## Computational information geometry on Bregman manifolds and submanifolds

#### Frank Nielsen

#### Sony Computer Science Laboratories, Inc.

Applied geometry for data sciences Part II Singapore 2<sup>nd</sup> June 2025



### Outline of the talk

- Bregman divergences with some extensions
- Geometry of Bregman balls
- Two applications on Bregman manifolds:
  - Jensen-Shannon centroid on a mixture family manifold
  - Chernoff information/point on an exponential family manifold

#### Bregman divergences (1960's)

• F:  $\Theta \subseteq \mathbb{R}^m \rightarrow \mathbb{R}$  a strictly convex and smooth real-valued function on a finite dim. Hilbert space <.,.>

**Bregman divergence**  $B_F: \Theta \times \text{RelInt}(\Theta) \rightarrow \mathbb{R}_{\geq 0}$ 

$$\mathsf{B}_{\mathsf{F}}(\theta_1:\theta_2) \!=\! \mathsf{F}(\theta_1) \!-\! \mathsf{F}(\theta_2) \!-\! < \theta_1 \!-\! \theta_2, \nabla \mathsf{F}(\theta_2) \!>$$



Lev M. Bregman (1941 - 2023) Photo: courtesy of Alexander Fradkov

Smooth measure of discrepancy, not a metric distance because it violates the triangle inequality, and is asymmetric when F is not quadratic function. Hence the delimiter notation ":" instead of  $B_F(\theta_1, \theta_2)$ 

BD interpreted as **remainder** of a first order Taylor expression of  $F(\theta_1)$  around  $\theta_2$ :  $F(\theta_1) = F(\theta_2) + \langle \theta_1 - \theta_2, \nabla F(\theta_2) \rangle + \underbrace{B_F(\theta_1 : \theta_2)}_{Taylor remainder}$ 

Example of remainder: Lagrange remainder (smooth C<sup>2</sup> generators):  $\nabla^2 \mathbf{F} \mathbf{SPD} \Rightarrow B_F(\theta_1 : \theta_2) \ge 0$ 

$$\mathsf{B}_{\mathsf{F}}(\theta_1:\theta_2) = \frac{1}{2} \ (\theta_2 - \theta_1)^\top \nabla^2 \mathsf{F}(\theta) \ (\theta_2 - \theta_1) \ge 0 \ , \ \theta \in [\theta_1 \ \theta_2]$$

#### BDs: Versatile and popular in OR, ML, IT, signal processing

Originally motivated for finding an intersection point in a set of convex objects using **Bregman projections**. (ex. of convex objects: halfspaces, balls, etc.)

BDs unify:

- squared Euclidean divergence  $F(\theta) = \frac{1}{2} \Sigma_i < \theta, \theta > \theta$
- Kullback-Leibler divergence  $F(\theta) = \Sigma_i \theta_i \log(\theta_i)$ (relative Shannon entropy)
- *Itakura-Saito divergence*  $F(\theta) = \Sigma_i \log(\theta_i)$ (relative Burg entropy)

 $\mathsf{B}_{\mathsf{F}}(\theta_1:\theta_2) = \mathsf{F}(\theta_1) - \mathsf{F}(\theta_2) - <\theta_1 - \theta_2, \nabla \mathsf{F}(\theta_2) >$ 



L22 ( $\beta = 2$ ), KLD ( $\beta \rightarrow 0$ ), ISD ( $\beta = 1$ ), belong to a *family* of  $\beta$ -divergences, learn ad hoc  $\beta \ge 0$ 

$$\begin{array}{ll} \mathsf{x},\mathsf{y} \! > \! 0, \ \beta \geq \! 0 & d_{\beta}(x|y) = \left\{ \begin{array}{ll} \frac{x}{y} - \log\left(\frac{x}{y}\right) - 1 & \beta = 0 \\ x(\log x - \log y) + (y - x) & \beta = 1 \\ \frac{x^{\beta} + (\beta - 1)y^{\beta} - \beta xy^{\beta - 1}}{\beta(\beta - 1)} & \beta \in \mathbb{R} \backslash \{0, 1\} \end{array} \right. \begin{array}{ll} \text{Bregman} & \text{Bregman} \\ \text{Generator: } \phi_{\beta}(x) = \left\{ \begin{array}{ll} -\log x + x - 1 & \beta = 0 \\ x\log x - x + 1 & \beta = 1 \\ \frac{x^{\beta}}{\beta(\beta - 1)} - \frac{x}{\beta - 1} + \frac{1}{\beta} & \text{otherwise.} \end{array} \right. \end{array} \right.$$



