

# Nonparametric Information Geometry http://www.giannidiorestino.it/GSI2013-talk.pdf

Giovanni Pistone

de Castro Statistics Initiative Collegio Carlo Alberto Moncalieri, Italy

August 30, 2013

## Abstract

The differential-geometric structure of the set of positive densities on a given measure space has raised the interest of many mathematicians after the discovery by C.R. Rao of the geometric meaning of the Fisher information. Most of the research is focused on parametric statistical models. In series of papers by author and coworkers a particular version of the nonparametric case has been discussed. It consists of a minimalistic structure modeled according the theory of exponential families: given a reference density other densities are represented by the centered log likelihood which is an element of an Orlicz space. This mappings give a system of charts of a Banch manifold. It has been observed that, while the construction is natural, the practical applicability is limited by the technical difficulty to deal with such a class of Banach spaces. It has been suggested recently to replace the exponential function with other functions with similar behavior but polynomial growth at infinity in order to obtain more tractable Banach spaces, e.g. Hilbert spaces. We give first a review of our theory with special emphasis on the specific issues of the infinite dimensional setting. In a second part we discuss two specific topics, differential equations and the metric connection. The position of this line of research with respect to other approaches is briefly discussed.

#### References in

- GP, GSI2013 Proceedings. A few typos corrected in arXiv:1306.0480;
- GP, arXiv:1308.5312

• If  $\mu_1$ ,  $\mu_2$  are equivalent measures on the same sample space, a statistical model has two representations

$$L_1(x;\theta)\mu_1(dx) = L_2(x;\theta)\mu_2(dx).$$

• Fisher's score is a valid option

$$s(x;\theta) = \frac{d}{d\theta} \ln L_i(x;\theta), \quad i = 1, 2,$$

and  $\mathsf{E}_{\theta}[s_{\theta}] = 0$ .

• Each density q equivalent to p is of the form

$$q(x) = \frac{\mathrm{e}^{v(x)} p(x)}{\mathsf{E}_{p} \left[ \mathrm{e}^{v} \right]} = \exp\left(v(x) - \ln\left(\mathsf{E}_{p} \left[ \mathrm{e}^{v} \right] \right)\right) p(x),$$

where v is a random variable such that  $E_p[e^v] < +\infty$ .

• To avoid borderline cases, we actually require

$$\mathsf{E}_{p}\left[\mathrm{e}^{ heta \mathbf{v}}
ight] < +\infty, \quad heta \in I ext{ open } \supset [0,1].$$

• Finally, we require  $E_p[v] = 0$ .



# Part I Exponential manifold Part II Vector bundles Part III Deformed exponential

# Part I Exponential manifold

# Sets of densities

## Definition

 $\mathcal{P}^1$  is the set of real random variables f such that  $\int f d\mu = 1$ ,  $\mathcal{P}_{\geq}$  the convex set of probability densities,  $\mathcal{P}_{>}$  the convex set of strictly positive probability densities:

$$\mathcal{P}_{>}\subset\mathcal{P}_{\geq}\subset\mathcal{P}^{1}$$

- We define the (differential) geometry of these spaces in a way which is meant to be a non-parametric generalization of Information Geometry
- We try to avoid the use of explicit parameterization of the statistical models and therefore we use a parameter free presentation of differential geometry.
- We construct a manifold modeled on an Orlicz space.
- We look for applications to applications intrisically non parametric, i.e. Statistical Physics, Information Theory, Optimization, Filtering.

## Banach manifold

#### Definition

- 1. Let  $\mathcal{P}$  be a set,  $\mathcal{E} \subset \mathcal{P}$  a subset, B a Banach space. A 1-to-1 mapping  $s: \mathcal{E} \to B$  is a chart if the image  $s(\mathcal{E}) = \mathcal{S} \subset B$  is open.
- 2. Two charts  $s_1: \mathcal{E}_1 \to B_1$ ,  $s_2: \mathcal{E}_2 \to B_2$ , are both defined on  $\mathcal{E}_1 \cap \mathcal{E}_2$ and are compatible if  $s_1(\mathcal{E}_1 \cap \mathcal{E}_2)$  is an open subset of  $B_1$  and the change of chart mapping

$$s_2 \circ s_1^{-1}$$
:  $s_1(\mathcal{E}_1 \cap \mathcal{E}_2) \xrightarrow{s_1^{-1}} \mathcal{E}_1 \cap \mathcal{E}_2 \xrightarrow{s_2} s_2(\mathcal{E}_1 \cap \mathcal{E}_2)$ 

is smooth.

- 3. An atlas is a set of compatible charts.
- Condition 2 implies that the model spaces  $B_1$  and  $B_2$  are isomorphic.
- In our case: P = P<sub>></sub>, the atlas has a chart s<sub>p</sub> for each p ∈ P<sub>></sub> such that s<sub>p</sub>(p) = 0 and two domains E<sub>p1</sub> and E<sub>p2</sub> are either equal or disjoint.

## Charts on $\mathcal{P}_{>}$



# Model space

#### Orlicz Φ-space

If  $\phi(y) = \cosh y - 1$ , the Orlicz  $\Phi$ -space  $L^{\Phi}(p)$  is the vector space of all random variables such that  $\mathsf{E}_p[\Phi(\alpha u)]$  is finite for some  $\alpha > 0$ .

#### Properties of the $\Phi$ -space

- 1.  $u \in L^{\Phi}(p)$  if, and only if, the moment generating function  $\alpha \mapsto \mathsf{E}_{p}[\mathrm{e}^{\alpha u}]$  is finite in a neighborhood of 0.
- 2. The set  $S_{\leq 1} = \{ u \in L^{\Phi}(p) | \mathsf{E}_{p} [\Phi(u)] \leq 1 \}$  is the closed unit ball of a Banach space with norm

$$\|u\|_{p} = \inf \left\{ \rho > 0 \left| \mathsf{E}_{p} \left[ \Phi \left( \frac{u}{\rho} \right) \right] \le 1 \right\}.$$

3.  $\|u\|_{\rho} = 1$  if either  $E_{\rho}[\Phi(u)] = 1$  or  $E_{\rho}[\Phi(u)] < 1$  and  $E_{\rho}\left[\Phi\left(\frac{u}{\rho}\right)\right] = \infty$  for  $\rho < 1$ . If  $\|u\|_{\rho} > 1$  then  $\|u\|_{\rho} \le E_{\rho}[\Phi(u)]$ . In particular,  $\lim_{\|u\|_{\rho}\to\infty} E_{\rho}[\Phi(u)] = \infty$ .

