2015 Workshop on High-Dimensional Statistical Analysis Dec.11 (Friday) ~15 (Tuesday) Humanities and Social Sciences Center, Academia Sinica, Taiwan

# Information Geometry and Spontaneous Data Learning

Shinto Eguchi

Institute Statistical Mathematics, Japan

This talk is based on a joint work with Osamu Komori and Atsumi Ohara, University of Fukui

# Outline

- Short review for Information geometry
- **Kolmogorov-Nagumo mean**
- $\phi$  -path in a function space
- Generalized mean and variance
- $\phi$  -divergence geometry
- U-divergence geometry
- Minimum *U*-divergence and density estimation

# A short review of IG

Nonparametric space

**Space of statistics** 

$$\mathcal{F}_{\mu} = \{ f : f(x) \ge 0, \int f(x) d\mu(x) = 1 \}$$
$$\mathcal{O}_{\mu} = \{ s(x) : \int |s(x)| d\mu(x) < \infty \}$$

where  $\mu$  is a fixed probability measure.

Information geometry is discussed on the product space  $\mathscr{F}_{\mu} \times \mathscr{O}_{\mu}$ 

$$\{ model \} \subseteq \mathscr{F}_{\mu}$$
$$IG \ \uparrow \downarrow$$
$$\{ estimators \} \subseteq \mathscr{O}_{\mu}$$

A canonical biliear form :  $\mathscr{F}_{\mu} \times \mathscr{Q}_{\mu} \to \mathbb{R}$  defined by

 $\langle f, s \rangle = \int f(x) s(x) d\mu(x)$ 

#### **Bartlett's identity**

**Parametric model**  $M = \{ f_{\theta}(x) \in \mathscr{F}_{\mu} : \theta = (\theta_1, \dots, \theta_d) \in \Theta \}$ 

**Bartlett's first identity** 

$$\left\langle f_{\theta}, \frac{\partial}{\partial \theta} \log f_{\theta} \right\rangle \to \operatorname{E}_{f_{\theta}} \left\{ \frac{\partial}{\partial \theta} \log f_{\theta}(X) \right\} = 0$$

**Bartlett's second identity** 

$$\operatorname{E}_{f_{\theta}} \{-\frac{\partial^{2}}{\partial \theta \partial \theta^{\mathrm{T}}} \log f_{\theta}(X)\} = \operatorname{Var}_{f_{\theta}} \{\frac{\partial}{\partial \theta} \log f_{\theta}(X)\}$$

#### **Metric and connections**

Let 
$$M = \{ f_{\theta}(x) \in \mathscr{F}_{\mu} : \theta = (\theta_1, \cdots, \theta_d) \in \Theta \}$$

**Information metric** 
$$G_{ij}(\theta) = \operatorname{Cov}_{f_{\theta}} \{S_i(X,\theta), S_j(X,\theta)\}$$
  
where  $S_i(x,\theta) = \frac{\partial}{\partial \theta_i} \log f_{\theta}(x)$ 

**Mixture connection** 

$$\Gamma_{ij,k}^{(\mathrm{m})}(\theta) = \int \frac{\partial^2 f_{\theta}(x)}{\partial \theta_i \partial \theta_i} S_k(x,\theta) \mathrm{d}\mu(x)$$

**Exponential connection** 

$$\Gamma_{ij,k}^{(e)}(\theta) = \mathbb{E}_{f_{\theta}} \left\{ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f_{\theta}(X,\theta) S_k(X,\theta) \right\}$$

Rao (1945), Dawid (1975), Amari (1982)

#### **Geodesic curves and surfaces in IG**

Let us fix f and g of  $\mathcal{F}_{\mu}$ . **m-geodesic curve**  $C^{(m)} = \{tf(x) + (1-t)g(x) : t \in (0, 1)\}$ **e-geodesic curve**  $C^{(e)} = \{ \exp\{t \log f(x) + (1-t) \log g(x) - \kappa_t \} : t \in (0, 1) \}$ Let us fix  $\{f_k\}_{k=1}^K$  in  $\mathcal{F}_{\mu}$ . **m-geodesic surface**  $M^{(m)} = \{\sum_{k=1}^{K} \pi_k f_k(x) : \pi \in S_K\}$  where  $S_{K-1}$  denotes a simplex. **e-geodesic surface**  $M^{(e)} = \{\exp\{\sum_{k=1}^{K} \pi_k \log f_k(x) - \kappa(\pi)\} : \pi \in S_K\}$  $f - C^{(e)}$  $C^{(m)}$ • π  $M^{(m)}$  $M^{(e)}$  $S_2$ 6

#### **Kullback-Leibler**

**K-L divergence** 
$$D_0(f,g) = \int f(x) \log \frac{f(x)}{g(x)} d\mu(x)$$

1. KL divergence is the expected log-likelihood ratio

$$D_0(f,g) = \mathcal{E}_f\left\{\log\frac{f(X)}{g(X)}\right\}$$

2. Maximum likelihood is minimum KL divergence.

Akaike (1974)

$$L_0(\theta) = \mathbb{E}_{\hat{f}_n} \left\{ \log f_{\theta}(X) \right\} = -D_0(\hat{f}_n, f_{\theta}) + \text{const}$$