Example: Bregman chord divergence, application: zero-order optimization in ML

The chord gap divergence and a generalization of the Bhattacharyya distance, IEEE ICASSP 2018

#### Bregman divergences in machine learning…

- Kullback-Leibler divergence between two probability densities: D<sub>KL</sub>[p(x):q(x)]= ∫ p(x) log (p(x)/q(x)) dμ(x) is difficult to calculate in closed form because of the integral ∫ ...
- But Kullback-Leibler divergence between two probability densities of a **natural exponential family** with densities  $p(x|\theta) \propto exp(\langle x, \theta \rangle)$ amount to a **reverse Bregman divergence**  $B_F^{rev}(\theta_1; \theta_2) := B_F(\theta_2; \theta_1)$  $D_{KL}[p(x|\theta_1): p(x|\theta_2)] = B_F^{rev}(\theta_1; \theta_2) = B_F(\theta_2; \theta_1)$

Bypass the  $\int, \nabla F$  in BD easy to calculate!  $\Rightarrow$  Easy calculations of KLDs

Azoury, Katy S., and Manfred K. Warmuth. "Relative loss bounds for on-line density estimation with the exponential family of distributions." *Machine learning* 43 (2001)

### Representational Bregman divergences (2009)

• Use a **representation function** R :

$$\begin{split} \mathsf{B}_{\mathsf{F},\mathsf{R}}(\lambda_1:\lambda_2) &:= \mathsf{B}_\mathsf{F}(\mathsf{R}(\lambda_1):\mathsf{R}(\lambda_2)) \\ &= \mathsf{F}(\mathsf{R}(\lambda_1)) - \mathsf{F}(\mathsf{R}(\lambda_2)) - < \mathsf{R}(\lambda_1) - \mathsf{R}(\lambda_2), \nabla \mathsf{F}(\mathsf{R}(\lambda_2)) > \\ & \text{Note that } \mathsf{F} \circ \mathsf{R} \text{ may not be a Bregman generator, i.e., not be strictly convex.} \end{split}$$

For example, consider the KLD between two densities of a generic exponential family (natural parameter from representation function)  $p_{\lambda}(x) \propto \tilde{p}_{\lambda}(x) = \exp(\langle \theta(\lambda), t(x) \rangle) h(x)$  include normal, Gamma/Beta, Wishart, Poisson, etc.  $\theta(\lambda)$ : natural parameter corresponding to  $\lambda$ , representation function R(.)= $\theta(.)$ 

 $\mathsf{D}_{\mathsf{KL}}[\mathsf{p}(\mathsf{x}|\lambda_1):\mathsf{p}(\mathsf{x}|\lambda_2)] = \mathsf{B}_{\mathsf{F}}^{\mathsf{rev}}(\theta_1(\lambda_1):\theta_2(\lambda_2)) = \mathsf{B}_{\mathsf{F}}(\theta_1(\lambda_2):\theta_2(\lambda_1))$ 

 $\mathsf{NEF} \, \mathsf{density} \, p(\mathbf{x} \big| \, \theta) \, \boldsymbol{\backsim} \, \mathsf{exp}(< \mathbf{x}, \theta >) \quad \mathsf{D}_{\mathsf{KL}}[p(\mathbf{x} | \theta_1) : p(\mathbf{x} | \theta_2)] = \mathsf{B}_{\mathsf{F}}^{\mathsf{rev}}(\theta_1 : \theta_2) = \mathsf{B}_{\mathsf{F}}(\theta_2 : \theta_1)$ 

#### Extended $\alpha$ -divergences are representational BDs

α-divergences extended to m-dimensional positive measures are representational Bregman divergences:

$$D_{\alpha}^{+}(q_{1}:q_{2}) = \begin{cases} \frac{4}{1-\alpha^{2}} \sum_{i=1}^{m} \left(\frac{1-\alpha}{2}q_{1} + \frac{1+\alpha}{2}q_{2} - q_{1}^{\frac{1-\alpha}{2}}q_{2}^{\frac{1+\alpha}{2}}\right), \alpha \in \mathbb{R} \setminus \{-1,1\} \\ D_{\mathrm{KL}}^{*}^{*}(q_{1}:q_{2}) = D_{\mathrm{KL}}^{+}(q_{2}:q_{1}) = \sum_{i=1}^{m} q_{2}^{i} \log \frac{q_{2}^{i}}{q_{1}^{i}} + q_{1}^{i} - q_{2}^{i} & \alpha = 1 \\ D_{\mathrm{KL}}^{+}(q_{1}:q_{2}) = \sum_{i=1}^{m} q_{1}^{i} \log \frac{q_{1}^{i}}{q_{2}^{i}} + q_{2}^{i} - q_{1}^{i} & \alpha = -1. \end{cases}$$