## Example: boolean state space

- In the case of a finite state space, the moment generating function is finite everywhere, but its computation can be challenging.
- Boolean case:  $\Omega = \{+1, -1\}^n$ , uniform density  $p(x) = 2^{-n}, x \in \Omega$ . A generic real function on  $\Omega$  has the form  $u(x) = \sum_{\alpha \in L} \hat{u}(\alpha) x^{\alpha}$ , with  $L = \{0, 1\}^n$ ,  $x^{\alpha} = \prod_{i=1}^n x_i^{\alpha_i}$ ,  $\hat{u}(\alpha) = 2^{-n} \sum_{x \in \Omega} u(x) x^{\alpha}$ .
- The moment generating function of *u* under the uniform density *p* is

$$\mathsf{E}_{p}\left[\mathrm{e}^{tu}\right] = \sum_{B \in \mathcal{B}(\hat{u})} \prod_{\alpha \in B^{c}} \cosh(t\hat{u}(\alpha)) \prod_{\alpha \in B} \sinh(t\hat{u}(\alpha)),$$

where  $\mathcal{B}(\hat{u})$  are those  $B \subset \text{Supp } \hat{u}$  such that  $\sum_{\alpha \in B} \alpha = 0 \mod 2$ .

$$\mathsf{E}_{\mathsf{P}}\left[\Phi(tu)\right] = \sum_{B \in \mathcal{B}_{0}(\hat{u})} \prod_{\alpha \in B^{c}} \cosh(t\hat{u}(\alpha)) \prod_{\alpha \in B} \sinh(t\hat{u}(\alpha)) - 1,$$

where  $\mathcal{B}_0(\hat{u})$  are those  $B \subset \text{Supp } \hat{u}$  such that  $\sum_{\alpha \in B} \alpha = 0 \mod 2$ and  $\sum_{\alpha \in \text{Supp } \hat{u}} \alpha = 0$ . Example : the sphere is not smooth in general

- $p(x) \propto (a+x)^{-\frac{3}{2}} e^{-x}$ , x, a > 0.
- For the random variable u(x) = x, the function

$$\mathsf{E}_{\rho}\left[\Phi(\alpha u)\right] = \frac{1}{\mathrm{e}^{\mathsf{a}}\,\Gamma\left(-\frac{1}{2},a\right)} \int_{0}^{\infty} (a+x)^{-\frac{3}{2}} \frac{\mathrm{e}^{-(1-\alpha)x} + \mathrm{e}^{-(1+\alpha)x}}{2} \, dx - 1$$

is convex lower semi-continuous on  $\alpha \in \mathbb{R}$ , finite for  $\alpha \in [-1, 1]$ , infinite otherwise, hence not smooth.



# Isomorphism of $L^{\Phi}$ spaces

#### Theorem

 $L^{\Phi}(p) = L^{\Phi}(q)$  as Banach spaces if  $\int p^{1-\theta} q^{\theta} d\mu$  is finite on an open neighborhood I of [0,1]. It is an equivalence relation  $p \smile q$  and we denote by  $\mathcal{E}(p)$  the class containing p. The two spaces have equivalent norms

#### Proof.

Assume  $u \in L^{\Phi}(p)$  and consider the convex function  $C: (s, \theta) \mapsto \int e^{su} p^{1-\theta} q^{\theta} d\mu$ . The restriction  $s \mapsto C(s, 0) = \int e^{su} p d\mu$  is finite on an open neighborhood  $J_p$  of 0; the restriction  $\theta \mapsto C(0, \theta) = \int p^{1-\theta} q^{\theta} d\mu$  is finite on the open set  $I \supset [0, 1]$ . hence, there exists an open interval  $J_q \ni 0$  where  $s \mapsto C(s, 1) = \int e^{su} q d\mu$  is finite.



## e-charts

## Definition (e-chart)

For each  $p\in\mathcal{P}_>$ , consider the chart  $s_p\colon\mathcal{E}(p) o L^{\Phi}_0(p)$  by

$$q\mapsto s_p(q)=\log\left(rac{q}{p}
ight)+D(p\|q)=\log\left(rac{q}{p}
ight)-\mathsf{E}_p\left[\log\left(rac{q}{p}
ight)
ight]$$

For  $u \in L_0^{\Phi}(p)$  let  $\mathcal{K}_p(u) = \ln \mathsf{E}_p[\mathrm{e}^u]$  the cumulant generating function of u and let  $\mathcal{S}_p$  the interior of the proper domain. Define

$$e_p \colon \mathcal{S}_p \ni u \mapsto e^{u - K_p(u)} \cdot p$$

 $e_p \circ s_p$  is the identity on  $\mathcal{E}(p)$  and  $s_p \circ e_p$  is the identity on  $\mathcal{S}_p$ .

Theorem (Exponential manifold)  $\{s_p: \mathcal{E}(p)|p \in \mathcal{P}_{>}\}$  is an affine atlas on  $\mathcal{P}_{>}$ .

## Cumulant functional

- The divergence  $q \mapsto D(p||q)$  is represented in the chart centered at p by  $K_p(u) = \log \mathsf{E}_p[\mathrm{e}^u]$ , where  $q = \mathrm{e}^{u K_p(u)} \cdot p$ ,  $u \in B_p = L_0^{\Phi}(p)$ .
- $K_p: B_p \to \mathbb{R}_{\geq} \cup \{+\infty\}$  is convex and its proper domain  $\text{Dom}(K_p)$  contains the open unit ball of  $T_p$ .
- $K_p$  is infinitely Gâteaux-differentiable on the interior  $S_p$  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:

$$d K_{p} uv = E_{q} [v]$$
  
$$d^{2} K_{p} u(v_{1}, v_{2}) = Cov_{q} (v_{1}, v_{2})$$
  
$$d^{3} K_{p} u(v_{1}, v_{2}, v_{3}) = Cov_{q} (v_{1}, v_{2}, v_{3})$$

# Change of coordinate

The following statements are equivalent:

1.  $q \in \mathcal{E}(p);$ 2.  $p \smile q;$ 3.  $\mathcal{E}(p) = \mathcal{E}(q);$ 4.  $\ln\left(\frac{q}{p}\right) \in L^{\Phi}(p) \cap L^{\Phi}(q).$ 

1. If  $p, q \in \mathcal{E}(p) = \mathcal{E}(q)$ , the change of coordinate

$$s_q \circ e_p(u) = u - \mathsf{E}_q\left[u
ight] + \ln rac{p}{q} - \mathsf{E}_q\left[\ln rac{p}{q}
ight]$$

is the restriction of an affine continuous mapping.

