統計数理研　2013-9-26

## 情報幾何入門

理研 脳科学総合研究センター

甘利俊一

# Manifold of Probability Distributions

Fisher metric  :  1929  H. Hotelling
                         1945  C. R. Rao

Invariance     :  1972  N. N. Chentsov
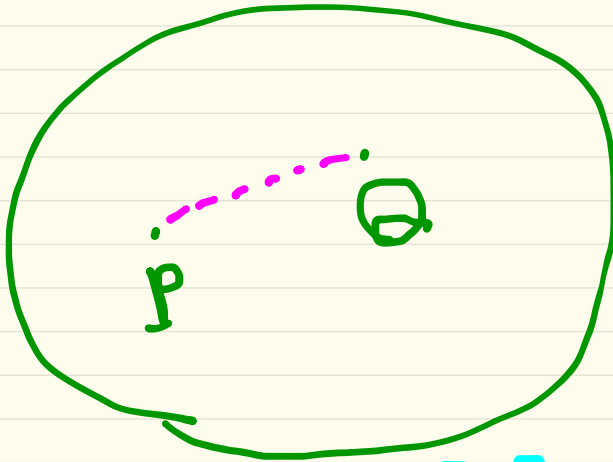  connections

Curvature  :  1978  B. Efron  (A. P. Dawid)

