$\begin{aligned} M^{-1} \mathcal{L}^*(e^{U-K_M(U)} \cdot M) &= \\ &= e^{U-K_M(U)} \sum_{i,j=1}^d \frac{\partial}{\partial x_i} \left[a_{ij}(x) \left(\frac{\partial}{\partial x_j} U(x) - x_j \right) \right] p(x) + \\ &= e^{U-K_M(U)} \sum_{i,j=1}^d a_{ij}(x) \left(\frac{\partial}{\partial x_i} U(x) - x_i \right) \left(\frac{\partial}{\partial x_j} U(x) - x_j \right) p(x). \end{aligned}$$

Part IV

Second order geometry

- L. Malagò and G. Pistone. [Combinatorial optimization with information geometry: Newton method.](#) *Entropy*, 16:4260–4289, 2014
- L. Malagò and G. Pistone. [Second-order optimization over the multivariate Gaussian distribution.](#) In F. Barbaresco and F. Nielsen, editors, *Geometric Science of Information*, number 9389 in LNCS, pages 349–358. Springer, 2015

Parallel transport

- *e-transport*:

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

- *m-transport*: for each $V \in {}^*B_q$

$${}^m\mathbb{U}_q^p: {}^*B_q \ni V \mapsto \frac{q}{p}V \in {}^*B_p$$

Properties

- $\langle U, {}^m\mathbb{U}_q^p V \rangle_p = \langle {}^e\mathbb{U}_p^q U, V \rangle_q$
- ${}^e\mathbb{U}_q^r {}^e\mathbb{U}_p^q = {}^e\mathbb{U}_p^r$
- ${}^m\mathbb{U}_q^r {}^m\mathbb{U}_p^q = {}^m\mathbb{U}_p^r$
- $\langle {}^e\mathbb{U}_p^q U, {}^m\mathbb{U}_p^q V \rangle_q = \langle U, V \rangle_p$
- $d^2K_p(q)[U, V] = \langle {}^e\mathbb{U}_p^q U, {}^e\mathbb{U}_p^q V \rangle_q = \langle {}^m\mathbb{U}_q^p {}^e\mathbb{U}_p^q U, V \rangle_p$.

Statistical exponential manifold and bundles

- The *exponential manifold* is the maximal exponential family \mathcal{E} with the affine atlas of global charts $(s_p: p \in \mathcal{E})$,

$$s_p(q) = \log \frac{q}{p} - \mathbb{E}_p \left[\log \frac{q}{p} \right].$$

- The *statistical exponential bundle* $T\mathcal{E}(M)$ is the manifold defined on the set

$$\{(p, V) | p \in \mathcal{E}, V \in B_p\}$$

by the affine atlas of global charts

$$\sigma_p: (q, V) \mapsto (s_p(q), {}^e\mathbb{U}_q^p V) \in B_p \times B_p, \quad p \in \mathcal{E}$$

- The *statistical predual bundle* ${}^*T\mathcal{E}(M)$ is the manifold defined on the set

$$\{(p, W) | p \in \mathcal{E}, W \in {}^*B_p\}$$

by the affine atlas of global charts

$${}^*\sigma_p: (q, W) \mapsto (s_p(q), {}^m\mathbb{U}_q^p W) \in B_p \times {}^*B_p, \quad p \in \mathcal{E}$$

Score and statistical gradient

Definition

$t \mapsto p(t)$ is a curve in $\mathcal{E}(p)$ and $F: \mathcal{E}(p) \rightarrow \mathbb{R}$.

- The *score* of the curve $t \mapsto p(t)$ is a curve in the statistical bundle $t \mapsto (p(t), Dp(t)) \in S\mathcal{E}(p)$ such that for all $X \in L^{(\cosh-1)*}(p)$ it holds

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

- The *statistical gradient* of F is a *section* of the statistical bundle, $p \mapsto (p, \text{grad } f(p)) \in eTofp$ such that for each regular curve $t \mapsto p(t)$, it holds

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

Everything applies if the tangent space is in $C_0^{1,(\cosh-1)}(p)$, but technical details have to be checked, e.d. ${}^m\mathbb{U}_p^q$

