

筑波大学講義

# 情報幾何の展開

甘利俊一

理化学研究所 名誉研究員

東京大学 名誉教授

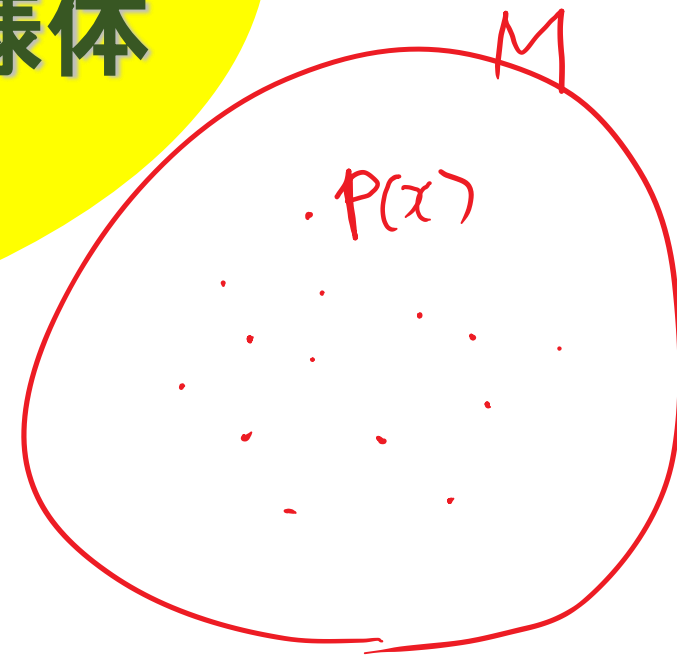
帝京大学特任教授

1. 統計推論の情報幾何: 不変性
2. 双対平坦空間: 凸解析とLegendre変換
3. 多層神経回路の統計神経力学
4. Wasserstein距離の情報幾何