Duality  :  1982  S. Amari
                         H. Nagaoka & S. Amari

 .... many others     G. Pistone          J. L. Koszul

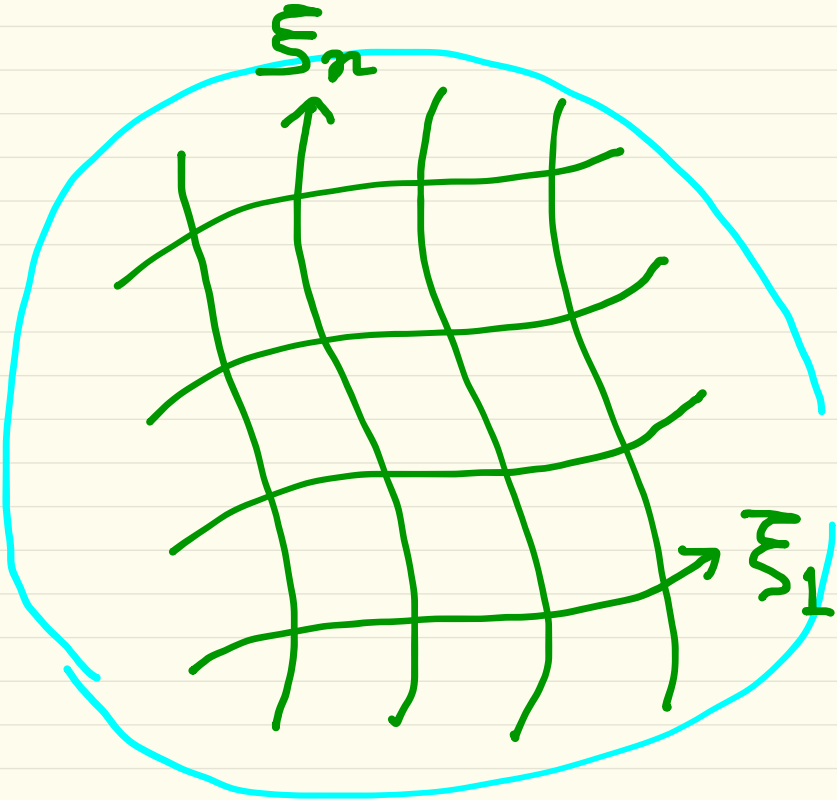 applications

# Manifold and Divergence

## Dually Flat Structure

$$D[P:Q] \geq 0$$

distance

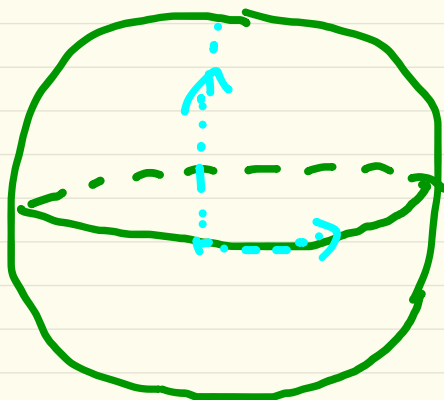$$D[P:Q] \neq D[Q:P]$$

# Manifold



$\xi_n$

$\xi_1$

coordinate system

$$\xi = (\xi_1, \cdots, \xi_n)$$

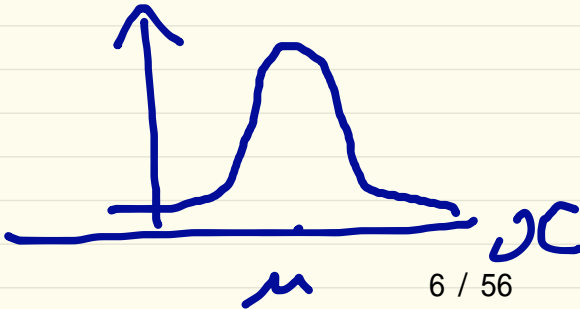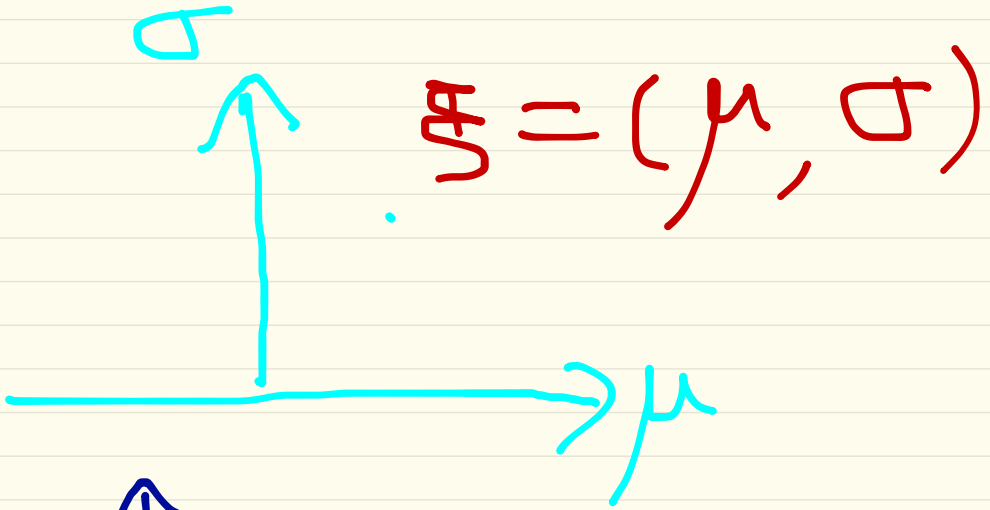# Euclidean space