2.  $u \mapsto u - \mathsf{E}_q[u]$  is an affine transport from  $B_p = L_0^{\Phi}(p)$  unto  $B_q = L_0^{\Phi}(q)$ .

## Summary

$$\begin{array}{c} p \smile q \end{array} \Longrightarrow \begin{array}{c} \mathcal{E}(p) \xrightarrow{s_p} \mathcal{S}_p \xrightarrow{l} \mathcal{B}_p \xrightarrow{l} \mathcal{L}^{\Phi}(p) \\ & \left\| \begin{array}{c} s_q \circ s_p^{-1} \\ s_q \circ s_p^{-1} \end{array} \right\| \\ \mathcal{E}(q) \xrightarrow{s_q} \mathcal{S}_q \xrightarrow{l} \mathcal{B}_q \xrightarrow{l} \mathcal{L}^{\Phi}(q) \end{array}$$

• If 
$$p \smile q$$
, then  $\mathcal{E}(p) = \mathcal{E}(q)$  and  $L^{\Phi}(p) = L^{\Phi}(q)$ .

• 
$$B_p = L_0^{\Phi}(p), \ B_q = L_0^{\Phi}(q)$$

•  $\mathcal{S}_p \neq \mathcal{S}_q$  and  $s_q \circ s_p^{-1} \colon \mathcal{S}_p \to \mathcal{S}_q$  is affine

$$s_q \circ s_p^{-1}(u) = u - \mathsf{E}_q[u] + \ln\left(\frac{p}{q}\right) - \mathsf{E}_q\left[\ln\left(\frac{p}{q}\right)\right]$$

The tangent application is d(s<sub>q</sub> ∘ s<sub>p</sub><sup>-1</sup>)(v) = v − E<sub>q</sub>[v] (does not depend on p)

# Duality

## Young pair (N-function)

- $\bullet \ \phi^{-1} = \phi_* \text{,}$
- $\Phi(x) = \int_0^{|x|} \phi(u) \ du$

• 
$$\Phi_*(y) = \int_0^{|y|} \phi_*(v) \, dv$$

•  $|xy| \leq \Phi(x) + \Phi_*(y)$ 



• 
$$L^{\Phi_*}(p) \times L^{\Phi}(p) \ni (v, u) \mapsto \langle u, v \rangle_p = \mathsf{E}_p[uv]$$

• 
$$\left| \langle u, v \rangle_{p} \right| \leq 2 \left\| u \right\|_{\Phi_{*}, p} \left\| v \right\|_{\Phi, p}$$

• 
$$(L^{\Phi_*}(p))' = L^{\Phi}(p)$$
 because  $\Phi_*(ax) \le a^2 \Phi_*(x)$  if  $a > 1$   $(\Delta_2)$ .

#### m-charts

For each  $p \in \mathcal{P}_{>}$ , consider a second type of chart on  $f \in \mathcal{P}^{1}$ :

$$\eta_p: f o \eta_p(f) = rac{f}{p} - 1$$

#### Definition (Mixture manifold)

The chart is defined for all  $f \in \mathcal{P}^1$  such that f/p - 1 belongs to  ${}^*B_p = L_0^{\Phi_+}(p)$ . The atlas  $(\eta_p : {}^*\mathcal{E}(p))$ ,  $p \in \mathcal{P}_>$  defines a manifold on  $\mathcal{P}^1$ .

If the sample space is not finite, such a map does not define charts on  $\mathcal{P}_{>},$  nor on  $\mathcal{P}_{\geq}.$ 

Example:  $N(\mu, \Sigma)$ , det  $\Sigma \neq 0 | I$ 

$$\mathcal{G} = \left\{ (2\pi)^{-\frac{n}{2}} (\det \Sigma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right) \middle| \mu \in \mathbb{R}^n, \Sigma \in \mathsf{Sym}^n_+ \right\}.$$

$$\ln\left(\frac{f(x)}{f_{0}(x)}\right) = -\frac{1}{2}\ln\left(\det\Sigma\right) - \frac{1}{2}(x-\mu)^{T}\Sigma^{-1}(x-\mu) + \frac{1}{2}x^{T}x$$

$$= \frac{1}{2}x^{T}(I-\Sigma^{-1})x + \mu^{T}\Sigma^{-1}x - \frac{1}{2}\mu^{T}\Sigma^{-1}\mu - \frac{1}{2}\ln\left(\det\Sigma\right)$$

$$E_{f_{0}}\left[\ln\left(\frac{f}{f_{0}}\right)\right] = \frac{1}{2}(n-\operatorname{Tr}\Sigma^{-1}) - \frac{1}{2}\mu^{T}\Sigma^{-1}\mu - \frac{1}{2}\ln\left(\det\Sigma\right)$$

$$u(x) = \ln\left(\frac{f(x)}{f_{0}(x)}\right) - E_{f_{0}}\left[\ln\left(\frac{f}{f_{0}}\right)\right]$$

$$= \frac{1}{2}x^{T}(I-\Sigma^{-1})x + \mu^{T}\Sigma^{-1}x - \frac{1}{2}(n-\operatorname{Tr}\Sigma^{-1})$$

$$K_{f_{0}}(u) = -\frac{1}{2}(n-\operatorname{Tr}\Sigma^{-1}) + \frac{1}{2}\mu^{T}\Sigma^{-1}\mu + \frac{1}{2}\ln\left(\det\Sigma\right)$$

Example:  $N(\mu, \Sigma)$ , det  $\Sigma \neq 0 | I |$ 

#### ${\mathcal G}$ as a sub-manifold of ${\mathcal P}_>$

$$\mathcal{G} = \left\{ x \mapsto \mathrm{e}^{u(x) - \mathcal{K}(u)} f_0(x) \Big| u \in \mathcal{H}_{1,2} \cap \mathcal{S}_{f_0} 
ight\}$$