**3.** KL divergence induces to the m-connection and e-connection

$$\Gamma_{ij,k}^{(m)}(\theta) = -\frac{\partial^{3}}{\partial\theta_{1i}\partial\theta_{1j}\partial\theta_{2k}} D_{0}(f_{\theta_{1}}, f_{\theta_{2}}) \Big|_{\substack{\theta_{1}=\theta\\\theta_{2}=\theta}}$$
  
$$\Gamma_{ij,k}^{(e)}(\theta) = -\frac{\partial^{3}}{\partial\theta_{1i}\partial\theta_{2j}\partial\theta_{2k}} D_{0}(f_{\theta_{2}}, f_{\theta_{1}}) \Big|_{\substack{\theta_{1}=\theta\\\theta_{2}=\theta}}$$
  
Eguchi (1983)

# **Pythagoras**

**Thm** Let us fix distinct densities f, g and h in  $\mathscr{F}_{\mu}$  $C^{(m)} = \{tf(x) + (1-t)g(x) : t \in (0, 1)\}$   $C^{(e)} = \{\exp\{t \log h(x) + (1-t)\log g(x) - \kappa_t\} : t \in (0, 1)\}$ 

If  $C^{(m)}$  and  $C^{(e)}$  orthogonally intersects at g, then

$$D_0(f,h) = D_0(f,g) + D_0(g,h)$$

Amari-Nagaoka (2001)



Pf 
$$G(\dot{C}^{(m)}, \dot{C}^{(e)}) = \int (f - g)(\log g - \log h - \dot{\kappa}_t)d\mu$$
  
=  $D_0(f, h) - \{D_0(f, g) + D_0(g, h)\}$ 

#### **Exponential model**

Let us fix 
$$\mathbf{t} = (t_1, \dots, t_K)$$
 where  $t_k \in \mathcal{Q}_{\mu}$  for  $k = 1, \dots, K$ 

**Exponential model**  $M^{(e)} = \{ f_{\theta}^{(e)}(x) := \exp\{\theta^{T} t(x) - \kappa(\theta)\} : \theta \in \Theta \}$  **Mean parameter**  $\eta = \mathbb{E}_{f_{\theta}^{(e)}}\{t(X)\} = \frac{\partial}{\partial \theta}\kappa(\theta)$  **For**  $\theta$   $G_{ij}(\theta) = \frac{\partial^{2}}{\partial \theta_{i}\partial \theta_{j}}\kappa(\theta)$   $\Gamma_{ij,k}^{(e)}(\theta) = 0$ **For**  $\eta$   $G_{ij}(\eta) = (G(\theta)^{-1})_{ji}$   $\Gamma_{ij,k}^{(m)}(\eta) = 0$  Amari (1982)

**Degenerated Bartlett identity** 

$$-\frac{\partial^2}{\partial\theta\partial\theta^{\mathrm{T}}}\log f_{\theta}^{(\mathrm{e})}(\boldsymbol{x}) = \operatorname{Var}_{f_{\theta}}(\frac{\partial}{\partial\theta}\log f_{\theta}^{(\mathrm{e})})$$

#### Minimum KL leaf

**Exponential model**  $M^{(e)} = \{ f_{\theta}^{(e)}(x) := \exp\{\theta^{T} t(x) - \kappa(\theta)\} : \theta \in \Theta \}$ Mean equal space  $\mathscr{Q}(g) = \{ f \in \mathscr{F}_{\mathcal{U}} : E_f \{ t(X) \} = E_g \{ t(X) \} \}$  $f \in \mathscr{Q}(f_{A^*}^{(e)}) \implies D_0(f, f_A^{(e)}) = D_0(f, f_{A^*}^{(e)}) + D_0(f_{A^*}^{(e)}, f_A^{(e)})$ Hence  $D_0(f, f_{\theta^*}) = \min_{\theta \in \Theta} D_0(f, f_{\theta}^{(e)})$  $D_{0}(f_{\theta^{*}}^{(e)}, f_{\theta}^{(e)})$   $D_{0}(f, f_{\theta}^{(e)})$   $D_{0}(f, f_{\theta}^{(e)})$  f



### **Pythagoras foliation**



11

#### $\log + \exp$



# Path geometry



{ m-geodesic, e-geodesic , ... }

#### Kolmogorov-Nagumo mean

K-N mean is 
$$\phi^{-1}\left(\frac{1}{n}\sum_{i=1}^{n}\phi(x_i)\right)$$
 for positive numbers  $\{x_1,...,x_n\}$ 

(1)  $\phi(s) = s$  sample mean :  $\frac{1}{n} \sum_{i=1}^{n} x_i$ 

(2) 
$$\phi(s) = \log(s)$$
 geometric mean :  $\exp\left(\frac{1}{n}\sum_{i=1}^{n}\log x_i\right) = (x_1 \cdots x_n)^{\frac{1}{n}}$   
(3)  $\phi(s) = s^{-1}$  harmonic mean :  $\frac{1}{\frac{1}{n}\sum_{i=1}^{n}\frac{1}{x_i}}$ 

