

Arithmer Cloud

AI Systems



速習情報幾何

Arithmer Cloud Div. Kimiaki Kinugasa

2018/10/25

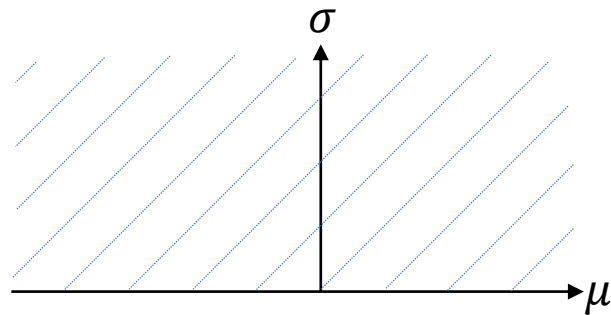
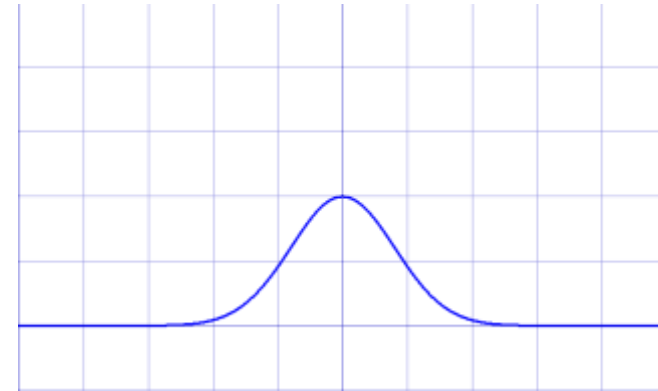
本スライドは、当社と国立大学法人東京大学との間の共同研究に基づき、2018年11月20日、同大学大学院数理科学研究科・理学部数学科において開催された離散数理モデリングセミナーに甘利俊一博士を招請するにあたり、事前の知識習得のために行った社内講義資料を公開用に一部修正したものです。

- **Kimiaki Kinugasa**
 - **Academic Background**
 - The University of Tokyo
 - Bachelor of Engineering
 - **Former Job**
 - Pal Software Service, Inc.
 - Programmer
 - **Current Job**
 - Cloud Architect
 - Cloud Developer
 - Team Manager

- ガウス分布の幾何構造
- 指数型分布族
- 双対平坦空間

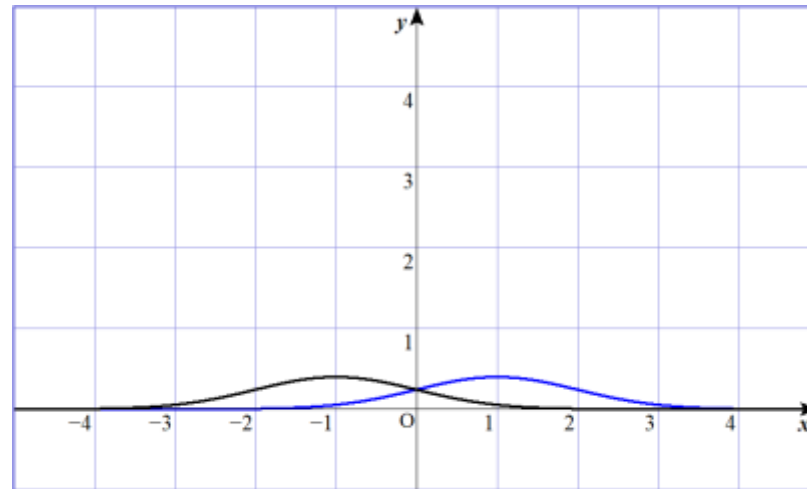
- **ガウス分布の幾何構造**
- 指数型分布族
- 双対平坦空間

$$\mu \in \mathbb{R}, \sigma > 0, x \in \mathbb{R}$$
$$p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\}$$