$$D_{\alpha}^{+}(q_{1}:q_{2}) = B_{F_{\alpha}}(R_{\alpha}(q_{1}):R_{\alpha}(q_{2}))$$

$$\begin{array}{ll} \text{Bregman generator:} & F_{\alpha}(r) = \sum_{i=1}^{m} f_{\alpha}(r_{i}), \quad f_{\alpha}(x) = \begin{cases} \frac{2}{1+\alpha} \left(\frac{1-\alpha}{2}x\right)^{\frac{2}{1-\alpha}}, & \alpha \neq 1\\ \log x, & \alpha = 1. \end{cases} \\ \text{Representation function:} & R_{\alpha}(q) = (r_{\alpha}(q_{1}), \ldots, r_{\alpha}(q_{m})), \quad r_{\alpha}(x) = \frac{2}{1-\alpha}x^{\frac{1-\alpha}{2}} \\ \text{Bregman divergence:} & \mathsf{B}_{\mathsf{F}}(\theta_{1}:\theta_{2}) = \mathsf{F}(\theta_{1}) - \mathsf{F}(\theta_{2}) - <\theta_{1} - \theta_{2}, \, \nabla \mathsf{F}(\theta_{2}) > \end{cases}$$

#### "The dual Voronoi diagrams with respect to representational Bregman divergences." IEEE ISVD 2009

### Convex duality via Legendre-Fenchel transform

- Legendre-Fenchel transform of a convex function F:  $F^*(\eta) = \sup_{\theta \in \Theta} \{ < \theta, \eta > -F(\theta) \}$
- Problem: some *tricky functions* with gradient map  $\nabla F$  domain not convex. Example:  $h(\xi_1, \xi_2) = [(\xi_1^2/\xi_2) + \xi_1^2 + \xi_2^2]/4$  on upper plane domain  $\Xi = (\xi_1, \xi_2)$
- Thus, we consider "nice convex functions" = Legendre-type functions (Θ,F(θ))
   (i) Θ open, and (ii) lim <sub>θ→∂Θ</sub> || ∇F(θ) || =∞

Then we get:

- **1** reciprocal gradient maps  $\eta = \nabla F(\theta)$  and  $\theta = \nabla F^*(\eta)$ ,  $\nabla F^* = (\nabla F)^{-1}$
- **2** conjugation yields  $(H,F^*(\eta))$  of Legendre type
- **3** biconjugation is an **involution**:  $(H,F^*(\eta))^* = (H^*=\Theta,F^{**}=F(\theta))$
- Convex conjugate:  $F^*(\eta) = \langle \nabla F^{-1}(\eta), \eta \rangle F(\nabla F^{-1}(\eta))$  since  $\eta = \nabla F(\theta)$

### Fenchel-Young divergences & convex duality

- Young inequality: F ( $\theta_1$ )+F\* ( $\eta_2$ )≥<  $\theta_1$ ,  $\eta_2$ > with equality when  $\eta_2 = \nabla F(\theta_1)$
- Build the Fenchel-Young divergence from the inequality: lhs-rhs ≥0

$$(\mathbf{Y}_{F,F^{*}}(\theta_{1}, \eta_{2})) = \mathbf{F}(\theta_{1}) + \mathbf{F^{*}}(\eta_{2}) - \langle \theta_{1}, \eta_{2} \rangle \geq 0$$

- Mixed parameterizations  $\theta$  and  $\eta$  :  $B_F(\theta_1:\theta_2) = Y_{F,F^*}(\theta_1, \eta_2)$
- Duality:  $B_F(\theta_1; \theta_2) = Y_{F, F^*}(\theta_1, \eta_2) = Y_{F^*,F}(\eta_2, \theta_1) = B_{F^*}(\eta_2, \eta_1)$
- Dual BDs + Dual FYs from involution F\*\*=F
- Note:  $B_F(\theta_1; \theta_2) = 0 \Leftrightarrow \theta_1 = \theta_2 \Leftrightarrow \eta_1 = \eta_2$  i.e.,  $\nabla F(\theta_1) = \nabla F(\theta_2)$

(FY initially called Legendre-Fenchel divergences…) 10

#### Bregman divergence vs Fenchel-Young divergence

Same parameterization  $B_F(\theta_1; \theta_2) = Y_{F, F^*}(\theta_1, \eta_2)$  mixed parameterization