K-N mean in  $\mathscr{F}_{\mu}$ 

#### Def. K-N mean

$$\phi^{-1}((1-t)\phi(f(x)) + t\phi(g(x))), \quad f,g \in \mathscr{F}_{\mu}, \ t \in [0,1]$$



Cf. Naudts (2009)

# *\phi*-path

### **Def.** $\phi$ -path connecting f and g

 $f_t(x,\phi) = \phi^{-1}((1-t)\phi(f(x)) + t\phi(g(x)) - \kappa_t),$ 

where  $\kappa_t$  is a normalizing factor to satisfy  $\int \phi^{-1}((1-t)\phi(f) + t\phi(g) - \kappa_t)d\mu = 1$ .

- **Thm** Assume that  $\phi$  is strictly increasing and concave. Then there exists such a normalizing constant  $\kappa_t$  as above.
- (Pf) Since  $\phi^{-1}$  is strictly inceasing and convex from the assumption  $\phi^{-1}((1-t)\phi(f) + t\phi(g)) \le (1-t)f + tg$   $\int \phi^{-1}((1-t)\phi(f) + t\phi(g))d\mu \le 1.$ Further  $\lim_{c \to \infty} \phi^{-1}(c) = \infty$

# **Examples of** $\phi$ -path

**Exm 0** 
$$\phi(s) = s \implies f_t(x, \phi) = (1-t)f(x) + tg(x)$$
 (m - geodesic)

Exm 1 
$$\phi_0(s) = \log s \Rightarrow$$
  
 $f_t(x, \phi_0) = \exp((1-t)\log f(x) + t\log g(x) - \kappa_t)$  (e-geodesic)  
where  $\kappa_t = \log \int \exp((1-t)\log f(x) + t\log g)d\mu$ .

Exm 2 
$$\phi_{\beta}(s) = \frac{s^{\beta} - 1}{\beta} \implies f_t(x, \phi_{\beta}) = \{(1 - t)f(x)^{\beta} + tg(x)^{\beta} - \kappa_t\}^{\frac{1}{\beta}}$$

Exm 3 
$$\phi_{\eta}(s) = \log(s+\eta) \Rightarrow$$
  
 $f_t(x,\phi_{\eta}) = \exp((1-t)\log(f(x)+\eta)+t\log(g(x)+\eta)-\kappa_t)$   
where  $\kappa_t = \log\left[\int \exp((1-t)\log(f+\eta)+t\log(g+\eta))d\mu-\eta\right]$ 

17

# **Identities of** *\phi***-density**

**Model** 
$$M = \{f_{\theta}(x) : \theta \in \Theta\}$$

$$\int \frac{1}{\phi'(f_{\theta})} \frac{\partial}{\partial \theta} \phi(f_{\theta}) d\mu = 0$$

2nd identity 
$$-\int \frac{\partial^2}{\partial \theta \partial \theta^{\mathrm{T}}} \phi(f_{\theta}) \frac{1}{\phi'(f_{\theta})} d\mu = \int \frac{-\phi''(f_{\theta})}{\{\phi'(f_{\theta})\}^2} \frac{\partial}{\partial \theta} f_{\theta} \frac{\partial}{\partial \theta^{\mathrm{T}}} \phi(f_{\theta}) d\mu$$

**because** 
$$\int \frac{\partial}{\partial \theta} f_{\theta} d\mu = 0 \qquad \int \frac{\partial^2}{\partial \theta \partial \theta^{\mathrm{T}}} f_{\theta} d\mu = 0$$

#### **Generalized mean and variance**

Def 
$$E_{f}^{(\phi)}\{a(X)\} = \int a(x)f^{(\phi)}(x)d\mu(x)$$
  
 $\operatorname{Cov}_{f}^{(\phi)}\{a(X), b(X)\} = \int w^{(\phi)}(f)(a - E^{(\phi)}a)(b - E^{(\phi)}b)^{\mathrm{T}}f^{(\phi)}d\mu$   
where  $f^{(\phi)}(x) = \frac{\overline{\phi'(f(x))}}{\int \frac{1}{\phi'(f)}d\mu} \quad w^{(\phi)}(f) = \frac{-\phi''(f)}{\{\phi'(f)\}^{3}}$ 

Note  $E_f^{(\phi)}$  is a linear functional on  $\mathscr{Q}_{\mu}$  $\operatorname{Cov}_f^{(\phi)}$  is a bilinear positive definite functional on  $\mathscr{Q}_{\mu}$ 

#### **Generalized mean and variance**

Exm 
$$\phi(s) = \log s \implies \operatorname{E}_{f}^{(\phi)}\{a(X)\} = \int a f \, d\mu$$
  