$$\xi_2$$

$$\xi_1$$

# Sphere

latitude
longitude

# Probability Distributions

## Gaussian distribution

$$P(x, \xi) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

$$\xi = (\mu, \sigma)$$

# Probability Simplex $S_n$

$$S_n = \{ P(x) \}$$

$$x = \{ 0, 1, 2, \cdots, n \}$$

$$P = (P_0, \cdots, P_n)$$

$$P_i = \text{Prob} \{ x = i \}$$

$P_3$

$P_2$

$S_3$
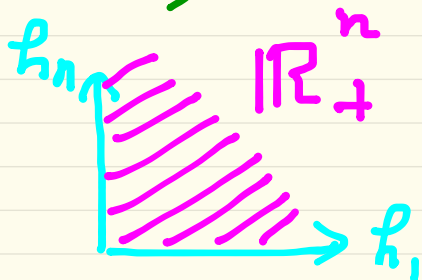
$P_0$

$P_1$

# histgram

$$h = (h_1, \cdots, h_n), \quad h_i > 0$$

$$\mathbb{R}_+^n$$



# vision

$$S(x, y)$$

$$S_{ij} = S(i, j)$$

$$\mathbb{R}_+^{n^2}$$

# Positiv-definite matrix

$$\{P\} \qquad P_{ij}, \ (i \le j)$$

# Neural networks

$$\{w_{ij}\}$$

# Coordinate transformation



$$\xi = f(\zeta)$$

$$\begin{cases} \text{invariant} \\ \text{convenient} \end{cases}$$

$$\xi_i = f_i(\zeta_1, \cdots, \zeta_n)$$

# Manifold with Divergence

$$D[P:Q] \qquad D[\xi_P : \xi_Q]$$



M

P    Q

$\xi$

1) $D[P:Q] \geq 0$ ; $0$ iff $P=Q$

2) $D[P:P+dP] = \frac{1}{2} g_{ij} d\xi^i d\xi^j$

Positive-definite

$$D[P:Q] \neq D[Q:P]$$

in general

# Riemannian metric

$$e_i = \frac{\partial}{\partial \xi_i} \qquad \langle e_i, e_j \rangle = g_{ij}$$



$$dS^2 = g_{ij} \, d\xi_i \, d\xi_j = \frac{1}{2} D[\xi : \xi + d\xi]$$

$T_{ijk}$ : <u>cubic, symmetric</u>

# アファイン接続



$$e_i' \sim e_i$$

$$\triangle e_i = \Pi e_i' - e_i$$

$$\nabla_{e_j} e_i = \Gamma_{ji}{}^k e_k$$

# Dual connections $(\Gamma, \Gamma^*)$ $(\nabla, \nabla^*)$



Parallel transport

$$\Pi A , \Pi^* B$$

$$\langle A, B \rangle = \langle \Pi A, \Pi^* B \rangle$$

$$(\langle A, B \rangle = \langle \Pi A, \Pi B \rangle) : \text{Levi-Civita} .$$

# Convex Function

$$z = \psi(\xi)$$



$$\psi(\alpha \xi_1 + (1-\alpha)\xi_2) < \alpha \psi(\xi_1)$$
$$+ (1-\alpha)\psi(\xi_2)$$

$$g_{ij} = \frac{\partial^2 \psi(\xi)}{\partial \xi_i \, \partial \xi_j} > 0$$

# Bregman Divergence



$$z = \psi(\xi') + \nabla\psi(\xi') \cdot (\xi - \xi')$$

$$D[\xi : \xi'] = \psi(\xi) - \psi(\xi')$$
$$- \nabla\psi(\xi') \cdot (\xi - \xi')$$

$$g_{ij} = \frac{\partial^2 \psi(\xi)}{\partial \xi_i \partial \xi_j} > 0$$

# Examples

$$\psi(\xi) = \frac{1}{2} \sum \xi_i^2 \Rightarrow \text{Euclid}$$

$$\psi(\xi) = -\sum \log \xi_i$$

$$\psi(p(x)) = \int p(x) \log p(x) \, dx$$

$$D[p(x) : q(x)] \quad \text{KL-divergence}$$

$$= \int p(x) \log \frac{p(x)}{q(x)} \, dx$$

指数型分布族
　　曲指数型分布族

$$P(x, \theta) = \exp\{\theta \cdot x - \psi(\theta)\}$$

$$\psi(\theta) : 凸函数$$

$$g_{ij} = \frac{\partial^2}{\partial\theta^i\partial\theta^j}\psi(\theta)$$

$$T_{ijk} = \frac{\partial^3}{\partial\theta^i\partial\theta^j\partial\theta^k}\psi(\theta)$$

函数空间

flat-divergence
⇒ exponential family

Banerjee et al.