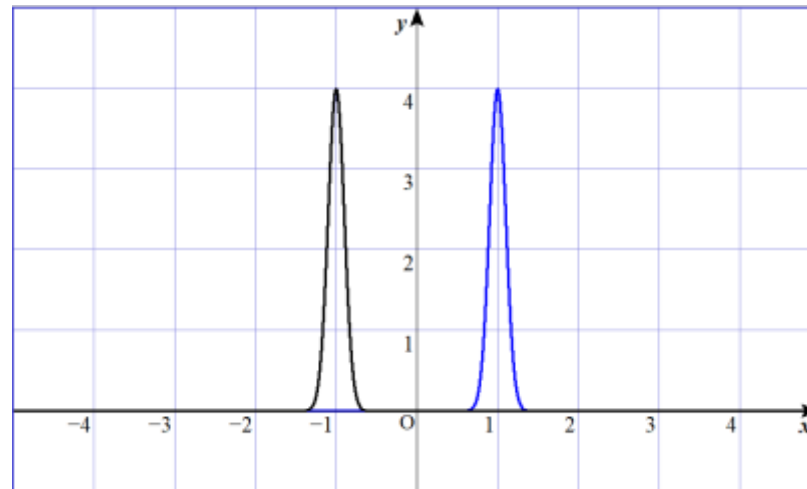


μ, σ 空間の幾何構造について考える

$(-1,1), (1,1)$ の距離は2

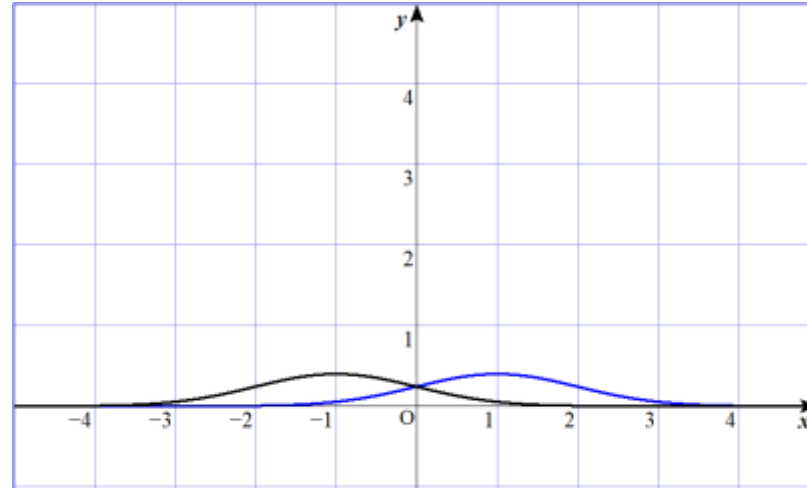


$(-1,0.1), (1,0.1)$ の距離も2

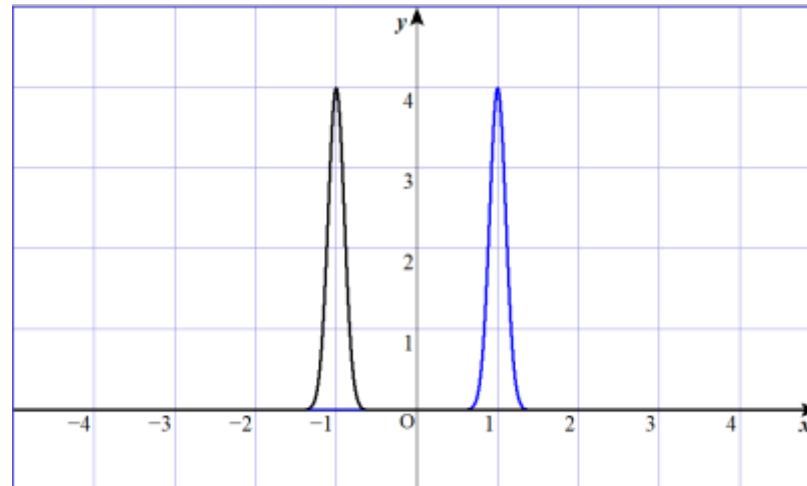


変な気がする . . .