 $\operatorname{Cov}_{f}^{(\phi)}\{a(X), b(X)\} = \int (a - \operatorname{E}a)(b - \operatorname{E}b)^{\mathrm{T}} f \, d\mu$ 

$$\phi(s) = \frac{1}{\beta} (s^{\beta} - 1) \implies E_{f}^{(\phi)} \{a(X)\} = \frac{\int a f^{1-\beta} d\mu}{\int f^{1-\beta} d\mu}$$
$$Cov_{f}^{(\phi)} \{a(X), b(X)\} = \frac{\int (a - E^{(\phi)}a)(b - E^{(\phi)}b)^{T} f^{1-2\beta} d\mu}{\int f^{1-\beta} d\mu}$$

20

# **Bartlett Identity**

Model

$$M = \{f_{\theta}(x)\}_{\theta \in \Theta}$$

**Bartlett identity** 

$$\begin{split} & \mathrm{E}_{f_{\theta}} \{ \frac{\partial}{\partial \theta} \log f_{\theta} \} = 0 \\ & \mathrm{E}_{f_{\theta}} \{ -\frac{\partial^{2}}{\partial \theta \partial \theta^{\mathrm{T}}} \log f_{\theta} \} = \mathrm{Var}_{f_{\theta}} \left( \frac{\partial}{\partial \theta} \log f_{\theta} \right) \end{split}$$

*φ*–Bartlett identities

$$\mathrm{E}_{f_{\theta}}^{(\phi)}\{\frac{\partial}{\partial \theta_{j}}\phi(f_{\theta})\}=0$$

$$\mathbf{E}_{f}^{(\phi)} \{ -\frac{\partial^{2}}{\partial\theta\partial\theta^{\mathrm{T}}} \phi(f_{\theta}) \} = \mathrm{Var}_{f}^{(\phi)} \{ \frac{\partial}{\partial\theta} \phi(f_{\theta}) \}$$

# **Tangent space of** $\mathscr{F}_{\mu}$

**Tangent space** $T_f^{(\phi)} = \{a(x) \in \mathscr{D}_\mu : \mathbb{E}_f^{(\phi)}\{a(X)\} = 0\}$ **Riemannian metric** $\langle a, b \rangle_f^{(\phi)} = \operatorname{Cov}_f^{(\phi)}\{a(X), b(X)\}$ 

Expectation gives the tangent space. Topological properties of  $T_f^{(\phi)}$  depend on  $\phi$ . If  $\phi = \log$ , then  $T_f^{(\phi)}$  is too large to do statistics on  $\mathscr{F}_{\mu}$ Cf. Pistone (1992)

#### **Parallel transport**

Def

**A vector field**  $\{A_t(x): t \in [0,1]\}$  is parallel along a curve  $C = \{f_t: t \in [0,1]\}$ 

$$\Leftrightarrow \qquad \frac{d}{dt}A_t(x) = \operatorname{Cov}_{f_t}^{(\phi)} \left\{ A_t, \frac{d}{dt}\phi(f_t) \right\}$$

A curve  $C = \{f_t : t \in [0,1]\}$  is  $\phi$ -geodesic

$$\Leftrightarrow \qquad \frac{d^2}{dt^2}\phi(f_t(x)) = \operatorname{Var}_{f_t}^{(\phi)}\left(\frac{d}{dt}\phi(f_t(X))\right)$$

Cf. Amari (1982).



# *\$*-geodesic

- **Thm** If  $C = \{f_t\}_{0 \le t \le 1}$  is the  $\phi$ -geodesic curve then *C* is the  $\phi$ -path connecting  $f_0$  with  $f_1$ .
- **Proof.** By definition,  $\frac{d^2}{dt^2}\phi(f_t(x)) = \operatorname{const} \operatorname{in} x \quad (\forall t \in (0, 1))$ which is solved by

$$\phi(f_t(x)) = (1-t)\phi(f_1(x)) + t\phi(f_0(x)) - \kappa_t.$$



# *\phi*-divergence

$$\phi$$
-cross entropy  $C^{(\phi)}(f,g) = -E_f^{(\phi)}\{\phi(g)\}$ 

$$\phi$$
-entropy  $H^{(\phi)}(f) = C^{(\phi)}(f, f)$ 

*\phi*-divergence

$$D^{(\phi)}(f,g) = C^{(\phi)}(f,g) - H^{(\phi)}(f) = E_f^{(\phi)}\{\phi(f) - \phi(g)\}$$

Note:  $\phi$ -divergence is KL-divergence if  $\phi = \log \phi$ 

### **Divergence geometry**

**Def.** *D* is said to be a divergence measure on  $\mathscr{F}_{\mu}$  if  $D(f,g) \ge 0 \quad (\forall f,g \in \mathscr{F}_{\mu})$  with equality iff f = g (a.e. $\mu$ )

Let  $M = \{f_{\theta} : \theta \in \Theta\}$  be a statistical model.