 $\mathbf{Y}_{\mathsf{F},\,\mathsf{F}^*}(\boldsymbol{\theta}_{1,}\boldsymbol{\eta}_{2}) = \mathbf{F}(\boldsymbol{\theta}_{1}) + \mathbf{F^*}(\boldsymbol{\eta}_{2}) - <\boldsymbol{\theta}_{1,}\boldsymbol{\eta}_{2} >$ 

 $\mathsf{B}_{\mathsf{F}}(\boldsymbol{\theta}_1:\boldsymbol{\theta}_2) \!=\! \mathsf{F}(\boldsymbol{\theta}_1) \!-\! \mathsf{F}(\boldsymbol{\theta}_2) \!-\! < \! \boldsymbol{\theta}_1 \!-\! \boldsymbol{\theta}_2, \nabla \mathsf{F}(\boldsymbol{\theta}_2) \!>$ 

#### Kullback-Leibler divergence between non-normalized exponential family densities

- Kullback-Leibler divergence between two **positive measures**:  $D_{KI} + [p_1(x):p_2(x)] = \int \{p_1(x) \log (p_1(x)/p_2(x)) + p_2(x) - p_1(x)\} d\mu(x)$
- Exponential family density:
  - Normalized:  $p(x \mid \theta) = exp(\langle x, \theta \rangle F(\theta)) d\mu(x)$
  - Non-normalized:  $q(x | \theta) = exp(\langle x, \theta \rangle) d\mu(x)$
- Hence,  $p(x|\theta) = q(x|\theta)/Z(\theta)$  with **partition function**  $Z(\theta) = exp(F(\theta))$  and **cumulant function**  $F(\theta) = \log Z(\theta)$
- When F is convex, Z=exp(F) is log-convex
- log-convex functions are convex functions: So both F and Z are convex functions
- KLD between normalized densities = reverse Bregman wrt F:

$$\mathsf{D}_{\mathsf{KL}}[\mathsf{p}_{\theta 1}(\mathsf{x}):\mathsf{p}_{\theta 2}(\mathsf{x})] = \mathsf{B}_{\mathsf{F}}^*[\theta_1:\theta_2] = \mathsf{B}_{\mathsf{F}}[\theta_2:\theta_1]$$

• KLD between non-normalized densities = reverse Bregman wrt Z:

 $\mathbf{D}_{\mathsf{KL}}^{+}[\mathbf{q}_{\theta 1}(\mathbf{x}):\mathbf{q}_{\theta 2}(\mathbf{x})] = \mathbf{B}_{\mathsf{Z}}^{*}[\theta_{1}:\theta_{2}] = \mathbf{B}_{\mathsf{Z}}[\theta_{2}:\theta_{1}]$ 

2312 12849

#### Duo Bregman divergences: Generalize BDs with <u>a pair of generators</u>



- Recover Bregman divergence when  $F_1(\theta) = F_2(\theta) = F(\theta)$  $B_F(\theta_1; \theta_2) = F(\theta_1) - F(\theta_2) - \langle \theta_1 - \theta_2, \nabla F(\theta_2) \rangle$
- Only **pseudo-divergence** because  $B_{F1,F2}(\theta'':\theta'')$  positive, not zero

## KLD between nested exponential families amount to duo Bregman pseudo-divergences $\begin{array}{c} q(x \mid \theta) & p(x \mid \theta) \\ p(x \mid \theta) & X_1 \\ \hline q(x \mid \theta) & X_2 \end{array}$

- Consider an exponential family on support X<sub>1</sub>:  $D_{KL}[p(x):q(x)] = \int p(x) \log (p(x)/q(x)) d\mu(x)$   $p(x \mid \theta) = \exp(\langle x, \theta \rangle - F_1(\theta)) d\mu(x)$ with cumulant function  $F_1(\theta) = \log \int_{X_1} \exp(\langle x, \theta \rangle) d\mu(x)$
- Another exponential family with **nested supports:**  $X_1 \subseteq X_2$  $q(x \mid \theta) = exp(\langle x, \theta \rangle - F_2(\theta)) d \mu(x)$

is an exponential family with  $F_2(\theta) = \log \int_{X_2} \exp(\langle x, \theta \rangle) d\mu(x) \ge F_1(\theta)$ 

• Then KLD amounts to a reverse duo Bregman pseudo-divergence:  $D_{KL}[p(x | \theta_1) : q(x | \theta_2)] = B_{F2,F1}^{rev}(\theta_1; \theta_2) = B_{F2,F1}(\theta_2; \theta_1)$ 