$(-1,1), (1,1)$ の距離は2



$(-1,0.1), (1,0.1)$ の距離も2



変な気がする・・・

μ, σ 空間が確率分布族の空間であることを無視している

μ, σ 空間の微小な2点間の距離は

$$ds^2 = (d\mu \quad d\sigma) \begin{pmatrix} g_{\mu\mu} & g_{\mu\sigma} \\ g_{\mu\sigma} & g_{\sigma\sigma} \end{pmatrix} \begin{pmatrix} d\mu \\ d\sigma \end{pmatrix}$$

と表せる

μ, σ 空間をユークリッド空間とみなすと

$$ds^2 = (d\mu \quad d\sigma) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} d\mu \\ d\sigma \end{pmatrix}$$

となる

リーマン計量 g として何を選ぶべきか?

μ, σ 空間の微小な2点間の距離は

$$ds^2 = (d\mu \quad d\sigma) \begin{pmatrix} g_{\mu\mu} & g_{\mu\sigma} \\ g_{\mu\sigma} & g_{\sigma\sigma} \end{pmatrix} \begin{pmatrix} d\mu \\ d\sigma \end{pmatrix}$$

と表せる

μ, σ 空間をユークリッド空間とみなすと

$$ds^2 = (d\mu \quad d\sigma) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} d\mu \\ d\sigma \end{pmatrix}$$

となる

リーマン計量 g として何を選ぶべきか?

結論

Fisher情報行列を使うべき

定義:

確率密度関数 $p(x; \theta)$ に対して、そのFisher情報行列を

$$E \left[\frac{\partial \log p}{\partial \theta} \left(\frac{\partial \log p}{\partial \theta} \right)^T \right]$$

と定義する。

ガウス分布の場合、Fisher情報行列は

$$\begin{pmatrix} g_{\mu\mu} & g_{\mu\sigma} \\ g_{\mu\sigma} & g_{\sigma\sigma} \end{pmatrix} = \begin{pmatrix} 1/\sigma^2 & 0 \\ 0 & 2/\sigma^2 \end{pmatrix}$$

となる。

定義:

確率密度関数 $p(x; \theta)$ に対して、そのFisher情報行列を

$$E \left[\frac{\partial \log p}{\partial \theta} \left(\frac{\partial \log p}{\partial \theta} \right)^T \right]$$

と定義する。

ガウス分布の場合、Fisher情報行列は

$$\begin{pmatrix} g_{\mu\mu} & g_{\mu\sigma} \\ g_{\mu\sigma} & g_{\sigma\sigma} \end{pmatrix} = \begin{pmatrix} 1/\sigma^2 & 0 \\ 0 & 2/\sigma^2 \end{pmatrix}$$

となる。

平面(ユークリッド空間)じゃない！

ガウス曲率: 空間の各点における曲がり具合を表現する量

2次元曲面の微小距離 ds^2 が1次微分形式 θ^1, θ^2 を用いて

$$ds^2 = \theta^1\theta^1 + \theta^2\theta^2$$

で与えられたるとき、

第1構造式

$$\begin{aligned}d\theta^1 &= \theta^2 \wedge \omega_2^1 \\d\theta^2 &= \theta^1 \wedge \omega_1^2\end{aligned}$$

但し、 $\omega_2^1 = -\omega_1^2$

第2構造式

$$d\omega_2^1 = K\theta^1 \wedge \theta^2$$

但し、 K はガウス曲率

$$ds^2 = (d\mu \quad d\sigma) \begin{pmatrix} g_{\mu\mu} & g_{\mu\sigma} \\ g_{\mu\sigma} & g_{\sigma\sigma} \end{pmatrix} \begin{pmatrix} d\mu \\ d\sigma \end{pmatrix}$$
$$\begin{pmatrix} g_{\mu\mu} & g_{\mu\sigma} \\ g_{\mu\sigma} & g_{\sigma\sigma} \end{pmatrix} = \begin{pmatrix} 1/\sigma^2 & 0 \\ 0 & 2/\sigma^2 \end{pmatrix}$$

