

IGAIA IV, Liblice, June 2016

In honor of Amari-sensei

Revisit to the Autoparallelity
and the Canonical Divergence
for Dually Flat Spaces

Hiroshi Nagaoka

The University of Electro-Communications, Tokyo.

- A talk on a similar subject was given in a workshop held in Nara, March 2012 without details.
- A related result (but in a different context) will be presented in the forthcoming IEEE ISIT, July 2016, with the title "A characterization of statistical manifolds on which the relative entropy is a Bregman divergence".

First, I would like to review some of the fundamental concepts of information geometry, and then state the main result and explain the proof of some nontrivial parts.

A dually flat space

A manifold S equipped with (g, ∇, ∇^*) such that

(1) g is a Riemannian metric.

(2) ∇ and ∇^* are flat affine connections which are dual w.r.t. g :

$$Xg(Y, Z) = g(\nabla_X Y, Z) + g(Y, \nabla_X^* Z).$$

(3) S is covered by a global ∇ -affine chart θ and a global

∇^* -affine chart η .

The canonical divergence

When θ and η are chosen to satisfy

$$g(\partial_i, \partial^j) = \delta_j^i, \quad \partial_i := \frac{\partial}{\partial \theta^i}, \quad \partial^j := \frac{\partial}{\partial \eta_j}$$

the canonical divergence $D: S \times S \rightarrow \mathbb{R}$ is defined by

$$(1) \quad D(p||q) = \varphi(p) + \psi(q) - \sum_i \eta_i(p) \theta^i(q)$$

by functions $\psi, \varphi: M \rightarrow \mathbb{R}$ satisfying

$$(2) \quad \eta_i = \partial_i \psi,$$

$$(3) \quad \theta^i = \partial^i \varphi,$$

$$(4) \quad \varphi + \psi = \sum_i \eta_i \theta^i$$

Another characterization of the canonical divergence

The canonical divergence is also defined as a function $D : S \times S \rightarrow \mathbb{R}$ such that for any $p, q, r \in S$

$$(1) D(p\|q) \geq 0$$

$$(2) D(p\|q) = 0 \Leftrightarrow p = q$$

$$(3) D(p\|q) + D(q\|r) - D(p\|r) = \sum_i \{\eta_i(p) - \eta_i(q)\} \{\theta^i(r) - \theta^i(q)\}$$

Note: (3) is the necessary and sufficient condition for mappings $\theta : S \rightarrow \mathbb{R}^{\dim S}$ and $\eta : S \rightarrow \mathbb{R}^{\dim S}$ to be a pair of dual affine charts.

Autoparallelity

- For a manifold S with an affine connection ∇ and a submanifold $M \subset S$,

M is ∇ -autoparallel (in S).

$\stackrel{\text{def}}{\iff} \forall X, Y$: vector fields on M , $\nabla_X Y$ is a vector field on M .

\iff The restriction $\nabla|_M$ becomes an affine connection on M .

$\iff M$ is ∇ -totally geodesic (when ∇ is torsion-free).

- We denote

$$M \stackrel{\nabla}{\subset} S \stackrel{\text{def}}{\iff} M \text{ is } \nabla\text{-autoparallel in } S$$

- When S is ∇ -flat with a ∇ -affine chart $\theta : S \rightarrow \mathbb{R}^{\dim S}$,

$$M \stackrel{\nabla}{\subset} S \iff \theta(M) \text{ forms (an open subset of) an affine subspace of } \mathbb{R}^{\dim S}.$$

$$\implies M \text{ is } \nabla|_M\text{-flat. (or we simply say } M \text{ is } \nabla\text{-flat.)}$$

Autoparallel submanifolds of a dually flat space

When

(1) (S, g, ∇, ∇^*) is dually flat with canonical divergence D and

(2) either $M \stackrel{\nabla}{\subset} S$ or $M \stackrel{\nabla^*}{\subset} S$,

then

(3) M is dually flat w.r.t. $(g|_M, \nabla|_M, \pi_M(\nabla^*))$ or $(g|_M, \pi_M(\nabla), \nabla^*|_M)$

(where π_M is the g -orthogonal projection onto M)

(4) The canonical divergence of M is $D|_{M \times M}$.