"Statistical divergences between densities of truncated exponential families with nested supports: Duo Bregman and duo Jensen divergences." *Entropy* 24.3 (2022)

### Curved Bregman divergences

 $F^*(\eta)$ 

11 -

Consider a domain U which maps to a subset of  $\Theta$  by  $\theta = c(u)$  with dim(U)<dim( $\Theta$ ):

 $B_{F,u}(u_1 : u_2) := B_F(C(u_1):C(u_2))$  is not Bregman when  $\{c(u) \mid u \in U\}$  not convex usually not a Bregman divergence unless c(.) is affine Example: Symmetrized Bregman divergences (Jeffreys-Bregman div.) are curved Bregman divergences:  $S_F(\theta_1, \theta_2) = \langle \theta_1 - \theta_2, \eta_1 - \eta_2 \rangle$ 

$$\begin{split} S_F(\theta_1:\theta_2) &= B_F(\theta_1:\theta_2) + B_F(\theta_2:\theta_1), \\ &= B_F(\theta_1:\theta_2) + B_{F^*}(\nabla F(\theta_1):\nabla F(\theta_2)) \\ &= \breve{B}_{F_{\xi}}(\xi(\theta_1):\xi(\theta_2)), \\ &= \langle \theta, \eta \rangle - F(\theta) \qquad F_{\xi}(\theta, \eta) := F(\theta) + F^*(\eta) \qquad \xi(\theta) = (\theta, \nabla F(\theta)) \\ \{(\theta, \nabla F(\theta)): \theta \in \Theta\} \qquad \text{m-dimensional submanifold in 2m-dimensional space} \end{split}$$

# Curved Bregman centroid is the Bregman projection of the full Bregman centroid

#### Theorem:

$$\arg\min_{u\in\mathcal{U}}\sum_{i=1}^{n}w_{i}B_{F}(\theta_{i}:\theta(u)) = \arg\min_{u\in\mathcal{U}}B_{F}(\bar{\theta}:\theta(u)) \qquad [Bregman \text{ projection}]$$
$$\theta_{i} = \theta(u_{i}) \qquad \bar{\theta} = \sum_{i}w_{i}\theta_{i}$$

#### Proof.

$$\begin{split} \min_{u \in \mathcal{U}} \sum_{i=1}^{n} w_{i} B_{F}(\theta_{i}:\theta(u)) &= \sum_{i=1}^{n} w_{i} (F(\theta_{i}) - F(\theta(u)) - \langle \theta_{i} - \theta(u), \nabla F(\theta(u)) \rangle), \\ &\equiv -F(\theta(u)) - \langle \bar{\theta} - \theta(u), \nabla F(\theta(u)) \rangle, \\ &\equiv F(\bar{\theta}) - F(\theta(u)) - \langle \bar{\theta} - \theta(u), \nabla F(\theta(u)) \rangle \\ &= B_{F}(\bar{\theta}:\theta(u)). \end{split}$$
  
"What is... an information projection?" Notices of the AMS 65.3 (2018): 321-324.



Right-sided Bregman ball: $\sigma_F(\theta, r) = \{\theta' \in \Theta : B_F(\theta':\theta) \leq r\}$ Left-sided Bregman ball: $\sigma_F^{\star}(\theta, r) = \{\theta' \in \Theta : B_F(\theta:\theta') \leq r\}$ 

Application: Boolean algebra of unions & intersections of Bregman balls

## Right Bregman ball and its complement

 $\mathcal{F} := \{ (\theta, y \ge F(\theta)) : \theta \in \Theta \subset \mathbb{R}^m \} \subset \mathbb{R}^{m+1}$ 



 $_{\mathbb{R}}$   $\downarrow$  means vertical projection S<sup>c</sup>: complement of set S

To any sphere, associate an hyperplane:

 $H_{\theta,r}: y = \langle \theta' - \theta, \nabla F(\theta) \rangle + F(\theta) + r$ 

Reciprocally, to an hyperplane cutting the function graph, associate a sphere

$$z = \langle \mathbf{x}, \mathbf{a} \rangle + b$$
  
Center:  $\mathbf{c} = \nabla^{-1} F(\mathbf{a})$   
Radius:  $\langle \mathbf{a}, \mathbf{c} \rangle - F(\mathbf{c}) + b$ 

 $\sigma^{c} = \mathbb{X} \setminus \sigma = \downarrow (H_{\sigma}^{-} \cap \partial \mathcal{F}) \qquad \sigma = \downarrow (H_{\sigma}^{+} \cap \partial \mathcal{F}) \qquad \sigma^{c} = \mathbb{X} \setminus \sigma = \downarrow (H_{\sigma}^{-} \cap \partial \mathcal{F})$ 