$$D \mid_{\Theta \times \Theta} \to (G_{ij}^{(D)}, \Gamma_{ij,k}^{(D)}, {}^*\Gamma_{ij,k}^{(D)}) \text{ on } \Theta$$

with the Riemannian metric on *M* :

$$G_{ij}^{(D)}(\theta) = -\frac{\partial^2}{\partial \theta_{1i} \partial \theta_{2j}} D(f_{\theta_1}, f_{\theta_2}) \big|_{\substack{\theta \mid 1 = \theta \\ \theta \mid 2 = \theta}}$$

the pair of affine connections on *M*:

$$\Gamma_{ij,k}^{(D)}(\theta) = -\frac{\partial^3}{\partial \theta_{1i} \partial \theta_{1j} \partial \theta_{2k}} D(f_{\theta_1}, f_{\theta_2}) \big|_{\substack{\theta_1 = \theta \\ \theta_2 = \theta}} (\theta \in \Theta)$$

$$^* \Gamma_{ij,k}^{(D)}(\theta) = -\frac{\partial^3}{\partial \theta_{1i} \partial \theta_{1j} \partial \theta_{2k}} D(f_{\theta_2}, f_{\theta_1}) \big|_{\substack{\theta_1 = \theta \\ \theta_2 = \theta}} (\theta \in \Theta)$$

26

# *\phi*-divergence geometry

**The metric** 
$$G_{ij}^{(\phi)}(\theta) = \int \frac{\partial}{\partial \theta_i} \phi(f_\theta) \frac{\partial}{\partial \theta_j} \frac{1}{\phi'(f_\theta)} d\mu$$

#### Affine connection pair

$$\Gamma_{ij,k}^{(\phi)}(\theta) = \int \frac{\partial}{\partial \theta_k} \phi(f_\theta) \frac{\partial^2}{\partial \theta_j \partial \theta_k} \frac{1}{\phi'(f_\theta)} d\mu$$
$$*\Gamma_{ij,k}^{(\phi)}(\theta) = \int \frac{\partial^2}{\partial \theta_i \partial \theta_j} \phi(f_\theta) \frac{\partial}{\partial \theta_k} \frac{1}{\phi'(f_\theta)} d\mu$$

### $\phi$ – Pythagorean theorem

Thm  $\phi$ -geodesic  $f_t(\mathbf{x}) = \phi^{-1}((1-t)\phi(g(\mathbf{x})) + t\phi(f(\mathbf{x})))$  $\phi^*$ -geodesic  $h_s(\mathbf{x})$  s.t.  $\phi'(h_s) = \frac{1}{(1-s)\frac{1}{\phi'(g)} + s\frac{1}{\phi'(h)}}$ 

$$G^{(\phi)}(\dot{f}_{t},\dot{h}_{s})|_{t=0,s=0} = 0 \implies D_{\phi}(f,g) = D_{\phi}(f,g) + D_{\phi}(g,h)$$



$$\begin{aligned} \mathbf{Pf} & G^{(\phi)}(\dot{f}_{t}, \dot{h}_{s})|_{t=0,s=0} \\ &= \int \frac{d}{dt} \phi(f_{t}) \frac{d}{ds} \frac{1}{\phi'(g_{s})} d\mu|_{t=0,s=0} \\ &= \int \{\phi(f) - \phi(g)\} \{\frac{1}{\phi'(g)} - \frac{1}{\phi'(h)}\} d\mu \\ &= D_{\phi}(f, h) - \{D_{\phi}(f, g) + D_{\phi}(g, h)\} \end{aligned}$$

# $\phi$ -Pythagorean foliation

If f satisfies  $E_{f}^{(\phi)}\{t(X)\} = E_{f_{\theta^{*}}^{(\phi)}}^{(\phi)}\{t(X)\}$ , then  $D_{\phi}(f, f_{\theta}^{(\phi)}) = D_{\phi}(f, f_{\theta^{*}}^{(\phi)}) + D_{\phi}(f_{\theta^{*}}^{(\phi)}, f_{\theta}^{(\phi)})$ 

 $\phi$ -mean equal space  $\mathscr{Q}(g) = \{f : \mathcal{E}_f^{(\phi)}\{t(X)\} = \mathcal{E}_g^{(\phi)}\{t(X)\}\}$ 





What is a statistical meaning of *\phi*-mean and *\phi*-variance?

 $\phi$  - independence,  $\phi$  - LLN,  $\phi$  - CLT,  $\phi$  - moment generating function, ..... <sub>30</sub>

# **U-divergence**

Let U be a function satisfying  $U'(s) = \phi^{-1}(s)$ 

**U-cross-entropy**  $C_U(f,g) = \int \{-f \phi(g) + U(\phi(g))\} d\mu$ 

*U*-entropy  $H_U(f) = C_U(f, f)$ 

*U*-divergence  $D_U(f,g) = C_U(f,g) - H_U(f)$ 

Note  $C_U(f,g) \ge H_U(f)$  or  $D_U(f,g) \ge 0$ 

Exm Let  $U_{\beta}(s) = \frac{1}{1+\beta} (1-\beta s)^{\frac{1+\beta}{\beta}}$ .