$\mathcal{P}(\mathcal{X})$ and its autoparallel submanifolds

— as a representative example —

- \mathcal{X} : an arbitrary finite set
- $S = \mathcal{P}(\mathcal{X}) := \{p \mid p : \mathcal{X} \rightarrow (0, 1), \sum_x p(x) = 1\}$
- $g = g^{(\text{F})}$: Fisher metric
- $\nabla = \nabla^{(\text{e})}$: exponential connection, e-connection
- $\nabla^* = \nabla^{(\text{m})}$: mixture connection, m-connection
- $D = D_{\text{KL}}$: $D_{\text{KL}}(p \parallel q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$
- We abbreviate $\nabla^{(\text{e})}$ - and $\nabla^{(\text{m})}$ - as **e-** and **m-**, respectively:
e.g., e-geodesic, m-autoparallel , etc.

$\mathcal{P}(\mathcal{X})$ and its autoparallel submanifolds; cont.

- $M \overset{\text{e}}{\subset} \mathcal{P}(\mathcal{X}) \iff M$ is an **exponential family** on \mathcal{X} ;
(**e-family** for short)

$$\log p_{\theta}(x) = C(x) + \sum_i \theta^i F_i(x) - \psi(\theta)$$

$\theta = (\theta^i)$: an e-affine chart

$$\eta_i(p) = \mathbb{E}_p[F_i(X)] = \sum_x F_i(x)p(x)$$

$\eta = (\eta_i)$: an m-affine chart

- $M \overset{\text{m}}{\subset} \mathcal{P}(\mathcal{X}) \iff M$ is a **mixture family** on \mathcal{X} ;
(**m-family** for short)

$$p_{\eta}(x) = C(x) + \sum_i \eta_i F^i(x)$$

$\eta = (\eta_i)$: an m-affine chart

$$\theta^i(p) = \sum_x F^i(x) \log p(x)$$

$\theta = (\theta^i)$: an e-affine chart \uparrow_0

Given a dually flat space (S, g, ∇, ∇^*) with canonical divergence D ,

- If $K \overset{\nabla}{\subset} S$ and $M \overset{\nabla|_K}{\subset} S$, then
 - $M \overset{\nabla}{\subset} S$,
 - M is dually flat w.r.t. $(g|_M, \nabla|_M, \pi_M(\nabla^*))$
with canonical divergence $D|_{M \times M}$.

- For simplicity, we write this implication as

$$M \overset{\nabla}{\subset} K \overset{\nabla}{\subset} S \implies M \overset{\nabla}{\subset} S \quad \text{and} \quad M \text{ is dually flat} \\ \text{with canonical divergence } D.$$

Note: the canonical divergence determines the dually flat structure.

- Similarly,

$$M \overset{\nabla^*}{\subset} K \overset{\nabla^*}{\subset} S \implies M \overset{\nabla^*}{\subset} S \quad \text{and} \quad M \text{ is dually flat} \\ \text{with canonical divergence } D.$$

- On the other hand,

$$\text{either } M \overset{\nabla^*}{\subset} K \overset{\nabla}{\subset} S \quad \text{or} \quad M \overset{\nabla}{\subset} K \overset{\nabla^*}{\subset} S \\ \implies \text{not } M \overset{\nabla}{\subset} S \text{ nor } M \overset{\nabla^*}{\subset} S \text{ in general,}$$

but M is still dually flat with canonical divergence D .