Lifting to potential Bregman generator graph

#### Intersection of two right Bregman balls



#### Union of two right Bregman balls



#### Example: Euclidean spheres <sup>™</sup> potential function: Paraboloid, L22



Top view displays the union of disks



Bregman manifolds: Geometry of convex conjugates

## Dual Hessian geometry

[Koszul'64, Shima'70's, Amari&Nagaoka'80's]

On geodesic triangles with right angles in a dually flat space, Progress in Information Geometry: Theory and Applications, Springer 2021 Dual geometry of Bregman manifolds: Convex conjugates (F, F\*) yield dual flat connections  $(M,F \rightarrow g(\theta) = \nabla^2 F(\theta), F \rightarrow \nabla, F^* \rightarrow \nabla^*)$  • A connection  $\nabla$  is flat





 $\mathsf{D}(\mathsf{P}_1,\mathsf{P}_2) = \mathsf{B}_{\mathsf{F}}(\boldsymbol{\theta}_1:\boldsymbol{\theta}_2) = \mathsf{Y}_{\mathsf{F},\,\mathsf{F}^*}(\,\boldsymbol{\theta}_{1,\,\boldsymbol{\eta}_2}) = \mathsf{Y}_{\mathsf{F}^*,\mathsf{F}}(\,\boldsymbol{\eta}_{2,\,\boldsymbol{\theta}_1}) = \mathsf{B}_{\mathsf{F}^*}(\,\boldsymbol{\eta}_{2,\,\boldsymbol{\eta}_1})$ 

 A connection ∇ is flat if there exists a coordinate system θ such that all Christoffel symbols vanish: Γ (θ) =0.

- θ is called ∇ –affine coordinate system
- \[
   \begin{aligned}
   -geodesic solves as line segments
   \]



"The many faces of information geometry." Not. Am. Math. Soc 69.1 (2022): 36-45.

### Dual geometry of smooth Legendre-type functions



#### Example: Bregman manifold of multivariate Gaussians **Cumulant function is convex:** $(M,g, \nabla, \nabla^*)$ $F_{\theta}(\theta) = \frac{1}{2} \left( d \log \pi - \log |\theta_M| + \frac{1}{2} \theta_v^{\top} \theta_M^{-1} \theta_v \right)$ $\mu_{\alpha}^{e} = \Sigma_{\alpha}^{e} \left( (1-\alpha) \Sigma_{1}^{-1} \mu_{1} + \alpha \Sigma_{2}^{-1} \mu_{2} \right)$ $\Sigma_{\alpha}^{e} = \left( (1-\alpha)\Sigma_{1}^{-1} + \alpha\Sigma_{2}^{-1} \right)^{-1}$ with respect to natural parameters: $\gamma_{p_{\mu_{1},\sigma_{1}},p_{\mu_{2},\Sigma_{2}}}^{e}(\alpha) =: p_{\mu_{\alpha}^{e},\Sigma_{\alpha}^{e}} = p_{(1-\alpha)\theta_{1}+\alpha\theta_{2}} \qquad \theta = (\Sigma^{-1}\mu, \frac{1}{2}\Sigma^{-1}) \qquad \theta = (\theta_{v},\theta_{M}) = \left(\Sigma^{-1}\mu, \frac{1}{2}\Sigma^{-1}\right) \qquad \theta = (\theta$ $\theta = (\theta_v, \theta_M) = \left(\Sigma^{-1}\mu, \frac{1}{2}\Sigma^{-1}\right)$ e-geodesic $\dot{\nabla} = \frac{\nabla^e + \nabla^m}{2}$ **m-geodesic** beware not mixture of Gaussians! $\nabla^e$ $\gamma^m_{p_{\mu_1,\sigma_1},p_{\mu_2,\Sigma_2}}(\alpha) =: p_{\mu^m_\alpha,\Sigma^m_\alpha} = p_{(1-\alpha)\eta_1+\alpha\eta_2}$ $\eta = (\mu, -\Sigma - \mu \mu^{\mathsf{T}})$ $p_{\mu_1,\Sigma_1}$

$$\mu_{\alpha}^{m} = (1-\alpha)\mu_{1} + \alpha\mu_{2} =: \bar{\mu}_{\alpha}$$
  
$$\Sigma_{\alpha}^{m} = (1-\alpha)\Sigma_{1} + \alpha\Sigma_{2} + (1-\alpha)\mu_{1}\mu_{1}^{\top} + \alpha\mu_{2}\mu_{2}^{\top} - \bar{\mu}_{\alpha}\bar{\mu}_{\alpha}^{\top}$$