$$P(x, \theta) = \exp\{\theta \cdot x - \psi(\theta)\}$$

$$\eta = \nabla\psi(\theta) = E[x]$$

$$D[\eta : \eta'] = \psi(\theta) + \varphi(\eta')$$
$$- \theta \cdot \eta \qquad \theta(\eta)$$

$$P(x, \theta) = \exp\{-D[x : \eta]\} \, b(x)$$

# Legendre Transformation

$$\eta = \nabla \psi(\theta)$$

$$\theta = \nabla \varphi(\eta)$$



$$\theta \Longleftrightarrow \eta$$

$$\psi(\theta) + \varphi(\eta) - \theta \cdot \eta = 0$$

$$D[\theta : \theta'] = \psi(\theta) + \varphi(\eta')$$
$$- \theta \cdot \eta'$$

$$\varphi(\eta) = \max_{\theta} \{ \theta \cdot \eta - \psi(\theta) \}$$

# Affine Coordinates
## flatt $(\theta, \eta)$

$\theta$ : flat

$\eta$ : dual flat

biorthogonal



$$\langle e_i, e^{*j} \rangle$$

$$= \delta_{ij}$$

# Riemannian metric

$$d\theta = \sum d\theta_i \, \mathbb{e}_i$$

$$d\eta = \sum d\eta_i \, \mathbb{e}_i^*$$

$$ds^2 = \langle d\theta, d\theta \rangle$$

$$= \sum g_{ij} \, d\theta_i \, d\theta_j$$

$$g_{ij} = \langle \mathbb{e}_i, \mathbb{e}_j \rangle$$

$$ds^2 = \sum g_{ij}^* \, d\eta_i \, d\eta_j$$

$$g_{ij}^* = \langle \mathbb{e}_i^*, \mathbb{e}_j^* \rangle$$

$$G = G^{*-1}$$



$\theta$    $\theta + d\theta$

$\eta$    $\eta + d\eta$

# Pythagorean Theorem

dual flat
— geodesic
— dual geodesic



$$D[P:Q] + D[Q:R] = D[P:R]$$

Euclidean space
$$\Psi(\theta) = \frac{1}{2} \sum \theta_i^2 \; : \quad \eta_i = \theta_i$$

Proof

# Projection Theorem

$$\hat{P} = \underset{Q \in M}{\arg\min}\, D[P:Q]$$



$$D^*[P:Q] = D[Q:P]$$

$$\uparrow \varphi(\eta)$$

$$\hat{P}^* = \underset{Q \in M}{\arg\min}\, D[Q:P]$$

# Applications :
## Statistical Inference

$M = \{P(x, \xi)\}$

$x_1, \cdots, x_N$



$\hat{\xi} = \text{argmin} \; D[\hat{P}_{emp} : Q]$

$\quad Q = P(x, \xi)$

D : KL-divergence : MLE (maximum likelihood estimator

# Invariance :

Probability distributions

$$D[P(x) : \tilde{g}(x)] \quad \text{Information monotone}$$

$$y = k(x)$$

$$P(x) \rightarrow \bar{P}(y)dy = P(x)dx$$

$$D[P(x) : \tilde{g}(x)] \geq \bar{D}[\bar{P}(y) : \bar{g}(y)]$$

$y$ : sufficient statistics

$$\Longleftrightarrow D = \bar{D}$$

$$P(x, \xi) = \bar{P}(y, \xi) \Gamma(x)$$

# f-divergence

$$f(1) = f'(1) = 0, \quad f''(1) = 1$$

$f(u):$ convex function

$$D_f[p(x) : g(x)] = \int p(x) f\left(\frac{g(x)}{p(x)}\right) dx$$

## invariant divergence

$$D_f[g(x) : p(x)] = D_{f^*}[p(x) : g(x)]$$

$$f^*(u) = u f\left(\frac{1}{u}\right)$$

$$f = -\log u \qquad KL\text{-divergence}$$

$S_n$ : discrete

$$\mathbb{P} = (P_0, P_1, \cdots, P_n). \quad P_i = \text{Prob}\{x = i\}$$

$$\underset{}{\text{o o o} \cdots\cdots \text{o}} \quad x$$

$$\underbrace{\text{o o o}}_{B_1} \, \underbrace{\quad}_{B_2} \, \underbrace{\text{o}}_{B_m} \quad y$$