For instance

- $M \overset{e}{\subset} K \overset{e}{\subset} \mathcal{P}(\mathcal{X}) \implies M \overset{e}{\subset} \mathcal{P}(\mathcal{X})$ (M is an e-family) and M is dually flat with canonical divergence D_{KL} .
- $M \overset{m}{\subset} K \overset{m}{\subset} \mathcal{P}(\mathcal{X}) \implies M \overset{m}{\subset} \mathcal{P}(\mathcal{X})$ (M is an m-family) and M is dually flat with canonical divergence D_{KL} .
- either $M \overset{m}{\subset} K \overset{e}{\subset} \mathcal{P}(\mathcal{X})$ or $M \overset{e}{\subset} K \overset{m}{\subset} \mathcal{P}(\mathcal{X})$
 \implies not $M \overset{e}{\subset} \mathcal{P}(\mathcal{X})$ nor $M \overset{m}{\subset} \mathcal{P}(\mathcal{X})$ in general,
but M is dually flat with canonical divergence D_{KL} .

Main theorem

Given a dually flat space (S, g, ∇, ∇^*) with canonical divergence D and a submanifold $M \subset S$, the following conditions are equivalent.

(1) M is dually flat with canonical divergence D .

(2) $M \stackrel{\nabla^*}{\subset} \exists K \stackrel{\nabla}{\subset} S$.

(3) $M \stackrel{\nabla}{\subset} \exists K \stackrel{\nabla^*}{\subset} S$.

(4) $\exists K_1 \stackrel{\nabla}{\subset} S$ and $\exists K_2 \stackrel{\nabla^*}{\subset} S$ such that

$$M = K_1 \cap K_2 \quad \text{and} \quad \forall p \in M, T_p(K_1)^\perp \perp T_p(K_2)^\perp$$

(5) $\exists K_1 \stackrel{\nabla}{\subset} S$ and $\exists K_2 \stackrel{\nabla^*}{\subset} S$ such that

$$M = K_1 \cap K_2 \quad \text{and} \quad \exists p \in M, T_p(K_1)^\perp \perp T_p(K_2)^\perp$$

We show that

(1) M is dually flat with canonical divergence D .

implies

(2) $M \stackrel{\nabla^*}{\subset} \exists K \stackrel{\nabla}{\subset} S$.

- Let $m := \dim M \leq n := \dim S$.
- Since (S, g, ∇, ∇^*) is dually flat, there exist
 - a ∇ - affine chart $\sigma : S \rightarrow \mathbb{R}^{n \times 1}$ (column vectors) and
 - a ∇^* -affine chart $\zeta : S \rightarrow \mathbb{R}^{1 \times n}$ (row vectors) such that

$$\forall p, q, r \in S$$

$$D(p||q) + D(q||r) - D(p||r) = (\zeta(p) - \zeta(q)) (\sigma(r) - \sigma(q)).$$

- Assume (1), so that there exists a pair of affine charts
 - $\theta : M \rightarrow \mathbb{R}^{m \times 1}$ and $\eta : M \rightarrow \mathbb{R}^{1 \times m}$ satisfying

$$\forall p, q, r \in M$$

$$D(p||q) + D(q||r) - D(p||r) = (\eta(p) - \eta(q)) (\theta(r) - \theta(q)).$$

- Fix a point $p_0 \in M$ arbitrarily, and let

$$V := \text{span} \{ \sigma(p) - \sigma(p_0) \mid p \in M \} \subset \mathbb{R}^{n \times 1}$$

and

$$K := \{ p \in S \mid \sigma(p) - \sigma(p_0) \in V \} \subset S.$$

Note: the definitions of V and K do not depend on p_0 .

- It is obvious that

$$M \subset K \overset{\nabla}{\subset} S.$$

- Let $k := \dim V = \dim K$. ($m \leq k \leq n$)

Then there exists a matrix $F : n \times k$ such that $V = \text{Im} F$.