•  $\mathcal{H}_{1,2}$  is the Hemite space of total degree 1 and 2, that is the vector space generated by the Hermite polynomials

$$X_1, \ldots, X_n, (X_1^2 - 1), \ldots, (X_n^2 - 1), X_1 X_2, \ldots, X_{n-1} X_n$$

• If the matrix S,  $S_{ii} = \beta_{ii} - \frac{1}{2}$ ,  $S_{ij} = S_{ji} = \frac{\beta_{ij}}{2}$  is negative definite then  $u \in S_{f_0}$ ,

$$u(x) = \beta_1 x_1 + \dots + \beta_n x_n + \beta_{11} (x_1^2 - 1) + \dots + \beta_{nn} (x_n^2 - 1) + \beta_{12} x_1 x_2 + \dots + \beta_{(n-1)n} x_{n-1} x_n$$

# Part 2: Vector Bundles

## Velocity

Given a one dimensional statistical model p<sub>θ</sub> ∈ E (p), θ ∈ I, I open interval, 0 ∈ I, the representation in the chart centered at p of the e-manifold is u<sub>θ</sub> = s<sub>P</sub>(p<sub>θ</sub>),

$$p_{\theta} = e_p(u_{\theta}) = e^{u_{\theta} - K_p(u_{\theta})} \cdot p.$$

- As  $\mathcal{E}(p) \subset {}^{*}\mathcal{E}(p)$ , there is a representation in the m-manifold  $\eta_{\theta} = \eta_{p}(p_{\theta})$  $\eta_{\theta} = \frac{p_{\theta}}{p} - 1 = e^{u_{\theta} - K_{p}(u_{\theta})} - 1$
- To compute the velocity along a one-parameter statistical model in the s<sub>p</sub> chart we use u<sub>θ</sub>.
- To compute the velocity along a one-parameter statistical model in the  $\eta_p$  chart we use  $\dot{\eta}_{\theta}$ .

## Relation between the two presentation

• We get in the first case

$$\dot{p}_{ heta} = p_{ heta}(\dot{u}_{ heta} - \mathsf{E}_{ heta}[\dot{u}_{ heta}])$$

so that

$$\dot{u}_{\theta} = \frac{\dot{p}_{\theta}}{p_{\theta}} - \mathsf{E}_{p} \left[ \frac{\dot{p}_{\theta}}{p_{\theta}} \right] = {}^{\mathsf{e}} \mathbb{U}_{p(t)}^{p} \frac{\dot{p}_{\theta}}{p_{\theta}}$$

• In the second case we get

$$\dot{\eta}_{p}(\theta) = \frac{\dot{p}_{\theta}}{p} = \frac{p_{\theta}}{p} \frac{\dot{p}_{\theta}}{p_{\theta}} = {}^{\mathsf{m}} \mathbb{U}_{p}^{p_{\theta}} \frac{\dot{p}_{\theta}}{p_{\theta}}$$

# Moving frame

 Both in the e-manifold and the m-manifold there is one chart centered at each density. The two representations u
<sub>θ</sub> and η
<sub>p</sub>(θ) are equal at θ = 0 and are transported back to θ:

$$\frac{\dot{p}_{\theta}}{p_{\theta}} = \dot{u}_{\theta} - \mathsf{E}_{\theta} \left[ \dot{u}_{\theta} \right] = \frac{p}{p_{\theta}} \dot{\eta}_{\theta}.$$

#### Fisher information

The random variable  $\dot{p}_{\theta}/p_{\theta}$  is the Fisher score at  $\theta$  of the one-parameter model  $p_{\theta}$ . The Fisher information at  $\theta$  is the  $L^2$ -norm of the velocity vector of the statistical model in the moving frame centered at  $\theta$ . Moreover,

$$\mathsf{E}_{\theta}\left[\left(\frac{\dot{p}_{\theta}}{p_{\theta}}\right)^{2}\right] = \mathsf{E}_{\theta}\left[\left(\dot{u}_{\theta} - \mathsf{E}_{\theta}\left[\dot{u}_{\theta}\right]\right)\left(\dot{\eta}_{\theta}\frac{p}{p_{\theta}}\right)\right] = \mathsf{E}_{p}\left[\dot{u}_{\theta}\dot{\eta}_{\theta}\right].$$

## Transport

- At each point p ∈ P> of the statistical manifold is attached the vector space B<sub>p</sub>. Similarly, \*B<sub>p</sub>.
- The derivatives of the change of reference maps provide two isomorphic mappings

e-transport 
$${}^{e}\mathbb{U}_{\rho}^{q} \colon B_{\rho} \ni u \mapsto u - \mathbb{E}_{q}[u] \in B_{q}$$
  
m-transport  ${}^{m}\mathbb{U}_{\rho}^{q} \colon {}^{*}B_{\rho} \ni v \mapsto \frac{p}{q}v \in {}^{*}B_{q}$ 

• The two transport are adjoint to each other. If  $u \in B_p$ ,  $v \in {}^*B_q$ .

$$\mathsf{E}_{q}\left[\left({}^{\mathsf{e}}\mathbb{U}_{p}^{q}u\right)v\right] = \mathsf{E}_{q}\left[\left(u - \mathsf{E}_{q}\left[u\right]\right)v\right] = \mathsf{E}_{p}\left[u\left(\frac{q}{p}v\right)\right] = \mathsf{E}_{p}\left[u\left({}^{\mathsf{m}}\mathbb{U}_{q}^{p}v\right)\right]$$

• If  $u \in \mathcal{S}_p$ ,  $q = e_p(u) \in \mathcal{E}(p)$ ,  $v, w \in B_p$ ,

$$d^{2} \mathcal{K}_{p}(u)(v, w) = \operatorname{Cov}_{q}(v, w) = \\ \mathsf{E}_{q}\left[\left({}^{\mathsf{e}} \mathbb{U}_{p}^{q} u\right)\left({}^{\mathsf{e}} \mathbb{U}_{p}^{q} v\right)\right] = \mathsf{E}_{p}\left[\left({}^{\mathsf{m}} \mathbb{U}_{p}^{p} {}^{\mathsf{e}} \mathbb{U}_{p}^{q} v\right) w\right]$$