<span style="color:red">course graining</span>

$$\mathbb{P} \longrightarrow \bar{\mathbb{P}} \qquad \bar{P}_\alpha = \text{Prob}\{y = \alpha\}$$

$$P(x) \to \bar{P}(y) \qquad\qquad = \text{Prob}\{x \in B_\alpha\}$$

$$= \sum_{i \in B_\alpha} P_i$$

# Information monotone

$$D[p : q] \geq \bar{D}[\bar{p} : \bar{q}]$$

$$y = k(x) \quad \text{sufficient}$$

$$\text{Prob}_p\{x \mid y \in B_\alpha\} = \text{Prob}_q\{x \mid y \in B_\alpha\}$$

$$D[p : q] = \bar{D}[\bar{p} : \bar{q}]$$

# decomposable divergence

$$D[P : Q] = \sum k(P_i, Q_i)$$

Theorem    decomp. invariant d

$$: D_f[P : Q] = \sum P_i f\left(\frac{Q_i}{P_i}\right)$$

# Flat & Invariant Divergence

↑        ↑

Bregman    $f-$
$$\begin{cases} \theta = \log \dfrac{P_i}{P_0} \\ \\ \eta = P_i \end{cases}$$

## KL-divergence

$$D_{KL}[P : Q] = \sum P(x) \log \frac{P(x)}{Q(x)}$$

inv      flat      KL-divergence

# $\alpha$-divergence

$$f_\alpha(u) = \frac{-4}{1-\alpha^2} u^{\frac{1+\alpha}{2}}$$

$$D_\alpha[P:q] = \sum_i \left\{ \frac{1-\alpha}{2} p_i + \frac{1+\alpha}{2} q_i - p_i^{\frac{1-\alpha}{2}} q_i^{\frac{1+\alpha}{2}} \right\}$$

$$= 1 - \sum_i p_i^{\frac{1-\alpha}{2}} q_i^{\frac{1+\alpha}{2}} \qquad : p.d.$$

$\alpha = -1$     KL-divergence

$\alpha = 1$     dual KL

$\alpha = 0$     Hellinger

$$\frac{1}{2} \sum_i \left( \sqrt{p_i} - \sqrt{q_i} \right)^2$$

# Flat & Invariant Div.
## in $\underline{\mathbb{R}_+^n}$ : positive measure


flat · invariant · $\alpha$-divergence

$$\theta_i = P_i^{\frac{1+\alpha}{2}} \quad , \quad \eta_i = P_i^{\frac{1-\alpha}{2}}$$

$$D_\alpha[P : Q] = D_{-\alpha}[Q : P]$$

$$\Psi(\theta) = \sum \theta_i^{\frac{2}{1+\alpha}} \quad , \quad \varphi(\eta) = \sum \eta_i^{\frac{2}{1-\alpha}}$$

$$\alpha \longleftrightarrow -\alpha \quad \text{duality}$$

# α-structure

α-mean

α-family of prob. distribution

α - projection

α - optimality

Tsallis $q$-entropy

$$H(P) = \frac{1}{1-q} \left( \sum p_i^q - 1 \right)$$

$$\alpha = 2q - 1 \qquad \int p_i^{\frac{1+\alpha}{2}}$$

# $\alpha$- mean

$$x, y > 0$$

$$m_f(x, y) = f^{-1}\left(\frac{f(x) + f(y)}{2}\right)$$

scale-free

$$m_f(cx, cy) = c\, m_f(x, y)$$

$$f_\alpha(u) = \begin{cases} u^{\frac{1-\alpha}{2}} \\ \log u, \quad \alpha = 1 \end{cases}$$

$\alpha = 1$ : geometric mean $\quad \sqrt{xy}$

$\alpha = -1$ : arithmetic mean $\quad \dfrac{x+y}{2}$

$\alpha = 0$ : $\quad \dfrac{1}{2}\left(\dfrac{1}{2}(x+y) + \sqrt{xy}\right)$

$\alpha = 3$ : harmonic mean $\quad \dfrac{2}{\dfrac{1}{x} + \dfrac{1}{y}}$

$\alpha = \infty$ : $\quad \min\{x, y\}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad$ fuzzy