- Define a ∇ -affine chart $\rho : K \rightarrow \mathbb{R}^{k \times 1}$ (column vectors) of K by

$$\forall p \in K, \quad \sigma(p) - \sigma(p_0) = F \rho(p).$$

- Let $\xi : K \rightarrow \mathbb{R}^{1 \times k}$ (row vectors) be defined by

$$\forall p \in K, \quad \xi(p) = \zeta(p) F.$$

Proof of (1) \implies (2), 3/5 slides

- For $\forall p, q, r \in K$, we have

$$\begin{aligned}
 (\xi(p) - \xi(q)) (\rho(r) - \rho(q)) &= (\zeta(p) - \zeta(q)) F (\rho(r) - \rho(q)) \\
 &= (\zeta(p) - \zeta(q)) (\sigma(r) - \sigma(q)) \\
 &= D(p||q) + D(q||r) - D(p||r),
 \end{aligned}$$

which implies that ξ is ∇^* -affine chart of K .

manifold	dimension	∇ -chart (column vectors)	∇^* -chart (row vectors)
S	n	σ	ζ
K	k	ρ	ξ
M	m	θ	η

- Lemma:

$$\forall f \in V (\subset \mathbb{R}^{n \times 1}), \exists a \in \mathbb{R}^{m \times 1}, \forall p \in M,$$

$$(\zeta(p) - \zeta(p_0)) f = (\eta(p) - \eta(p_0)) a. \quad (\#)$$

\therefore) It suffices to show $(\#)$ in the case when $f = \sigma(q) - \sigma(p_0)$ for an arbitrary $q \in M$, since V is the linear span of such f 's. In this case, for $\forall p \in M$ we have

$$\begin{aligned} (\zeta(p) - \zeta(p_0)) f &= (\zeta(p) - \zeta(p_0)) (\sigma(q) - \sigma(p_0)) \\ &= D(p||p_0) + D(p_0||q) - D(p||q) \\ &= (\eta(p) - \eta(p_0)) (\theta(q) - \theta(p_0)), \end{aligned}$$

which means that $(\#)$ holds by setting $a := \theta(q) - \theta(p_0)$.

- Since $V = \text{Im}F$, the previous lemma implies the existence of a matrix $A : m \times k$ such that

$$\forall p \in M, (\zeta(p) - \zeta(p_0)) F = (\eta(p) - \eta(p_0)) A. \quad (b)$$

Proof of (1) \implies (2), 5/5 slides

- Since the LHS of (b) is $(\zeta(p) - \zeta(p_0)) F = \xi(p) - \xi(p_0)$, we have

$$\forall p \in M, \quad \xi(p) = \eta(p)A + b \quad (\natural)$$

where $b := \xi(p_0) - \eta(p_0)A$.

- (\natural) means that $\xi(M)$ forms an affine subspace of the ξ -coordinate space \mathbb{R}^k and hence $M \stackrel{\nabla^*}{\subset} K$. (QED)

manifold	dimension	∇ -chart (column vectors)	∇^* -chart (row vectors)
S	n	σ	ζ
K	k	ρ	ξ
M	m	θ	η

$$(2) \quad M \overset{\nabla^*}{\subset} \exists K \overset{\nabla}{\subset} S.$$

$$\implies (4) \quad \exists K_1 \overset{\nabla}{\subset} S \text{ and } \exists K_2 \overset{\nabla^*}{\subset} S \text{ such that}$$

$$M = K_1 \cap K_2 \quad \text{and} \quad \forall p \in M, T_p(K_1)^\perp \perp T_p(K_2)^\perp$$

- Assume that $M \overset{\nabla^*}{\subset} K \overset{\nabla}{\subset} S$ and fix a point $p_0 \in M$ arbitrarily. Let $\sigma : S \rightarrow \mathbb{R}^{n \times 1}$ and $\zeta : S \rightarrow \mathbb{R}^{1 \times n}$ be a pair of dual affine charts.
- Since $K \overset{\nabla}{\subset} S$, there exists a matrix $F : n \times k$ such that

$$K = \{p \in S \mid \sigma(p) - \sigma(p_0) \in \text{Im}F\}.$$

Let $\xi : K \rightarrow \mathbb{R}^{1 \times k}$ be defined by $\xi(p) = \zeta(p)F$ for $\forall p \in K$.

Then ξ becomes a ∇^* chart of K .