より、

$$ds^2 = \frac{1}{\sigma^2} d\mu^2 + \frac{2}{\sigma^2} d\sigma^2$$

従って、

$$\theta^1 = \frac{1}{\sigma} d\mu, \theta^2 = \frac{\sqrt{2}}{\sigma} d\sigma$$

と表せる。

$$ds^2 = (d\mu \quad d\sigma) \begin{pmatrix} g_{\mu\mu} & g_{\mu\sigma} \\ g_{\mu\sigma} & g_{\sigma\sigma} \end{pmatrix} \begin{pmatrix} d\mu \\ d\sigma \end{pmatrix}$$
$$\begin{pmatrix} g_{\mu\mu} & g_{\mu\sigma} \\ g_{\mu\sigma} & g_{\sigma\sigma} \end{pmatrix} = \begin{pmatrix} 1/\sigma^2 & 0 \\ 0 & 2/\sigma^2 \end{pmatrix}$$

より、

$$ds^2 = \frac{1}{\sigma^2} d\mu^2 + \frac{2}{\sigma^2} d\sigma^2$$

従って、

$$\theta^1 = \frac{1}{\sigma} d\mu, \theta^2 = \frac{\sqrt{2}}{\sigma} d\sigma$$

と表せる。第1構造式、第2構造式から計算すると、

$$\omega_2^1 = -\frac{1}{\sqrt{2}} \theta^1$$
$$K = -\frac{1}{2}$$

$$ds^2 = (d\mu \quad d\sigma) \begin{pmatrix} g_{\mu\mu} & g_{\mu\sigma} \\ g_{\mu\sigma} & g_{\sigma\sigma} \end{pmatrix} \begin{pmatrix} d\mu \\ d\sigma \end{pmatrix}$$

$$\begin{pmatrix} g_{\mu\mu} & g_{\mu\sigma} \\ g_{\mu\sigma} & g_{\sigma\sigma} \end{pmatrix} = \begin{pmatrix} 1/\sigma^2 & 0 \\ 0 & 2/\sigma^2 \end{pmatrix}$$

より、

$$ds^2 = \frac{1}{\sigma^2} d\mu^2 + \frac{2}{\sigma^2} d\sigma^2$$

従って、

$$\theta^1 = \frac{1}{\sigma} d\mu, \theta^2 = \frac{\sqrt{2}}{\sigma} d\sigma$$

と表せる。第1構造式、第2構造式から計算すると、

$$\omega_2^1 = -\frac{1}{\sqrt{2}} \theta^1$$

$$K = -\frac{1}{2}$$

ガウス分布は負の定曲率空間になる！美しい！

ガウス分布のリーマン計量は定まった

$$\begin{pmatrix} g_{\mu\mu} & g_{\mu\sigma} \\ g_{\mu\sigma} & g_{\sigma\sigma} \end{pmatrix} = \begin{pmatrix} 1/\sigma^2 & 0 \\ 0 & 2/\sigma^2 \end{pmatrix}$$

これを使えば、2点間の最短距離は以下のように定義できる。

$$d((\mu_1, \sigma_1), (\mu_2, \sigma_2)) = \min_{(\mu(t), \sigma(t)) \in C} \int_0^1 \sqrt{\frac{1}{\sigma^2(t)} \left(\frac{d\mu}{dt}(t) \right)^2 + \frac{2}{\sigma^2(t)} \left(\frac{d\sigma}{dt}(t) \right)^2} dt$$

但し $C = \{(\mu(0), \sigma(0)) = (\mu_1, \sigma_1), (\mu(1), \sigma(1)) = (\mu_2, \sigma_2)\}$ を満たす曲線の集合

最短距離は求めるのが大変。もっと便利なものは?

定義:

確率密度関数 p, q のKLダイバージェンスを

$$D(p||q) = \int p \log \frac{p}{q} dx$$

と定義する。

- 二つの確率分布間の分離度を測るもの
- 最尤推定に使われたりする
- $D(p||q) \geq 0$
- $D(p||q) = 0 \Leftrightarrow p = q$
- $D(p||q) \neq D(q||p)$

$$\begin{aligned} D(p(\mu_1, \sigma_1) \| p(\mu_2, \sigma_2)) &= \int p(\mu_1, \sigma_1) \log \frac{p(\mu_1, \sigma_1)}{p(\mu_2, \sigma_2)} dx \\ &= \log \frac{\sigma_2}{\sigma_1} - \frac{1}{2} + \frac{\sigma_1^2 + \mu_1^2 - 2\mu_1\mu_2 + \mu_2^2}{2\sigma_2} \end{aligned}$$

$$\begin{aligned}
 D(p(\mu_1, \sigma_1) \| p(\mu_2, \sigma_2)) &= \int p(\mu_1, \sigma_1) \log \frac{p(\mu_1, \sigma_1)}{p(\mu_2, \sigma_2)} dx \\
 &= \log \frac{\sigma_2}{\sigma_1} - \frac{1}{2} + \frac{\sigma_1^2 + \mu_1^2 - 2\mu_1\mu_2 + \mu_2^2}{2\sigma_2}
 \end{aligned}$$