$\alpha = -\infty \quad \max\{x, y\}$

$\quad m_\alpha(x, y) \geq m_{\alpha'}(x, y), \quad \alpha \leq \alpha'$

pessimistic mean : optimistic mean

$\alpha$-family of Prob. distr.

$$\xi_1(x), \cdots, \xi_m(x)$$

$\Rightarrow$

$$P_\alpha(x; w) = C f_\alpha^{-1}\left\{ \sum_{i=1}^{m} w_i f_\alpha(\xi_i(x)) \right\}$$

$\alpha = -1 \qquad P_\alpha(x) = \sum w_i \xi_i(x)$

mixture family

$\alpha = 1 \qquad P_\alpha(x) = \exp\left\{ \sum w_i \log \xi_i(x) - \psi \right\}$

exp. family

$$\alpha\text{-integration of } g_1(x), \cdots, g_m(x)$$

$$P(x) = f_\alpha^{-1} \left\{ \sum w_i \, f_\alpha \{ P_i(x) \} \right\}$$

$$R[P(x)] = \sum w_i \, D[g_i(x) : P(x)]$$

$$\min \; R_\alpha [P(x)]$$

: $\alpha$-integration

mixture of
   experts

$g_2 \quad \cdot \quad g_1$

$g_3 \quad P(x) \quad g_4$

# flat divergence (non-invariant)

$(\alpha, \beta)$-divergence : $\mathbb{R}^n_+$

$$D_{\alpha\beta}[P:Q] = \frac{1}{\alpha\beta(\alpha+\beta)} \sum \left\{ \alpha P_i^{\alpha+\beta} + \beta Q_i^{\alpha+\beta} \right.$$

$$\left. - (\alpha+\beta) P_i^{\alpha} Q_i^{\beta} \right\}$$

$$\theta_i = \frac{1}{\alpha} P_i^{\alpha}, \quad \eta_i = \frac{1}{\beta} P_i^{\beta}$$

$$\psi(\theta) = c \sum \theta_i^{\frac{\alpha+\beta}{\alpha}}, \quad \varphi(\eta) = c \sum \eta_i^{\frac{\alpha+\beta}{\beta}}$$

# $(u, v)$-divergence

$u(s), v(s)$ : **monotone incr.**

$$\theta_i = u(p_i), \quad \eta_i = v(p_i)$$

$$\psi(\theta) = \sum \int^{\theta_i} v(p) u'(p) \, dp$$

$$\varphi(\eta) = \sum \int^{\eta_i} u(p) v'(p) \, dp$$

$$D[p:\theta] = \sum_i \left[ \int^{p_i} v(p) u'(p) \, dp + \int^{\delta_i} u v' \, dp \right.$$
$$\left. - u(p_i) v(\delta_i) \right]$$

# Center of a cluster

$$\boldsymbol{x}^* = \arg\min \sum_i D[\boldsymbol{x}, \boldsymbol{x}_i]$$

**K-means clustering**

# Total Bregman Divergence

$$TD[\boldsymbol{x}:\boldsymbol{y}] = \frac{D[\boldsymbol{x}:\boldsymbol{y}]}{\sqrt{1+\|\nabla\psi\|^2}}$$



- •rotational invariance
- •conformal geometry

Figure: $d_f(x, y)$ (dotted red line) is BD, $\delta_f(x, y)$ (bold green line) is TBD, and the two arrows indicate the coordinate system. Note that $d_f(x, y)$ changes with rotation unlike $\delta_f(x, y)$ which is invariant to rotation.