- Since $M \overset{\nabla^*}{\subset} K$, there exists a linear subspace W of $\mathbb{R}^{1 \times k}$ such that

$$\begin{aligned} M &= \{p \in K \mid \xi(p) - \xi(p_0) \in W\} \\ &= \{p \in K \mid \zeta(p)F - \xi(p_0) \in W\} = K_1 \cap K_2 \end{aligned}$$

where $K_1 := K$ and $K_2 := \{p \in S \mid \zeta(p)F - \xi(p_0) \in W\}$.

- $K_1 \overset{\nabla}{\subset} S$ and $K_2 \overset{\nabla^*}{\subset} S$ are obvious, and

$$\forall p \in M, T_p(K_1)^\perp \perp T_p(K_2)^\perp$$

can be verified by some linear algebraic argument.

Let us observe that the set of stationary markov joint distributions forms a dually flat space which can be regarded as an example of our theorem.

Markov joint distributions

- Let \mathcal{X} be an arbitrary finite set and let $\mathcal{X}^n = \underbrace{\mathcal{X} \times \cdots \times \mathcal{X}}_n$.
- $S := \mathcal{P}(\mathcal{X}^n)$, the set of positive n -joint distributions.
- An n -joint distribution $p \in S$ is **markov**

$\stackrel{\text{def}}{\iff} X_1 - X_2 - \cdots - X_n$ is a Markov process

for $X^n = (X_1, \dots, X_n) \sim p$.

$\iff \exists u_1, \dots, u_{n-1} : \mathcal{X}^2 \rightarrow \mathbb{R}^+, \quad \forall x^n = (x_1, \dots, x_n) \in \mathcal{X}^n,$

$$p(x^n) = \prod_{t=1}^{n-1} u_t(x_t, x_{t+1}).$$

$\iff \exists \pi \in \mathcal{P}(\mathcal{X}), \quad \exists w_1, \dots, w_{n-1} \in \mathcal{P}(\mathcal{X} | \mathcal{X}),$

$\forall x^n = (x_1, \dots, x_n) \in \mathcal{X}^n,$

$$p(x^n) = \pi(x_1) \prod_{t=1}^{n-1} w_t(x_{t+1} | x_t).$$

Markov joint distributions, cont.

- $K_{\text{mar}} := \{p \in S \mid p \text{ is markov}\}$.
- K_{mar} is an exponential family.
(\because the Hammersley-Clifford theorem)
- For $\mathcal{X} = \{0, 1, \dots, m-1\}$,

$$N := \dim K_{\text{mar}} = n(m-1) + (n-1)(m-1)^2,$$

and $K_{\text{mar}} = \{p_\rho \mid \rho = (\rho^{it}; \rho^{ijt}) \in \mathbb{R}^N\}$ where

$$\begin{aligned} \log p_\rho(x^n) &= \sum_{t=1}^n \sum_{i=1}^{m-1} \rho^{it} \delta_i(x_t) \\ &\quad + \sum_{t=1}^{n-1} \sum_{i=1}^{m-1} \sum_{j=1}^{m-1} \rho^{ijt} \delta_{ij}(x_t, x_{t+1}) - \psi(\rho). \end{aligned}$$

Stationary joint distributions

- For an n -joint distribution $p \in S = \mathcal{P}(\mathcal{X}^n)$, its marginal distributions $p_t^{(k)} \in \mathcal{P}(\mathcal{X}^k)$ for $k \in \{1, \dots, n-1\}$ and $t \in \{1, \dots, n-k+1\}$ are defined by

$$p_1^{(n-1)}(x^{n-1}) = \sum_{x'} p(x^{n-1}, x'),$$

$$p_2^{(n-1)}(x^{n-1}) = \sum_{x'} p(x', x^{n-1}),$$

$$p_1^{(n-2)}(x^{n-2}) = \sum_{x'} p_1^{(n-1)}(x^{n-2}, x'),$$

$$p_2^{(n-2)}(x^{n-2}) = \sum_{x'} p_1^{(n-1)}(x', x^{n-2}) = \sum_{x'} p_2^{(n-1)}(x^{n-2}, x'),$$

$$p_3^{(n-2)}(x^{n-2}) = \sum_{x'} p_2^{(n-1)}(x', x^{n-2}),$$

.....