KLダイバージェンスは微分するとリーマン計量が出てくる

$$\begin{aligned}
 g_{\mu\mu} &= \frac{\partial^2 D(p(\mu_1, \sigma_1) \| p(\mu_2, \sigma_2))}{\partial \mu_1 \partial \mu_1} \bigg|_{\mu_1 = \mu_2 = \mu, \sigma_1 = \sigma_2 = \sigma} = \frac{1}{\sigma^2} \\
 g_{\mu\sigma} &= \frac{\partial^2 D(p(\mu_1, \sigma_1) \| p(\mu_2, \sigma_2))}{\partial \mu_1 \partial \sigma_1} \bigg|_{\mu_1 = \mu_2 = \mu, \sigma_1 = \sigma_2 = \sigma} = 0 \\
 g_{\sigma\sigma} &= \frac{\partial^2 D(p(\mu_1, \sigma_1) \| p(\mu_2, \sigma_2))}{\partial \sigma_1 \partial \sigma_1} \bigg|_{\mu_1 = \mu_2 = \mu, \sigma_1 = \sigma_2 = \sigma} = \frac{2}{\sigma^2}
 \end{aligned}$$

$$\begin{aligned}
 D(p(\mu_1, \sigma_1) \| p(\mu_2, \sigma_2)) &= \int p(\mu_1, \sigma_1) \log \frac{p(\mu_1, \sigma_1)}{p(\mu_2, \sigma_2)} dx \\
 &= \log \frac{\sigma_2}{\sigma_1} - \frac{1}{2} + \frac{\sigma_1^2 + \mu_1^2 - 2\mu_1\mu_2 + \mu_2^2}{2\sigma_2}
 \end{aligned}$$

KLダイバージェンスは微分するとリーマン計量が出てくる

$$\begin{aligned}
 g_{\mu\mu} &= \frac{\partial^2 D(p(\mu_1, \sigma_1) \| p(\mu_2, \sigma_2))}{\partial \mu_1 \partial \mu_1} \bigg|_{\mu_1=\mu_2=\mu, \sigma_1=\sigma_2=\sigma} = \frac{1}{\sigma^2} \\
 g_{\mu\sigma} &= \frac{\partial^2 D(p(\mu_1, \sigma_1) \| p(\mu_2, \sigma_2))}{\partial \mu_1 \partial \sigma_1} \bigg|_{\mu_1=\mu_2=\mu, \sigma_1=\sigma_2=\sigma} = 0 \\
 g_{\sigma\sigma} &= \frac{\partial^2 D(p(\mu_1, \sigma_1) \| p(\mu_2, \sigma_2))}{\partial \sigma_1 \partial \sigma_1} \bigg|_{\mu_1=\mu_2=\mu, \sigma_1=\sigma_2=\sigma} = \frac{2}{\sigma^2}
 \end{aligned}$$

KLダイバージェンスは、リーマン計量の情報を含んでいる

- ガウス分布の幾何構造
- 指数型分布族
- 双対平坦空間

定義:

n 次元モデルの確率密度関数が

$$p(x; \theta) = C(x) \exp \left(\sum_{i=1}^n \theta^i F_i(x) - \psi(\theta) \right) = C(x) \exp \left(\theta^i F_i(x) - \psi(\theta) \right)$$

と表せるとき、このモデルを指数型分布族という。

ガウス分布は2次元の指数型分布族

$$C(x) = 1, F_1(x) = x, F_2(x) = x^2, \theta^1 = \frac{\mu}{\sigma^2}, \theta^2 = -\frac{1}{2\sigma^2}$$
$$\psi(\theta) = -\frac{(\theta^1)^2}{4\theta^2} + \frac{1}{2} \log \left(-\frac{\pi}{\theta^2} \right)$$