# Clustering : *t*-center

$$E = \left\{ x_1, \cdots, x_m \right\}$$

**T-center of E**

$$\boldsymbol{x}^* = \arg \min \sum_i TD[\boldsymbol{x}, \boldsymbol{x}_i]$$

$\bullet \, y$

$E$

$\times \boldsymbol{x}^*$

# *t*-center $x^*$

$$\nabla \psi \left( x^* \right) = \frac{\sum w_i \nabla \psi \left( x_i \right)}{\sum w_i}$$

$$w_i = \frac{1}{\sqrt{1 + \left\| \nabla \psi \left( x_i \right) \right\|^2}}$$

# Total Bregman divergence (Vemuri)

$$\mathrm{TBD}(p:q) = \frac{\varphi(p) - \varphi(q) - \nabla\varphi(q)\cdot(p-q)}{\sqrt{1 + |\nabla\varphi(q)|^2}}$$

# Conformal change of divergence

$$\tilde{D}\left(p:q\right)=\sigma\left(p\right)D\left[p:q\right]$$

$$\tilde{g}_{ij}=\sigma\left(p\right)g_{ij}$$

$$\tilde{T}_{ijk}=\sigma(T_{ijk}+s_k g_{ij}+s_j g_{ik}+s_i g_{jk})$$

$$s_i=\partial_i \log \sigma$$

# *t*-center is robust

$$E^* = \{ \boldsymbol{x}_1, \cdots, \boldsymbol{x}_n ; \boldsymbol{y} \}$$

$$\tilde{\boldsymbol{x}}^* = \boldsymbol{x}^* + \varepsilon \boldsymbol{z}\left( \boldsymbol{x}^* ; \boldsymbol{y} \right), \quad \varepsilon = \frac{1}{n}$$

influence function $\boldsymbol{z}\left( \boldsymbol{x}^* ; \boldsymbol{y} \right)$

$$|\boldsymbol{z}| < c \ \text{ as } \ |\boldsymbol{y}| \rightarrow \infty \ : \text{robust}$$

# Positive-Definite Matrices

$(\alpha, \beta)$-divergence in $P = \{P > 0\}$

$$D_{\alpha, \beta}[P : Q] = tr\left\{ \frac{\alpha}{\alpha+\beta} P^{\alpha+\beta} + \frac{\beta}{\alpha+\beta} Q^{\alpha+\beta} \right.$$
$$\left. - P^{\alpha} Q^{\beta} \right\}$$

$D_{u,v}$

# flat divergence in $S_n$

: escort probability &
conformal geometry

$$\left( \begin{array}{c} \alpha\text{-exp.} \\ \text{fam.} \end{array} \right.$$

$$h_\alpha(P) = \sum P_i^{\frac{1+\alpha}{2}}$$

$$\tilde{D}_\alpha[P:Q] = \frac{2}{1-\alpha} \frac{1}{h_\alpha(Q)} \left[ 1 - \sum P_i^{\frac{1+\alpha}{2}} Q_i^{\frac{1-\alpha}{2}} \right]$$

$$\theta_i = \frac{2}{1-\alpha} \left[ P_i^{\frac{1-\alpha}{2}} - P_0^{\frac{1-\alpha}{2}} \right]$$

$$\eta_i = \frac{1}{h_\alpha(P)} P_i^{\frac{1+\alpha}{2}} \approx \tilde{P}$$

$$\tilde{g}_{ij}(\theta) = \frac{1}{h_\alpha(\theta)} g_{ij}(\theta)$$

conformal : unique (α-escort)

# Divergence and Geometry

$$D[\xi : \xi'] \qquad \nabla_\xi = \frac{\partial}{\partial \xi_i} , \; \nabla'_{\xi'} = \frac{\partial}{\partial \xi'_i}$$

## Riemmanian Metric

$$g_{ij} = - \nabla_{\xi_i} \nabla'_{\xi'_j} D[\xi : \xi']_{\xi = \xi'} \qquad : \text{positive-definite}$$
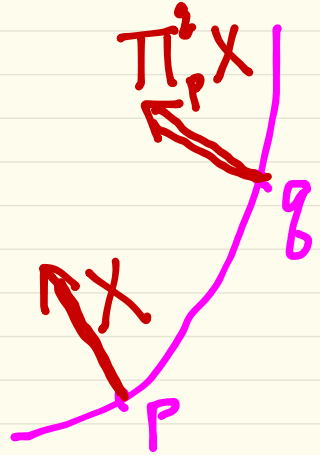
## cubic tensor