- An n -joint distribution $p \in S$ is **stationary**

$$\stackrel{\text{def}}{\iff} p_1^{(n-1)} = p_2^{(n-1)}$$

$$\iff \forall k \in \{1, \dots, n-1\}, p_t^{(k)} \text{ does not depend on } t.$$

Stationary joint distributions, cont.

- $K_{\text{sta}} := \{p \in \mathcal{S} \mid p \text{ is stationary}\}$.
- K_{sta} is a mixture family.

Stationary markov joint distributions

- $M := K_{\text{mar}} \cap K_{\text{sta}}$.
- $p \in M \iff \exists w \in \mathcal{P}(\mathcal{X}|\mathcal{X})$ s.t. $\forall x^n \in \mathcal{X}^n$,

$$p(x^n) = p_w^{(n)}(x^n) := \pi_w(x_1) \prod_{t=1}^{n-1} w(x_{t+1}|x_t)$$

where π_w is the stationary distribution of w :

$$\pi_w(x) = \sum_{x'} w(x|x')\pi_w(x').$$

- It does not hold that $\forall p \in M, T_p(K_{\text{mar}})^\perp \perp T_p(K_{\text{sta}})^\perp$.

- Let

$$K_{\text{sta},2} := \left\{ p \in \mathcal{S} \mid p_t^{(2)} \text{ does not depend on } t \right\} \quad (\supset K_{\text{sta}}).$$

Then

- $K_{\text{sta},2}$ is a mixture family.
- $M = K_{\text{mar}} \cap K_{\text{sta},2}$.
- $\forall p \in M, T_p(K_{\text{mar}})^\perp \perp T_p(K_{\text{sta},2})^\perp$.
- Therefore, we have
 - M is dually flat with canonical divergence D_{KL} .
 - $M \stackrel{\text{m}}{\subset} K_{\text{mar}}$ and $M \stackrel{\text{e}}{\subset} K_{\text{sta},2}$.

- $D = D_{\text{KL}}|_{M \times M}$ is represented as

$$D(p_{w_1}^{(n)} \| p_{w_2}^{(n)}) = D(\pi_{w_1} \| \pi_{w_2}) + (n-1) \sum_x \pi_{w_1}(x) D(w_1(\cdot|x) \| w_2(\cdot|x))$$

- M is not an e-family nor m-family.

Stationary markov joint distributions, cont. cont.

- Recall that $K_{\text{mar}} = \{p_\rho \mid \rho \in \mathbb{R}^N\}$ with

$$\begin{aligned} \log p_\rho(x^n) &= \sum_{t=1}^n \sum_{i=1}^{m-1} \rho^{it} \delta_i(x_t) \\ &\quad + \sum_{t=1}^{n-1} \sum_{i=1}^{m-1} \sum_{j=1}^{m-1} \rho^{ijt} \delta_{ij}(x_t, x_{t+1}) - \psi(\rho). \end{aligned}$$

- There exists a subset $U \subset \mathbb{R}^N$ such that $M = \{p_\rho \mid \rho \in U\}$.
- A pair of dual affine charts of M is given by

– e-affine chart $\theta = (\theta^i, \theta^{ij})$:

$$\theta^i = \sum_{t=1}^n \rho^{it}, \quad \theta^{ij} = \sum_{t=1}^{n-1} \rho^{ijt} \quad (\forall \rho \in U)$$

– m-affine chart: $\eta = (\eta_i, \eta_{ij})$:

$$\eta_i(p) = E_p [\delta_i(X_t)], \quad \eta_{ij}(p) = E_p [\delta_{ij}(X_t, X_{t+1})]$$

(do not depend on t)

Thank you for listening!