# Tangent bundle

### Definition

1. For each maximal exponential model  ${\mathcal E}$  the tangent bundle T  ${\mathcal E}$  is the set

$$\mathsf{T}\,\mathcal{E} = \{(p, u) | p \in \mathcal{E}, u \in B_p\}$$

endowed with the atlas of charts

$$s_{\rho} \colon \mathsf{T} \mathcal{E} \ni (q, v) \mapsto (s_{\rho}(q), {^{\mathrm{e}}\mathbb{U}}_{q}^{\rho}v) \in \mathcal{S}_{\rho} \times B_{\rho} \subset B_{\rho} \times B_{\rho}$$

2. 
$$T\mathcal{P}_{>} = \cup_{\mathcal{E}} T\mathcal{E}$$

#### Velocity field

- 1. If p(t),  $t \in I$  is a smooth curve in  $\mathcal{E}$ , then  $(p(t), \delta p(t))$ ,  $t \in I$  is a curve in  $T \mathcal{E}$ ,  $\delta p(t) = \dot{p}(t)/p(t)$ .
- 2. If  $E: \mathcal{E} \to \mathbb{R}$  and  $E_p = E \circ e_p \colon \mathcal{S}_p \to \mathbb{R}$  is  $C^1$ ,  $p \in \mathcal{E}$ , then

$$\frac{d}{dt}E(p(t)) = \frac{d}{dt}E_p(u(t)) = dE_p(u(t))^{e}\mathbb{U}_{p(t)}^p\delta p(t)$$

## e-covariant-derivative

#### Definition (Exponential covariant derivative)

- Let E: E → R be smooth. For each vector field G of T E the covariant derivative D<sub>G</sub>E: E → R is defined by D<sub>G</sub>E(p) = dE<sub>p</sub>(0)G(p).
- 2. If the linear functional  $dE_p \in \mathcal{L}(B_p)$  is representable in  ${}^*B_p$ , the representative is the (natural) gradient  $\nabla E$  and  $D_G E(p) = \mathbb{E}_p [(\nabla E(p)) G(p)]$
- 3. Let *F* be a vector field in T  $\mathcal{E}$ . In the chart centered at *p* we have  $F_p(u) = {}^{e}\mathbb{U}_{e_p(u)}^p F \circ e_p(u)$ . It is a mapping  $F_p : S_{\rightarrow}B_p$  we assume differentiable. The covariant derivative  $D_GF$  is the vector field of T  $\mathcal{E}$  defined by  $D_GF(p) = dF_p(0)G(p)$ .

## Example: KL divergence

- For each  $q \in \mathcal{E}(p_0)$ ,  $E \colon \mathcal{E}(p_0) \ni q \mapsto \mathsf{D}(q \| p_0) = \mathsf{E}_q \left[ \ln \left( \frac{q}{p_0} \right) \right]$ .
- In the chart centered at any  $p \in \mathcal{E}(p_0)$ ,  $q = e^{u K_p(u)} \cdot p$ ,  $p_0 = e^{u_0 - K_p(u_0)} \cdot p$ ,

$$E_{\rho}(u) = \mathsf{E}_{q} \left[ u - K_{\rho}(u) - u_{0} + K_{\rho}(u_{0}) \right] = dK_{\rho}(u)(u - u_{0}) - K_{\rho}(u) + K_{\rho}(u_{0})$$

• 
$$dE_p(0)v = d^2K_p(0)(-u_0,v) + \underline{dK_p(\theta)v} - \underline{K_p(\theta)v} = -E_p(u_0v)$$

$$D_{G}E(p) = \mathsf{E}_{p}\left[\left(\mathsf{ln}\left(\frac{p}{p_{0}}\right) - \mathsf{E}_{p}\left[\mathsf{ln}\left(\frac{p}{p_{0}}\right)\right]\right)G(p)\right]$$
$$= \mathsf{E}_{p}\left[\left(\mathsf{ln}\left(\frac{p}{p_{0}}\right) - E(p)\right)G(p)\right]$$

• 
$$\nabla E(p) = \ln\left(\frac{p}{p_0}\right) - E(p) \in B_p$$

## Pre-tangent bundle

### Definition

1. For each maximal exponential model  ${\mathcal E}$  the pretangent bundle  ${}^*{\sf T}\,{\mathcal E}$  is the set

$$^{*}\mathsf{T}\mathcal{E} = \{(p, v) | p \in \mathcal{E}, v \in ^{*}B_{p}\}$$

endowed with the atlas of charts

$$s_{\rho} \colon \ {}^{*}\mathsf{T} \ \mathcal{E} \ni (q, \nu) \mapsto (s_{\rho}(q), {}^{\mathsf{m}}\mathbb{U}_{q}^{\rho}\nu) \in \mathcal{S}_{\rho} \times {}^{*}B_{\rho} \subset B_{\rho} \times {}^{*}B_{\rho}$$

 $2. \ ^*T \mathcal{P}_{>} = \cup_{\mathcal{E}} \ ^*T \mathcal{E}$ 

## m-covariant-derivative

#### Definition (Mixture covariant derivative)

Let *F* be a vector field in  ${}^*T\mathcal{E}$ . In the chart centered at *p* we have  $F_p(u) = {}^m \mathbb{U}_{e_p(u)}^p F \circ e_p(u)$ . It is a mapping  $F_p : S_{\rightarrow} {}^*B_p$  that we assume to be differentiable. Let *G* be a vector field in  $T\mathcal{E}$ . The covariant derivative  $D_GF$  is the vector field of  ${}^*T\mathcal{E}$  defined by  $D_GF(p) = dF_p(0)G(p)$ .