$$T_{ijk} = \nabla_{\xi_i} \nabla_{\xi_j} \nabla'_{\xi'_k} D - \nabla'_{\xi'_i} \nabla'_{\xi'_j} \nabla_{\xi_k} D$$

$$D \longrightarrow \{M, g, \dot{T}\} : T \text{ symmetric}$$

$$\{M, g\}$$

$$\Pi^2_p X$$

covariant derivatives

$$\nabla \Longleftrightarrow \nabla^*$$

$$X$$

$$q$$

$$p$$

$$\Pi^2_p X, \; \Pi^{*}_p{}^q X \quad : \text{parallel}$$

transport

# Duall Affine Connections

$$\Gamma_{ijk} = \{i, j ; k\} - \frac{1}{2} T_{ijk}$$

$$\Gamma^{*}_{ijk} = \{i, j ; k\} + \frac{1}{2} T_{ijk}$$

$$\Gamma^{\alpha}_{ijk} = \{i, j ; k\} - \frac{\alpha}{2} T_{ijk}$$

$$\pm \alpha \text{ duality}$$

$$D_{Z} \langle X, Y \rangle = \langle \nabla_{Z} X, Y \rangle + \langle X, \nabla^{*}_{Z} Y \rangle$$

$$\langle X, Y \rangle_{p} = \langle \Pi^{2}_{p} X, \Pi^{*}{}^{2}_{p} Y \rangle_{2}$$

# Two geodesics

$$\nabla_{\dot\xi} \dot\xi (t) = 0$$

$$\nabla_{\dot\xi}^* \dot\xi (t) = 0$$

$$\xi (t)$$
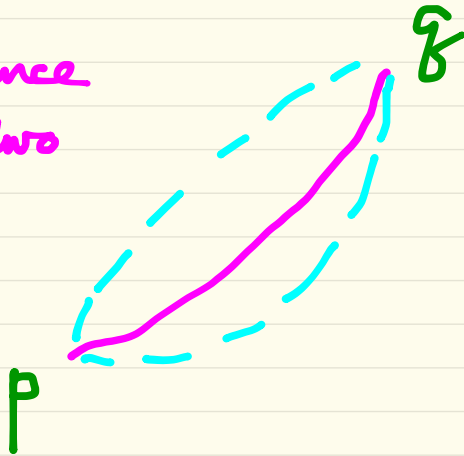
$$\ddot\xi_i + \sum \Gamma_{kji} \dot\xi_k \dot\xi_j = 0$$

Euclidean :    minimal distance
           straight

Riemannian : Levi-Civita connection

dual geometry :

   minimal distance
     splits into two

# Flat Manifold &
## Canonical Divergence

dually flat : $R = 0 \Leftrightarrow R^* = 0$

$\exists$ affine coordinates : $\theta, \eta$

$\exists$ convex functions : $\Psi(\theta), \varphi(\eta)$

canonical divergence

$$D[\theta : \theta'] = \Psi(\theta) + \varphi(\eta') - \theta \cdot \eta'$$

manifold of prob. distributions

invariance & flat $\Rightarrow$ KL-divergence

flat-divergence
⇒ exponential family

Banerjee et al.

$$p(x, \theta) = \exp\{\theta \cdot x - \psi(\theta)\}$$

$$\eta = \nabla \psi(\theta) = E[x]$$

$$D[\eta : \eta'] = \psi(\theta) + \varphi(\eta')$$
$$- \theta \cdot \eta \qquad \theta(\eta)$$

$$P(x, \theta) = \exp\{-D[x : \eta]\} \, b(x)$$

# Le Theorem

$\{M, g, T\}:$

realization in probability model

$$M = \{p(x, \xi)\}$$

embedding
curvature

invariance $\Rightarrow$

uniqueness of

$g_{ij}, T_{ijk}$

## Hessian manifold (Shima)

$$g_{ij}(\xi) = \nabla_{\xi_i} \nabla_{\xi_j} \psi(\xi)$$

$$\{M, g\} \Rightarrow \{M, g, T\}$$

dually flat ?

given $g_{ij} \Rightarrow^{\exists} T_{ijk}$ : flatten