指数型分布族の期待値

$$E[F_i(x)] = \frac{\partial \psi}{\partial \theta^i}(\theta)$$

指数型分布族のFisher情報行列は

$$E \left[\frac{\partial \log p}{\partial \theta} \left(\frac{\partial \log p}{\partial \theta} \right)^T \right] = \frac{\partial^2 \psi}{\partial \theta^2}(\theta)$$

Fisher情報行列は定義より非負定値行列
従って、 $\psi(\theta)$ は凸関数

指数型分布族の期待値

$$E[F_i(x)] = \frac{\partial \psi}{\partial \theta^i}(\theta)$$

指数型分布族のFisher情報行列は

$$E \left[\frac{\partial \log p}{\partial \theta} \left(\frac{\partial \log p}{\partial \theta} \right)^T \right] = \frac{\partial^2 \psi}{\partial \theta^2}(\theta)$$

Fisher情報行列は定義より非負定値行列
従って、 $\psi(\theta)$ は凸関数

今後、Fisher情報行列を計量として使いたいのので、
 $\frac{\partial^2 \psi}{\partial \theta^2}(\theta)$ は常に正定値であると仮定する。

- ガウス分布の幾何構造
- 指数型分布族
- 双対平坦空間

指数型分布族

$$p(x; \theta) = C(x) \exp(\theta^i F_i(x) - \psi(\theta))$$

に対して、 θ を自然パラメータまたは正準パラメータという。

$F_i(x)$ の期待値もパラメータと考えることができる。

$$\eta_i = E[F_i(x)] = \frac{\partial \psi}{\partial \theta^i}(\theta)$$

これを期待値パラメータという。

θ と η は以下の関係性がある。(Legendre変換)

$$\phi(\eta) := \max_{\theta} (\theta^i \eta_i - \psi(\theta))$$

$$\eta_i = \frac{\partial \psi}{\partial \theta^i}(\theta), \theta^i = \frac{\partial \phi}{\partial \eta_i}(\eta)$$

$$\psi(\theta) = \max_{\eta} (\theta^i \eta_i - \phi(\eta))$$

$\phi(\eta)$ も凸関数になる。

$$\begin{aligned} D(p(\theta) \| p(\theta')) &= \psi(\theta') - \psi(\theta) - \mathbb{E}_{\theta}[F_i(x)](\theta'^i - \theta^i) \\ &= \psi(\theta') - \psi(\theta) - \frac{\partial \psi}{\partial \theta^i}(\theta)(\theta'^i - \theta^i) \end{aligned}$$

これは、凸関数 $\psi(\theta)$ に関するBregman ダイバージェンス

$$D_{\psi}(\theta' \| \theta) := \psi(\theta') - \psi(\theta) - \frac{\partial \psi}{\partial \theta^i}(\theta)(\theta'^i - \theta^i)$$

に一致する。

KLダイバージェンスは ϕ のBregman ダイバージェンスでも表せる

$$\begin{aligned}
 D(p(\theta)||p(\theta')) &= D_\psi(\theta'||\theta) \\
 &= \psi(\theta') - \psi(\theta) - \frac{\partial \psi}{\partial \theta^i}(\theta)(\theta'^i - \theta^i) \\
 &= \theta'^i \eta'_i - \phi(\eta') - (\theta^i \eta_i - \phi(\eta)) - \eta_i(\theta'^i - \theta^i) \\
 &= \phi(\eta) - \phi(\eta') - \theta'^i(\eta_i - \eta'_i) \\
 &= \phi(\eta) - \phi(\eta') - \frac{\partial \phi}{\partial \eta'_i}(\eta')(\eta_i - \eta'_i) \\
 &= D_\phi(\eta||\eta')
 \end{aligned}$$

指数型分布族に対して、双対な2つの座標系が得られた。
それぞれの座標系における直線は以下の式を満たす。

$$\begin{array}{ll} \theta \text{座標系の直線 (m-測地線)} \theta(t) & \eta \text{座標系の直線 (e-測地線)} \eta(t) \\ \frac{d^2 \theta^i}{dt^2} = 0 & \frac{d^2 \eta_i}{dt^2} = 0 \end{array}$$