Then power entropy  $C_{U_{\beta}}(f,g) = -\frac{1}{\beta} \int f g^{\beta} d\mu + \frac{1}{\beta+1} \int g^{\beta+1} d\mu$ <br/>power divergence  $D_{U_{\beta}}(f,g) = \frac{1}{\beta} \int f (f^{\beta} - g^{\beta}) d\mu + \frac{1}{\beta+1} \int (f^{\beta+1} - g^{\beta+1}) d\mu$  31

#### **U-divergence geometry**

The metric associated with *U*-divergence:

$$G_{ij}^{(U)}(\theta) = \int \frac{\partial}{\partial \theta_i} f_{\theta} \frac{\partial}{\partial \theta_j} \phi(f_{\theta}) d\mu$$

Affine connections associated with U-divergence:

$$\Gamma_{ij,k}^{(U)}(\theta) = \int \frac{\partial^2}{\partial \theta_i \partial \theta_j} f_{\theta} \frac{\partial}{\partial \theta_k} \phi(f_{\theta}) d\mu$$
$$^* \Gamma_{ij,k}^{(U)}(\theta) = \int \frac{\partial}{\partial \theta_k} f_{\theta} \frac{\partial^2}{\partial \theta_j \partial \theta_k} \phi(f_{\theta}) d\mu$$

**Thm** (i) *U*-geodesic is mixture geodesic.

(ii)  $U^*$ -geodesic is  $\phi$  geodesic

#### *U*-geometry $\neq \phi$ -geometry

 $\phi$ -metric on a model M

$$G_{ij}^{(\phi)}(\theta) = \int \frac{\partial}{\partial \theta_i} \phi(f_\theta) \frac{\partial}{\partial \theta_j} \frac{1}{\phi'(f_\theta)} d\mu$$

U-metric on a model M

\*

$$G_{ij}^{(U)}(\theta) = \int \frac{\partial}{\partial \theta_i} f_{\theta} \frac{\partial}{\partial \theta_j} \phi(f_{\theta}) d\mu$$

*<i>ф*\*-connection

$$\Gamma_{ij,k}^{(\phi)}(\theta) = \int \frac{\partial^2}{\partial \theta_i \partial \theta_j} \phi(f_\theta) \frac{\partial}{\partial \theta_k} \frac{1}{\phi'(f_\theta)} d\mu$$

U\*-connection

$${}^{*}\Gamma^{(U)}_{ij,k}(\theta) = \int \frac{\partial^{2}}{\partial \theta_{i} \partial \theta_{j}} \phi(f_{\theta}) \frac{\partial}{\partial \theta_{k}} f_{\theta} d\mu$$

# **Triangle with** *D*<sub>*U*</sub>

The mixture geodesic  $f_t(\mathbf{x}) = (1-t)g(\mathbf{x}) + t f(\mathbf{x})$  $\phi$ -geodesic  $h_s(\mathbf{x}) = \phi^{-1}((1-s)\phi(g(\mathbf{x})) + s\phi(h(\mathbf{x})) - \kappa_s)$ 

 $G^{(U)}(\dot{f}_t, \dot{h}_s)|_{t=0,s=0} = 0 \implies D_U(f,g) = D_U(f,g) + D_U(g,h)$ 



$$\begin{aligned} \mathbf{f} & G^{(U)}(\dot{f}_{t},\dot{h}_{s})|_{t=0,s=0} \\ &= \int \frac{d}{dt} f_{t} \frac{d}{ds} \xi(g_{s}) d\mu|_{t=0,s=0} \\ &= \int (f-g) \{ \phi(h) - \phi(g) \} d\mu \\ &= D_{U}(f,h) - \{ D_{U}(f,g) + D_{U}(g,h) \} \end{aligned}$$

#### **U-estimation**

Let  $g(\mathbf{x})$  be a data density function with statistical model  $f_{\theta}(\mathbf{x})$ 

**U-loss function** 

$$L_U(\boldsymbol{\theta}) = C_U(f, f_{\boldsymbol{\theta}}) = - \operatorname{E}_f \{ \phi(f_{\boldsymbol{\theta}}) \} + \int U(\phi(f_{\boldsymbol{\theta}})) d\mu$$

**U-empirical loss function** 

$$L_U^{\text{emp}}(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \phi(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)) + \int U(\phi(f_{\boldsymbol{\theta}})) d\boldsymbol{\mu}$$

**U-estimator for**  $\theta$ 

$$\hat{\boldsymbol{\theta}}_U = \argmin_{\boldsymbol{\theta} \in \Theta} L_U^{emp}(\boldsymbol{\theta})$$

35

### U-estimator under *\phi*-model

$$\phi$$
-model  $M^{(\phi)} = \{f_{\theta}^{(\phi)}(\mathbf{x}) = \phi^{-1}(\theta^T \mathbf{t}(\mathbf{x}) - \kappa_{\theta}) : \theta \in \Theta\}$ 