#### Theorem

If G, H are smooth vector fields in T  $\mathcal{E}$  and F is a smooth vector field in \*T  $\mathcal{E}$ , then  $\langle G, F \rangle(p) = E_p[G(p)F(p)]$  is a smooth real function on  $\mathcal{E}$ and

$$D_H \langle G, F \rangle = \langle D_H G, F \rangle + \langle G, D_H F \rangle$$



- For each  $p \in \mathcal{P}_{>}$ ,  $\int (\sqrt{p})^2 d\mu = 1$ .
- The embedding  $I: p \mapsto \sqrt{p}$  is represented in the chart centered at p by

$$I_p \colon u \mapsto \mathrm{e}^{u - K_p(u)} \ p \mapsto \mathrm{e}^{\frac{1}{2}u - \frac{1}{2}K_p(u)} \ \sqrt{p}$$

• The mapping  $u \mapsto e^{\frac{1}{2}u - \frac{1}{2}K_p(u)} \in L^2(p)$  is analytic around 0 with differential

$$d(u\mapsto \mathrm{e}^{\frac{1}{2}u-\frac{1}{2}K_p(u)})v\Big|_{u=0}=\frac{1}{2}v\in B_p$$

- As  $\int \left(\frac{1}{2}v\sqrt{p}\right)^2 = \frac{1}{4}E_p\left[v^2\right]$ , the embedding *I* is smooth, with differential  $dI(p)v = \frac{1}{2}v\sqrt{p}$
- $\frac{1}{2}v\sqrt{p}$  belongs to the tangent at  $\sqrt{p}$  of the sphere.

## Parallel transport on the sphere



- $S_{\mu} = \{f \in L^2(\mu) | ||f||_2 = 1\}$  is the unit sphere of  $L^2(\mu)$ . It is an Hilbert manifold, sub-manifold of  $L^2(\mu)$ .
- $H_f = \{h \in L^2(\mu) | \int fh \ d\mu = 0\}$  is the tangent space of  $S_\mu$  at  $f \in S_\mu$ .
- Given  $f, g \in S_{\mu}$ , consider the circle  $S_{\mu} \cap \text{Span}(f, g)$ .

•  $U_f^g: H_f \to H_g$  is an isometric isomorphism if  $f, g \in S$  and

$$U_f^g h = h - \left(1 + \int fg \ d\mu\right)^{-1} (f+g) \int gh \ d\mu$$

• If  $h \perp f, g$ , then  $U_f^g h = h$ , otherwise it is rotated.

## Hilbert transport



Proposition

- 1. The mapping  $\mathbb{U}_p^q$  of is an isometry of  $H_p$  onto  $H_q$ .
- 2. In particular, if  $u_1, \ldots, u_n$  is an orthonormal system in  $H_p$ , then  $\mathbb{U}_p^q u_i$  in an orthonormal system in  $H_q$ .
- 3.  $\mathbb{U}_q^p \mathbb{U}_p^q u = u$  if  $u \in H_p$ . Generally,  $\mathbb{U}_q^r \mathbb{U}_p^q \neq \mathbb{U}_p^r$ .

# Hilbert bundle

### Definition

1. For each maximal exponential model  ${\mathcal E}$  the Hilbert bundle  $H{\mathcal E}$  is the set

$$H\mathcal{E} = \{(p, v) | p \in \mathcal{E}, v \in H_p\}$$

endowed with the atlas of charts

$$s_{p} \colon H \mathcal{E} \ni (q, v) \mapsto (s_{p}(q), \mathbb{U}_{q}^{p}v) \in \mathcal{S}_{p} imes H_{p} \subset B_{p} imes H_{p}$$

2.  $H\mathcal{P}_{>} = \cup_{\mathcal{E}} H\mathcal{E}$ 

Hilbert bundle vs tangent and pretangent bundle

$$\mathbb{U}_{q}^{p}v = \sqrt{\frac{q}{p}}v - \left(1 + \mathsf{E}_{p}\left[\sqrt{\frac{q}{p}}\right]\right)^{-1} \left(1 + \sqrt{\frac{q}{p}}\right)\mathsf{E}_{p}\left[\sqrt{\frac{q}{p}}v\right]$$

• Assume 
$$q, p \in \mathcal{E}$$
,  $q = e^{u - K_p(u)} \cdot p$ . Define  $p(t) = e^{tu - K_p(tu)} \cdot p$ ,  
 $p(1) = q$ . Assume  $v \in B_q \subset H_q$ .

• 
$$\sqrt{\frac{q}{p}} = e^{\frac{1}{2}u - \frac{1}{2}K_p(u)} = e^{-\frac{1}{2}K_p(u) + K_p(u/2)} \frac{p(1/2)}{p}.$$

• 
$$\mathsf{E}_{p}\left[\sqrt{\frac{q}{p}}\right] = \mathrm{e}^{-\frac{1}{2}K_{p}(u)+K_{p}(u/2)}.$$

• 
$$\sqrt{\frac{q}{p}}\mathbf{v} = e^{-\frac{1}{2}K_{p}(u)+K_{p}(u/2)}\frac{p(1/2)}{p}\mathbf{v} = e^{-\frac{1}{2}K_{p}(u)+K_{p}(u/2)} \left({}^{\mathsf{m}}\mathbb{U}_{p(1/2)}^{p}{}^{e}\mathbb{U}_{q}^{p(1/2)}\mathbf{v} + \frac{p(1/2)}{p}\mathsf{E}_{p(1/2)}[\mathbf{v}]\right).$$

• 
$$\mathsf{E}_{p}\left[\sqrt{\frac{q}{p}}\mathbf{v}\right] = \mathrm{e}^{-\frac{1}{2}\mathcal{K}_{p}(u) + \mathcal{K}_{p}(u/2)}\mathsf{E}_{p(1/2)}[v].$$

• Now we compute 
$$\left< \mathbb{U}_q^p v, w \right>_p$$
,  $w \in B_p$ 



#### Definition

Let G, F be vector fields in  $H\mathcal{P}_{>}$ , i.e.  $F(p), G(p) \in H_p\mathcal{P}_{>}$ . We define  $D_GF$  to be the vector field defined by  $D_GF(p) = \frac{d}{dt} \mathbb{U}_{p(t)}^p F(t)\Big|_{t=0}$ , where p(t) is a curve such that p(0) = p and  $\delta p(0) = G(p)$ .