θ 座標系の直線は η 座標系でどう表せるか？

η 座標系の直線は θ 座標系でどう表せるか？

$$\frac{d^2 \theta^i}{dt^2} = 0 \Leftrightarrow \frac{d^2 \eta_i}{dt^2} + \underbrace{\frac{\partial^2 \psi}{\partial \theta^i \partial \theta^l} \frac{\partial^3 \phi}{\partial \eta_j \partial \eta_k \partial \eta_l}}_{\Gamma_i^{mjk}} \frac{d\eta_j}{dt} \frac{d\eta_k}{dt} = 0$$

$$\frac{d^2 \eta_i}{dt^2} = 0 \Leftrightarrow \frac{d^2 \theta^i}{dt^2} + \underbrace{\frac{\partial^2 \phi}{\partial \eta_i \partial \eta_l} \frac{\partial^3 \psi}{\partial \theta^j \partial \theta^k \partial \theta^l}}_{\Gamma_{jk}^{ei}} \frac{d\theta^j}{dt} \frac{d\theta^k}{dt} = 0$$

指数型分布族

$$p(x; \theta) = C(x) \exp(\theta^i F_i(x) - \psi(\theta))$$

のなす空間は次のよう構造をもつ

θ 座標系で表すと,
リーマン計量

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

2種類の直線(測地線)

$$\frac{d^2 \theta^i}{dt^2} = 0, \frac{d^2 \theta^i}{dt^2} + \Gamma_i^{ejk} \frac{d\theta^j}{dt} \frac{d\theta^k}{dt} = 0$$

η 座標系で表すと
リーマン計量

$$g^{ij} = \frac{\partial^2 \phi}{\partial \eta_i \partial \eta_j}$$

2種類の直線(測地線)

$$\frac{d^2 \eta_i}{dt^2} + \Gamma_i^{mjk} \frac{d\eta_j}{dt} \frac{d\eta_k}{dt} = 0, \frac{d^2 \eta_i}{dt^2} = 0$$

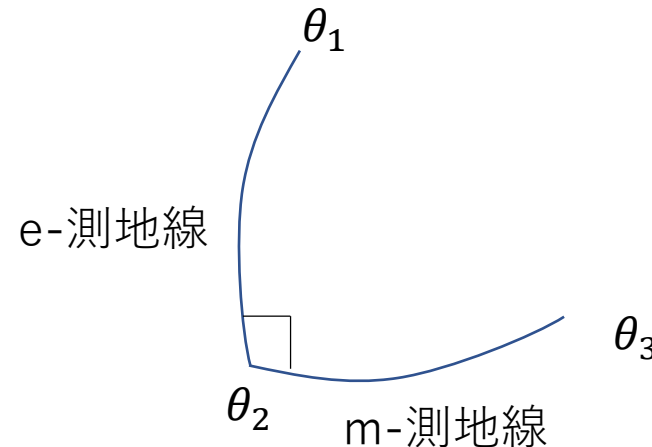
このような空間を双対平坦空間という。
(厳密な定義はしていない)

指数型分布族

$$p(x; \theta) = C(x) \exp(\theta^i F_i(x) - \psi(\theta))$$

の異なる3点 $\theta_1 = (\theta_1^i), \theta_2 = (\theta_2^i), \theta_3 = (\theta_3^i)$ に対して

θ_1, θ_2 を結ぶ e-測地線と θ_2, θ_3 を結ぶ m-測地線は θ_2 で直交する
 $\Leftrightarrow D(\theta_1 || \theta_3) = D(\theta_1 || \theta_2) + D(\theta_2 || \theta_3)$



- Amari, S. Information Geometry and Its Applications (Applied Mathematical Sciences); Springer, 2016.
- 甘利俊一; 長岡浩司; 情報幾何の方法; 岩波書店, 1993.
- 藤原彰夫 情報幾何学の基礎; 牧野書店, 2015.
- Amari, S. Information geometry of the EM and em algorithms for neural networks; Neural Networks 1995, 8, 1379-1408