**U-empirical loss function** 

$$L_{U}(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^{n} \phi(f_{\boldsymbol{\theta}}^{(\phi)}(\boldsymbol{x}_{i})) + \int U(\phi(f_{\boldsymbol{\theta}}^{(\phi)})) d\mu$$
$$= -\{ \boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\bar{t}} - \Phi(\boldsymbol{\theta}) \}$$
where  $\boldsymbol{\bar{t}} = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{t}(\boldsymbol{x}_{i}), \quad \Phi(\boldsymbol{\theta}) = \kappa(\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{t} - \kappa(\boldsymbol{\theta})) d\mu$ 

**U-estimator for** 
$$\boldsymbol{\theta}$$
  $\hat{\boldsymbol{\theta}}_U = \underset{\boldsymbol{\theta}\in\Theta}{\operatorname{arg solve}} \left\{ \frac{\partial}{\partial \boldsymbol{\theta}} \Phi(\boldsymbol{\theta}) = \bar{\boldsymbol{t}} \right\} = \left( \frac{\partial \Phi}{\partial \boldsymbol{\theta}} \right)^{-1} (\bar{\boldsymbol{t}})$ 

*U*-estimator under *\phi*-model has analogy with MLE under exponential model

### **Potential function** $\Phi$

**Def** We call 
$$\Phi(\theta) = \kappa(\theta) + \int U(\theta^{T} t - \kappa(\theta)) d\mu$$
  
the potential function on  $\phi$ -model

Note 
$$\frac{\partial}{\partial \theta} \Phi(\theta) = \frac{\partial}{\partial \theta} \kappa(\theta) + \int \{t(x) - \frac{\partial}{\partial \theta} \kappa(\theta)\} f_{\theta}^{(U)}(x) d\mu(x)$$
$$= \mathbb{E}_{f_{\theta}^{(\phi)}} \{t(X)\}$$

We define the mean parameter  $\eta$  by

$$\boldsymbol{\eta} = \frac{\partial}{\partial \boldsymbol{\theta}} \Phi(\boldsymbol{\theta})$$

**Cf.** 
$$\phi$$
-mean parameter  $E_{f_{\theta}}^{(\phi)} \{ t(X) \} = \frac{\partial}{\partial \theta} \kappa(\theta)$ 

**Thm** U-estimator for  $\eta$  is given by the sample mean  $\hat{\eta}_U = \bar{t} = \frac{1}{n} \sum_{i=1}^n t(x_i),$ 

# **Pythagoras foliation**

**Thm** If a pdf 
$$\hat{f}$$
 satisfies  $E_{\hat{f}}\{t(X)\} = \bar{t}$ , then  
 $D_U(\hat{f}, f_{\theta}) = D_U(\hat{f}, f_{\hat{\theta}_U}) + D_U(f_{\hat{\theta}_U}, f_{\theta}).$ 

Pf Consider  

$$\hat{f}_{t} = t\hat{f} + (1-t)f_{\hat{\theta}_{U}}^{(U)}, h_{s} = \phi^{-1}(s\phi(f_{\theta}^{(U)}) + (1-s)\phi(f_{\hat{\theta}_{U}}^{(U)}))$$
Then  

$$G^{(U)}(\dot{f}_{t}, \dot{h}_{s}) = \int (\hat{f} - f_{\hat{\theta}_{U}}^{(U)})\{\phi(f_{\hat{\theta}_{U}}^{(U)}) - \phi(f_{\theta}^{(U)})\}d\mu$$

$$= (\hat{\theta}_{U} - \theta)^{\mathrm{T}}(\mathrm{E}_{\hat{f}}\{t(X)\} - \mathrm{E}_{f_{\hat{\theta}_{U}}^{(U)}}\{t(X)\}) = 0$$

### **Pythagoras foliation**

$$D_{U}(\hat{f}, f_{\theta}^{(U)}) = D_{U}(\hat{f}, f_{\theta_{U}}^{(U)}) + D_{U}(f_{\theta_{U}}^{(U)}, f_{\theta}^{(U)})$$

 $\mathcal{F}_{\mu} = \bigcup_{f \in M^{(U)}} \mathcal{Q}(f) \quad \text{where} \quad \mathcal{Q}(f) = \{ g : E_g \{ t(X) \} = E_f \{ t(X) \} \}$ 



#### **U-Boost learning for density estimation**

**U-loss function** 
$$L_U(f) = -\frac{1}{n} \sum_{i=1}^n \phi(f(\mathbf{x}_i)) + \int_{\mathbb{R}^p} U(\phi(f(\mathbf{x}))) d\mathbf{x}$$

**Dictionary of density functions** 

$$\mathcal{D} = \{ g_{\lambda}(\mathbf{x}) : g_{\lambda}(\mathbf{x}) \ge 0, \int g_{\lambda}(\mathbf{x}) d\mathbf{x} = 1, \ \lambda \in \Lambda \}$$

Learning space = *\phi*-model

$$\mathcal{D}_{\phi}^{*} = \phi^{-1} \big( \operatorname{co}(\phi(\mathcal{D})) \big) = \{ \phi^{-1} \big( \sum_{\lambda \in \Lambda} \pi_{\lambda} \phi(g_{\lambda}(\boldsymbol{x})) \big) \}$$