#### Theorem

- 1.  $D_G F$  is a covariant derivative.
- 2. Let  $F_1$ ,  $F_2$ , G be vector fields in  $HP_>$  such that the ordinary products  $GF_1$  and  $GF_2$  are vector fields in  $HP_>$  and

 $D_{G}\mathsf{E}_{p}[F_{1}(p)F_{2}(p)] = \mathsf{E}_{p}[D_{G}F_{1}(p)F_{2}(p)] + \mathsf{E}_{p}[F_{1}(p)D_{G}F_{2}(p)],$ 

# Vector field

#### Definition

A vector field F of the the m-bundle  ${}^*T(p)$ ,  $p \in \mathcal{P}_>$ , is a mapping which is defined on some domain  $D \subset \mathcal{P}_>$  and it is a section of the m-bundle, that is  $F(p) \in {}^*T(p)$ , for all  $p \in D \subset \mathcal{P}_>$ .

### Example

1. For a given 
$$u \in T_p$$
 and all  $q \in \mathcal{E}(p)$ 

$$F: q \mapsto u - \mathsf{E}_q[u]$$

2. For all strictly positive density  $q \in \mathcal{P}_{>}(\mathbb{R}) \cap C^{1}(\mathbb{R})$ 

$$F: q \mapsto rac{q'}{q}$$

3. For all strictly positive  $q\in\mathcal{P}_>(\mathbb{R})\cap C^2(\mathbb{R})$ 

$$F: q \mapsto q''/q$$

## Evolution equation

#### Definition

A one-parameter statistical model in  $\mathcal{P}_{>}$ ,  $p(\theta)$ ,  $\theta \in I$ , solves the evolution equation associated to the vector field F if  $p(\theta) = e^{u(\theta) - K_p(u(\theta))} \cdot p$ :

- 1. the curve  $\theta \mapsto u(\theta) \in T(p)$  is continuous in  $L^2$ ;
- 2. the curve  $\theta \mapsto p(\theta)/p 1 \in {}^*T(p)$  is continuously differentiable;
- 3. for all  $\theta \in I$  it holds

$$rac{\dot{p}( heta)}{p( heta)} = F(p( heta))$$

## Evolution equation in the moving frame

## Theorem (???)

Assume F is locally maximal monotone. Then the equation  $\dot{p}/p + F(p) = 0$  has a solution which is unique.

- The evolution equation above is written with respect to the moving frame at  $p_{\theta}$  because  $\dot{p}(\theta)/p(\theta)$  is the representation of the velocity vector in  ${}^{*}T(p(\theta))$ .
- However, with respect to a fixed frame at p, we should have written

 $\begin{cases} \dot{u}_{\theta} = F(p(\theta)) - \mathsf{E}_{p}\left[F(p(\theta))\right] & \text{e-connection, assuming } \dot{u}_{\theta} \in T_{p_{\theta}} \\ \dot{l}_{\theta} = \frac{p}{p(\theta)}F(p(\theta)) & \text{m-connection} \end{cases}$ 

# Exponential family

### Example

- 1. The exponential model  $p_{\theta} = e^{\theta F} / \Lambda(\theta)$  is a solution of the equation  $\frac{\dot{p}_{\theta}}{p_{\theta}} = F \mathsf{E}_{p_{\theta}}[F].$
- 2. The second example follows by considering  $\Omega = \mathbb{R}$  and taking for domain the class of  $C^2$  positive densities q such that  $F(q) = -q'/q \in {}^*T(f)$ . We can therefore consider the evolution equation  $\dot{p}_{\theta}/p_{\theta} = -F(p_{\theta})$ .

Given any f in the domain, the statistical model  $p_{\theta}(x) = f(x - \theta)$  is such that the score is

$$\frac{\dot{p}_{\theta}(x)}{p_{\theta}(x)} = -\frac{f'(x-\theta)}{p(f-\theta)} = F(p(\cdot-\theta))(x)$$

and therefore is a solution of the evolution equation. The classical Pearson classes of distributions are related to this equation.

## Heat equation

The heat equation  $\frac{\partial}{\partial t}p(t,x) - \frac{\partial^2}{\partial x^2}p(t,x) = 0$  is an interesting example of evolution equation in  $\mathcal{M}_{>}$ . In fact, we can consider the vector field

$$F(p)(x) = \frac{\frac{\partial^2}{\partial x^2}p(x)}{p(x)}$$

Upon division of both sides of of the heat equation by p(t, x), we obtain an equation of our type, whose solution is the solution of the heat equation. Moreover, the heat equation has a variational form. For each  $v \in D$ 

$$\mathsf{E}_{p}[F(p)v] = \int p''(x)v(x) \ dx = -\int p'(x)v'(x) \ dx = -\mathsf{E}_{p}\left[\frac{p'}{p}v'\right]$$

from which the weak form of the evolution equation follows. as

$$\mathsf{E}_{p_{\theta}}\left[\frac{\dot{p}_{\theta}}{p_{\theta}}v\right] + \mathsf{E}_{p_{\theta}}\left[F_{0}(p_{\theta})v\right] = 0 \quad v \in D$$

where  $F_0$  is the vector field associated to the translation model.

## Example: Decision geometry (Dawid&Lauritzen) I

If the sample space is  $\mathbb R$  and  $p,q\in\mathcal P_>$ , write  $q=\mathrm e^{u-K_p(u)}\cdot p$ , so that

$$\log q - \log p = u - K_p(u).$$

Assume *u* belongs to the Sobolev space

$$W^{\Phi,1} = \left\{ u \in L_0^{\Phi}(p) \middle| \nabla u \in L_0^{\Phi}(p) 
ight\}.$$

It follows

$$d(p,q) = \frac{1}{4} \mathsf{E}_{p} \left[ \left\| \nabla \log q - \nabla \log p \right\|^{2} \right]$$
$$= \frac{1}{4} \mathsf{E}_{p} \left[ \left\| \nabla u \right\|^{2} \right].$$

# Example: Decision geometry (Dawid&Lauritzen) II

For  $u, v \in W^{\Phi,1}$  we have a bilinear form

$$\begin{aligned} \langle u, v \rangle_{\rho} &= \mathsf{E}_{\rho} \left[ \nabla u \nabla v \right] = \int u_{x}(x) v_{x}(x) \rho(x) dx \\ &= -\int \nabla (u_{x}(x) \rho(x)) v(x) dx \\ &= -\int (\Delta u(x) \rho(x) + \nabla u(x) \nabla \rho(x)) v(x) dx \\ &= \mathsf{E}_{\rho} \left[ (-\Delta u - \nabla \log \rho \nabla u) v \right] \end{aligned}$$

We have

$$\mathsf{E}_{\rho}\left[\nabla u\nabla v\right] = \mathsf{E}_{\rho}\left[F_{\rho}uv\right], \qquad F_{\rho}u \in {}^{*}W^{\Phi,1}$$

i.e a classical setting for evolution equations  $\partial_t u_t = F_p(u_t)$ .