Taylor formula in the Statistical Bundle

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

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

$$\text{with } R_2(f, p(\cdot)) = \int_0^1 (1-t) \left(\left. \frac{d^2}{dt^2} F(p(t)) \right|_{t=0} - \left. \frac{d^2}{dt^2} F(p(t)) \right|_{t=0} \right) dt$$

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

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

- where $\left. \frac{d}{dt} \langle \text{grad } F(p(t)), Dp(t) \rangle_{p(t)} \right|_{t=0}$ depends on $p(\cdot)$

Accelerations

- $p(t) = e^{U(t) - K_p(U(t))} \cdot p$, $U \in B_p$.
- Let us define the *acceleration* at t of a curve $t \mapsto p(t) \in \mathcal{E}(M)$.
The velocity is defined to be
 $t \mapsto (p(t), Dp(t)) = (p(t), \frac{d}{dt} \log(p(t))) \in T\mathcal{E}(M)$
- The *exponential acceleration* is $t \mapsto (p(t), {}^e D^2 p(t)) \in T\mathcal{E}(M)$ with

$${}^e D^2 p(t) = \left. \frac{d}{ds} {}^e \mathbb{U}_{p(s)}^{p(t)} Dp(s) \right|_{s=t} = \ddot{U}(t) - \mathbb{E}_{p(t)} [\ddot{U}(t)]$$

- The *mixture acceleration* is

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

Computation

$$\begin{aligned} {}^m D^2 p(t) &= \frac{d}{ds} \left. {}^m \mathbb{U}_{\rho(s)}^{\rho(t)} Dp(s) \right|_{s=t} \\ &= \frac{d}{ds} \left. {}^m \mathbb{U}_{\rho(s)}^{\rho(t)} \frac{d}{ds} (U(s) - K_p(U(s))) \right|_{s=t} \\ &= \frac{d}{ds} \left. {}^m \mathbb{U}_{\rho(s)}^{\rho(t)} (\dot{U}(s) - dK_p(U(s))\dot{U}(s)) \right|_{s=t} \\ &= \frac{d}{ds} \left. \frac{\rho(s)}{\rho(t)} (\dot{U}(s) - \mathbb{E}_{\rho(t)} [\dot{U}(s)]) \right|_{s=t} \\ &= \frac{\dot{\rho}(s)}{\rho(t)} (\dot{U}(s) - \mathbb{E}_{\rho(t)} [\dot{U}(s)]) + \frac{\rho(s)}{\rho(t)} (\ddot{U}(s) - \mathbb{E}_{\rho(t)} [\ddot{U}(s)]) \Big|_{s=t} \\ &= \frac{\dot{\rho}(t)}{\rho(t)} (\dot{U}(t) - \mathbb{E}_{\rho(t)} [\dot{U}(t)]) + \frac{\rho(t)}{\rho(t)} (\ddot{U}(t) - \mathbb{E}_{\rho(t)} [\ddot{U}(t)]) \\ &= (\dot{U}(t) - \mathbb{E}_{\rho(t)} [\dot{U}(t)])^2 + (\ddot{U}(t) - \mathbb{E}_{\rho(t)} [\ddot{U}(t)]) \\ &= \frac{\ddot{\rho}(t)}{\rho(t)} \end{aligned}$$

Autoparallel curves

- The *exponential acceleration* is

$${}^eD^2p(t) = \ddot{U}(t) - \mathbb{E}_{p(t)} [\ddot{U}(t)]$$

- The e-acceleration of a 1d-exponential family $t \mapsto e^{tU - K_M(tU)}$. M is zero, $t\ddot{U} = 0$.
- The *mixture acceleration* is

$${}^mD^2p(t) = \frac{\ddot{p}(t)}{p(t)}$$

- The m-acceleration of a 1d-mixture family $t \mapsto p(0) + t(p(1) - p(0))$ is zero,

$$\frac{p(0) + t(\ddot{p}(1) - p(0))}{p(0) + t(p(1) - p(0))} = 0$$

Taylor's formulæ

The computation of $\left. \frac{d}{dt} \langle \text{grad } F(p(t)), Dp(t) \rangle_{p(t)} \right|_{t=0}$ reduces to one term by choosing an autoparallel segment connecting p and q

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

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

- If $t \mapsto p(t)$ is the *exponential geodesic* connecting $p = p(0)$ to $q = p(1)$,

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