• Let 
$$f(\mathbf{x}, \boldsymbol{\pi}) = \phi^{-1} \Big( \sum_{\lambda \in \Lambda} \pi_{\lambda} \phi(g_{\lambda}(\mathbf{x})) \Big)$$
. Then  $f(\mathbf{x}, (0, \dots, 1, \dots, 0)) = g_{\lambda}(\mathbf{x})$   
•  $\mathcal{D}_{\phi}^{*} \supseteq \mathcal{D}$ 

**Goal: find** 
$$f^* = \operatorname{argmin}_{f \in \mathcal{D}_U^*} L_U(f)$$

#### **U-Boost algorithm**

(A) Find 
$$f_1 = \underset{g \in \mathcal{D}}{\operatorname{arg\,min}} L_U(g)$$

(B) Update 
$$f_k \longrightarrow f_{k+1} = \xi^{-1}((1-\alpha_{k+1})\xi(f_k) + \alpha_{k+1}\xi(g_{k+1}))$$
 st  

$$(\alpha_{k+1}, g_{k+1}) = \underset{(\alpha, g) \in (0,1) \times \mathcal{D}}{\operatorname{arg\,min}} L_U(\xi^{-1}((1-\alpha)\xi(f_k) + \alpha\xi(g)))$$

(C) Select *K*, and 
$$\hat{f} = \xi^{-1}((1 - \alpha_K)\xi(f_{K-1}) + \alpha_K\xi(g_K))$$

Example 2. Power entropy 
$$f^*(x) = \left(\sum_k \pi_k g_k(x)^\beta\right)^{\frac{1}{\beta}}$$
  
If  $\beta = 0$ ,  $f^*(x) = \exp\left(\sum_k \pi_k \log g_k(x)\right) = \prod_k g_k(x)^{\pi_k}$  Friedman et al (1984)  
If  $\beta = 1$ ,  $f^*(x) = \sum_k \pi_k g_k(x)$  Klemela (2007)

Inner step in the convex hull

$$\mathcal{D} = \{t(x,\lambda) : \lambda \in \Lambda\} \longrightarrow \mathcal{D}_U^* = \phi^{-1}(\operatorname{co}\phi(\mathcal{D}))$$

**Goal**:  $f^* = \underset{f \in \mathcal{D}_U^*}{\operatorname{argmin}} L_U(f) \quad f^*(x) = \phi^{-1}(\hat{\pi}_1 \phi(\hat{f}_1(x)) + \dots + \hat{\pi}_{\hat{k}} \phi(\hat{f}_{\hat{k}}(x)))$ 



 $g_1$ 

# **Non-asymptotic bound**

**Theorem.** Assume that a data distribution has a density g(x) and that

(A) 
$$\sup_{(\psi,\xi,\varphi)\in\operatorname{co}(\xi(\mathscr{D}))\times\mathscr{D}\times\mathscr{D}}\int U''(\psi)\{\phi(\xi)-\phi(\varphi)\}^2\leq b_U$$

Then we have

$$\mathbb{E}_{g}D_{U}(g,\hat{f}_{K}) \leq \mathrm{FA}(g,\mathcal{D}_{U}^{*}) + \mathrm{EE}(g,\mathcal{D}) + \mathbb{E}(K,\mathcal{D}),$$

where

$$FA(g, \mathcal{D}_{U}^{*}) = \inf_{f \in \mathcal{D}_{U}^{*}} D_{U}(g, f) \qquad (Functional approximation)$$

$$EE(g, \mathcal{D}) = 2 \mathbb{E}_{g} \{ \sup_{f \in \mathcal{D}} \left| \frac{1}{n} \sum_{i=1}^{n} \phi(f(\mathbf{x}_{i})) - \mathbb{E}_{p} \phi(f) \right| \} \quad (Estimation error)$$

$$IE(K, \mathcal{D}) = \frac{c^{2} b_{U}^{2}}{K + c - 1} \quad (c: step-length constant) \quad (Iteration effect)$$

Remark. Trade between  $FA(p, \mathcal{D}_U^*)$  and  $EE(p, \mathcal{D})$ 





# Conclusion



#### **Future problems**

#### **Tangent space**

$$T_{f}^{(\phi)} = \{a(x) \in H_{f}^{(\phi)} : \mathcal{E}_{f}^{(\phi)}\{a(X)\} = 0\}$$

where 
$$\langle a, b \rangle_f^{(\phi)} = \operatorname{Cov}_f^{(\phi)} \{a(X), b(X)\}$$

#### **Path space**

$$\mathcal{P}_{f}^{(\phi)} = \{ P^{(\phi)}(f,g) \colon g \in \mathcal{F}_{\Lambda} \}$$

where  $P^{(\phi)}(f,g) = \{ \phi^{-1}((1-t)\phi(f) + t\phi(g) - \kappa_t ) : t \in [0,1] \}$ 

#### **Future problems**

 $\phi$  - independence,  $\phi$  - LLN,  $\phi$  - CLT,  $\phi$  - moment generating function, .....

Cf. q - independec....

# Thank you