#### **Bregman divergence = reverse Kullback-Leibler divergence**

$$\frac{1}{2} \left( \operatorname{tr}(\Sigma_2^{-1} \Sigma_1) - \log \frac{\operatorname{det}(\Sigma_2)}{\operatorname{det}(\Sigma_1)} - d + (\mu_2 - \mu_1)^\top \Sigma_2^{-1} (\mu_2 - \mu_1) \right) \xrightarrow{25}$$

#### Curved exponential families: Submanifolds



**Theorem** (Curved Bregman centroid/barycenter) Let  $\theta_i = \theta(u_i)$ 's be n weighted parameters of  $\mathcal{U}$  with weight vector  $w \in \Delta_{n-1}$  (the (n-1)-dimensional standard simplex). Then the barycenter in  $\mathcal{U}$  with respect to the curved Bregman divergence amounts to the Bregman projection of the center of mass  $\overline{\theta} = \sum_i w_i \theta_i$  (right Bregman barycenter) onto  $\mathcal{U}$ :

$$\arg\min_{u\in\mathcal{U}}\sum_{i=1}^{n}w_{i}B_{F}(\theta_{i}:\theta(u)) = \arg\min_{u\in\mathcal{U}}B_{F}(\bar{\theta}:\theta(u))$$

$$\overset{\theta(u_{2}) \text{ submanifold}}{\bar{\theta}}$$

$$\overset{\theta(u_{2}) \text{ submanifold}}{\overline{\theta}}$$

U



Note: submanifold topology can be non-trivial

Bijection between regular exponential families and regular Bregman divergences:

$$\log p_F(x;\theta) = -B_{F^*}(t(x):\eta) + F^*(t(x))$$

Curved BD centroid ↔ MLE of curved exp. fam.

k-MLE: A fast algorithm for learning statistical mixture models, IEEE ICASSP 2012

#### Scaled skewed Jensen divergences & Bregman divergences

 $\forall \alpha \in (0,1), \quad J_{F,\alpha}(\theta_1:\theta_2) := (1-\alpha)F(\theta_1) + \alpha F(\theta_2) - F((1-\alpha)\theta_1 + \alpha\theta_2)$ 



#### The Burbea-Rao and Bhattacharyya centroids, IEEE Transactions on Information Theory 57.8 (2011)

## Example 1 of Bregman manifolds:

## Mixture family manifolds (F=-S is Shannon negentropy)

#### Jensen-Shannon centroid for mixture families

Jensen-Shannon divergence Bounded symmetrization of KLD  $JS(p,q) := \frac{1}{2} \left( KL\left(p : \frac{p+q}{2}\right) + KL\left(q : \frac{p+q}{2}\right) \right)$ 

• Jensen-Shannon divergence between two mixtures amounts to a Jensen divergence:  $JS(p_1, p_2) = J_F(\theta_1, \theta_2)$  for  $p_1 = m_{\theta_1}$  and  $p_2 = m_{\theta_2}$ , where

$$J_F(\theta_1:\theta_2) = \frac{F(\theta_1) + F(\theta_2)}{2} - F\left(\frac{\theta_1 + \theta_2}{2}\right).$$

• Task: Given a set of discrete distributions (categorical distributions, normalized histograms), calculate its Jensen-Shannon centroid:

$$\begin{split} \min_{p} \sum_{i} \mathrm{JS}(p_{i}, p), \\ \min_{\theta} \sum_{i} J_{F}(\theta_{i}, \theta), \\ \min_{\theta} \sum_{i} \frac{F(\theta_{i}) + F(\theta)}{2} - F\left(\frac{\theta_{i} + \theta}{2}\right), \\ \equiv \min_{\theta} \frac{1}{2}F(\theta) - \frac{1}{n}\sum_{i} F\left(\frac{\theta_{i} + \theta}{2}\right) := E(\theta). \end{split}$$

Need to minimize a **difference of convex functions** DCA or **ConCave Convex algorithm or DCA**!