Example: binary case (L. Malagò) Consider the optimization of  $\theta \mapsto \mathsf{E}_{\theta}[F]$  along a binary exponential family

$$p_{ heta} = \mathsf{E}\left(\sum_{j=1}^{d} heta_j T_j - \mathcal{K}( heta)
ight) \cdot p, \quad T_j^2 = 1.$$

- $\partial_j \mathsf{E}_{\theta} [F] = \mathsf{Cov}_{\theta} (F, T_j) = \mathsf{E}_{\theta} [FT_j] \mathsf{E}_{\theta} [F] \mathsf{E}_{\theta} [T_j]$
- θ is is a critical point if Cov<sub>θ</sub> (F, T<sub>j</sub>) = 0, j = 1,..., d. This is not possible if F is a linear combination of the T<sub>j</sub>'s or if the remainder after projection on the tangent space is small enough.
- At the critical point the Hessian matrix is

$$\partial_i \partial_j = \operatorname{Cov}_{\theta} (F, T_i T_j)$$

which is not zero if F is a linear combination of the  $T_j$ 's and the interactions  $T_i T_j$ 's.

• The diagonal elements of the Hessian matrix are  $\partial_i^2 = \operatorname{Cov}_{\theta}(F, T_i^2) = \operatorname{Cov}_{\theta}(F, 1) = 0$ . The Hessian matrix is not sign-defined at the critical point.

# Part III Deformed exponentials

## Kaniadakis' exponential

Definition  $\left(\exp_{\kappa} \text{ with } \kappa = \frac{1}{2}\right)$  $\exp_{\kappa}(y) = \left(\frac{1}{2}y + \sqrt{1 + \frac{1}{4}y^2}\right)^2$   $\ln_{\kappa}(x) = x^{1/2} - x^{-1/2}$ 



#### Deformed logarithm formalism

$$\begin{aligned} \ln_{\kappa}(x) &= \int_{1}^{x} \frac{du}{\chi(u)} \\ \chi(u) &= 2 \frac{u^{3/2}}{u+1} \\ \frac{d}{dy} \exp_{\kappa}(y) &= \chi(\exp_{\kappa}(y)) \end{aligned}$$



# Musielak-Orlicz space

## Definition $(L^{\kappa}(p))$

For each  $p \in \mathcal{P}_{>}$ , the set

$$L^{\kappa}(\boldsymbol{p}) = \left\{ \boldsymbol{v} \middle| \int \exp_{\kappa} \left( \alpha \left| \boldsymbol{v} \right| + \ln_{\kappa} \left( \boldsymbol{p} \right) \right) \ d\mu < +\infty, \text{some } \alpha > 0 \right\}$$

is a Banach space with closed unit ball

$$S_\kappa = \left\{ v \in L^\kappa(p) igg| \int \exp_\kappa \left( |v| + \ln_\kappa\left(p
ight) 
ight) \ d\mu - 1 \leq 1 
ight\}$$

#### Integrability

From the convexity

$$\exp_{\kappa}\left(lpha\left| m{v} 
ight| + \ln_{\kappa}\left( m{p} 
ight) 
ight) - \exp_{\kappa}\left( \ln_{\kappa}\left( m{p} 
ight) 
ight) \geq \exp_{\kappa}'(\ln_{\kappa}\left( m{p} 
ight))\left| m{v} 
ight|,$$

, i.e.

$$\exp_{\kappa}\left( \alpha \left| \mathbf{v} \right| + \ln_{\kappa}\left( \mathbf{p} 
ight) \right) - \mathbf{p} \geq \chi(\mathbf{p}) \left| \mathbf{v} \right|,$$

then  $L^{\kappa}(p) \subset L^{1}(\chi(p) \cdot d\mu)$ .

## Deformed exponential arc

#### Example

Assume p = 1, so that  $\ln_{\kappa}(p) = 0$ . As for  $y \ge 0$ 

$$y^2 \le \left(\frac{1}{2}y + \sqrt{1 + \frac{1}{4}y^2}\right)^2 \le (1 + y)^2$$

we have

$$\mathcal{L}^{\kappa}(1) = \left\{ v \left| \int \exp_{\kappa} \left( lpha \left| v \right| 
ight) \ d\mu \leq +\infty, ext{for some } lpha > 0 
ight\} = \mathcal{L}^{2}(\mu)$$

#### Definition

As  $\exp_\kappa$  is convex, for  $p,q\in\mathcal{P}_>$ 

$$\int \exp_{\kappa}\left(\left(1- heta
ight) \ln_{\kappa}\left(p
ight)+ heta \ln_{\kappa}\left(q
ight)
ight) \; d\mu < +\infty, \quad heta \in [0,1].$$

If it holds for  $\theta \in I$ , open  $I \supset [0, 1]$ , then  $p \smile q$ .

## Deformed model

#### Example

If p = 1, we want to check

$$\int \exp_{\kappa}\left( heta \ln_{\kappa}\left( q
ight) 
ight) \; d\mu < +\infty, \quad heta > 1,$$

which is equivalent to

$$p^{1/2} + p^{-1/2} \in L^2(\mu)$$

#### lf

$$q=\exp_{\kappa}\left(u-\mathcal{K}_{p}(u)+\ln_{\kappa}\left(p
ight)
ight),\quad u\in L^{\kappa}(\mu)\cap L^{1}_{0}(\chi(p))$$

then  ${{\ln }_{\kappa }}\left( p 
ight) - {{\ln }_{\kappa }}\left( q 
ight) \in {L^{\kappa }}(\mu )$  and

$$K_{p}(u) = \mathsf{E}_{\chi(p)} \left[ \mathsf{ln}_{\kappa} \left( p \right) - \mathsf{ln}_{\kappa} \left( q \right) \right] = D_{\kappa} \left( p \| q \right)$$
$$u = \mathsf{ln}_{\kappa} \left( q \right) - \mathsf{ln}_{\kappa} \left( p \right) + D_{\kappa} \left( p \| q \right)$$