> F is Shannon negentropy (convex)



## Example 2 of Bregman manifolds:

**Exponential family manifolds** (F is cumulant function aka log-partition function)



#### Chernoff point & information-geometry

Unique intersection point of the exponential geodesic with the dual mixture bisector



Here 2D probability simplex of the family of categorical distributions with 3 choices

In the beginning of IG…

[Hotelling 1930, Rao 1945]

$$I(\theta) = [I_{ij}(\theta)], \quad I_{ij}(\theta) = \operatorname{Cov}(X_i, X_j) = E_{\theta} \left[ \frac{\partial}{\partial_{\theta_i}} \log p_{\theta}(x) \ \frac{\partial}{\partial_{\theta_j}} \log p_{\theta}(x) \right] = -E_{\theta} \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p_{\theta}(x) \right]$$

## Fisher-Rao manifolds

### Riemannian geometry

![](_page_33_Picture_5.jpeg)

Length element ds

1854

![](_page_33_Picture_7.jpeg)

![](_page_33_Picture_8.jpeg)

![](_page_33_Picture_9.jpeg)

![](_page_33_Picture_10.jpeg)

Photo 1956

Tractability of Fisher-Rao distance: Yet the open case of the multivariate normal family!

$$I_{ij}(\theta) = \left(\frac{\partial\mu}{\partial\theta_i}\right)^{\top} \Sigma^{-1} \frac{\partial\mu}{\partial\theta_j} + \frac{1}{2} \operatorname{tr} \left(\Sigma^{-1} \frac{\partial\mu}{\partial\theta_i} \Sigma^{-1} \frac{\partial\mu}{\partial\theta_j}\right) \qquad \operatorname{ds}_{\mathcal{N}}^2(\mu, \Sigma) = \mathrm{d}\mu^{\top} \Sigma^{-1} \mathrm{d}\mu + \frac{1}{2} \operatorname{tr} \left(\left(\Sigma^{-1} \mathrm{d}\Sigma\right)^2 \right)^2$$
  

$$\operatorname{Geodesic ODE:} \quad \left\{ \begin{array}{cc} \ddot{\mu} - \dot{\Sigma} \Sigma^{-1} \dot{\mu} & = 0, \\ \ddot{\Sigma} + \dot{\mu} \dot{\mu}^{\top} - \dot{\Sigma} \Sigma^{-1} \dot{\Sigma} & = 0. \end{array} \right. \qquad \begin{array}{c} \operatorname{Solve ODE with} \\ \operatorname{initial values (IV) or} \\ \operatorname{boundary values (BV)} \end{array}$$

Fisher longth.

[BV: Kobayashi 2023]

 $\gamma(0), \gamma(1)$ 

Non-constant sectional curvatures which can also be positive! (geodesics are always unique when negative sectional curvatures)

**Bivariate normal** (represented by ellipsoids)

![](_page_34_Picture_4.jpeg)

 $\begin{matrix} [\mathsf{IV: Eriksen 1987}]\\ \gamma(0), \dot{\gamma}(0) \in T_{\gamma(0)} \end{matrix}$ 

### Fisher-Rao geodesics with boundary

![](_page_35_Figure_1.jpeg)

 $\gamma(0),\gamma(1)$ 

$$\begin{cases} \ddot{\mu} - \dot{\Sigma} \Sigma^{-1} \dot{\mu} &= 0, \\ \ddot{\Sigma} + \dot{\mu} \dot{\mu}^{\mathsf{T}} - \dot{\Sigma} \Sigma^{-1} \dot{\Sigma} &= 0. \end{cases}$$

**Red ellipsoids** are the boundary conditions: That is bivariate normal distributions  $(\mu_0, \Sigma_0)$  and  $(\mu_1, \Sigma_1)$ 

#### [BV: Kobayashi 2023]

Technically, MVN Fisher-Rao geodesic: Riemannian submersion of a horizontal geodesic of a Riemannian symmetric space in 2d+1 dimension

![](_page_36_Figure_0.jpeg)

#### A Python library for geometric computing on Bregman Manifolds **pyBregMan** https://franknielsen.github.io/pyBregMan/

![](_page_36_Figure_2.jpeg)

![](_page_36_Figure_3.jpeg)

### Thank you!

#### Some references

- NF and Richard Nock. "The dual Voronoi diagrams with respect to representational Bregman divergences." Sixth International Symposium on Voronoi Diagrams. IEEE, 2009.
- Boissonnat, Jean-Daniel, FN, and Richard Nock. "Bregman Voronoi diagrams." Discrete & Computational Geometry 44 (2010): 281-307.
- NF. "Statistical divergences between densities of truncated exponential families with nested supports: Duo Bregman and duo Jensen divergences." *Entropy* 24.3 (2022)
- NF and Richard Nock. "Generalizing skew Jensen divergences and Bregman divergences with comparative convexity." *IEEE Signal Processing Letters* 24.8 (2017)
- NF. "Curved representational Bregman divergences and their applications." arXiv preprint arXiv:2504.05654