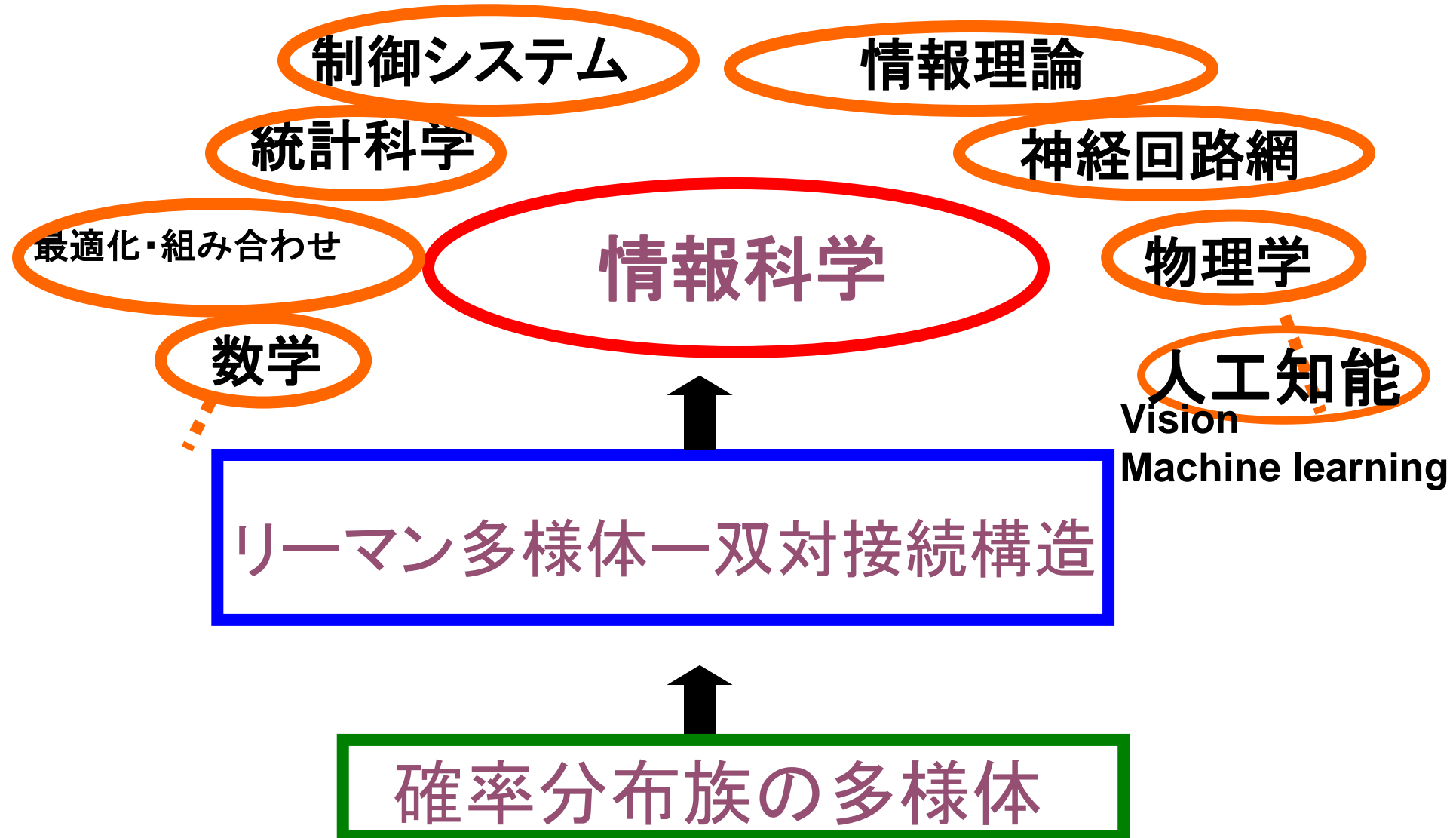
# 情報幾何

-- 確率分布族のなす多様体

$$M = \{p(\mathbf{x})\}$$

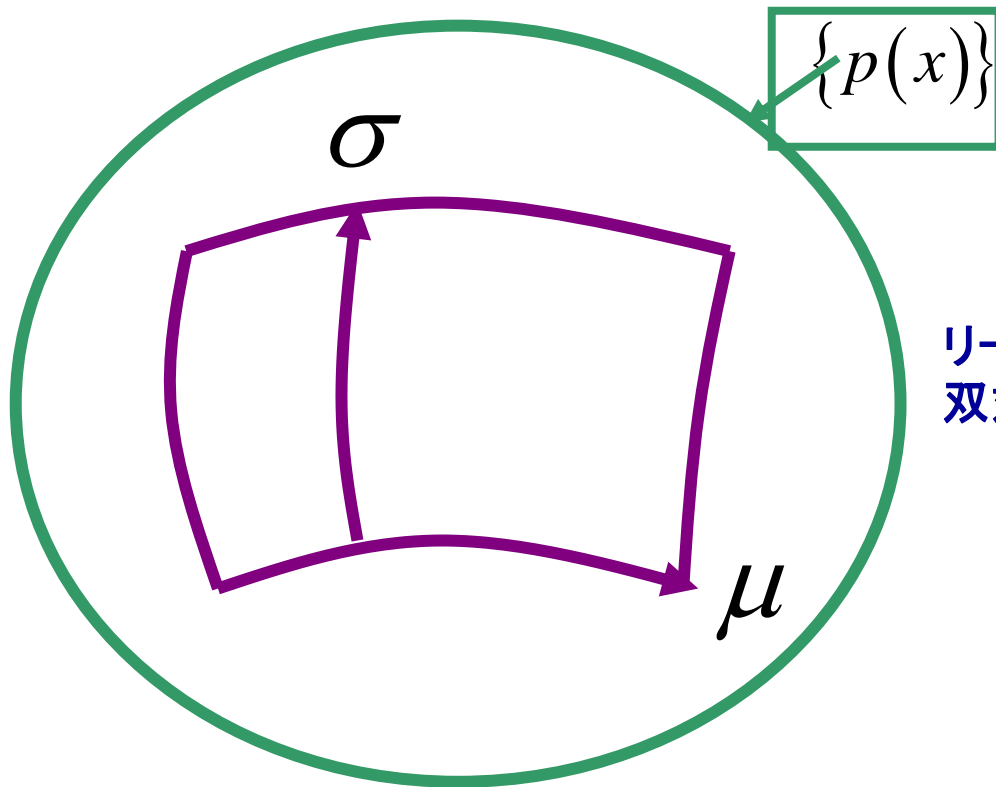


# 情報幾何



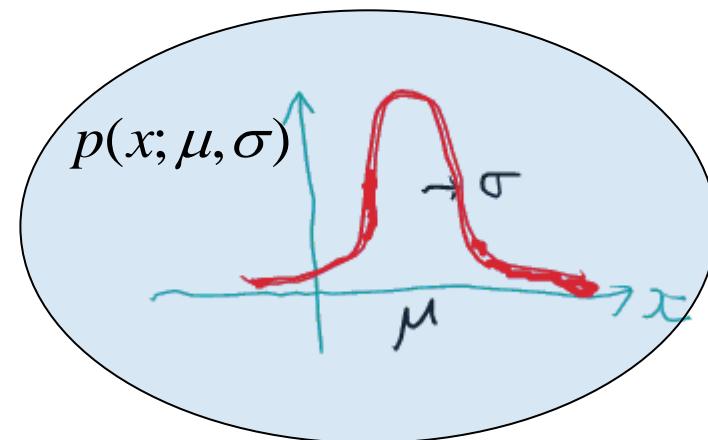
# 情報幾何とは？

$$S = \{p(x; \mu, \sigma)\} \quad p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$



$$S = \{p(x; \theta)\}$$

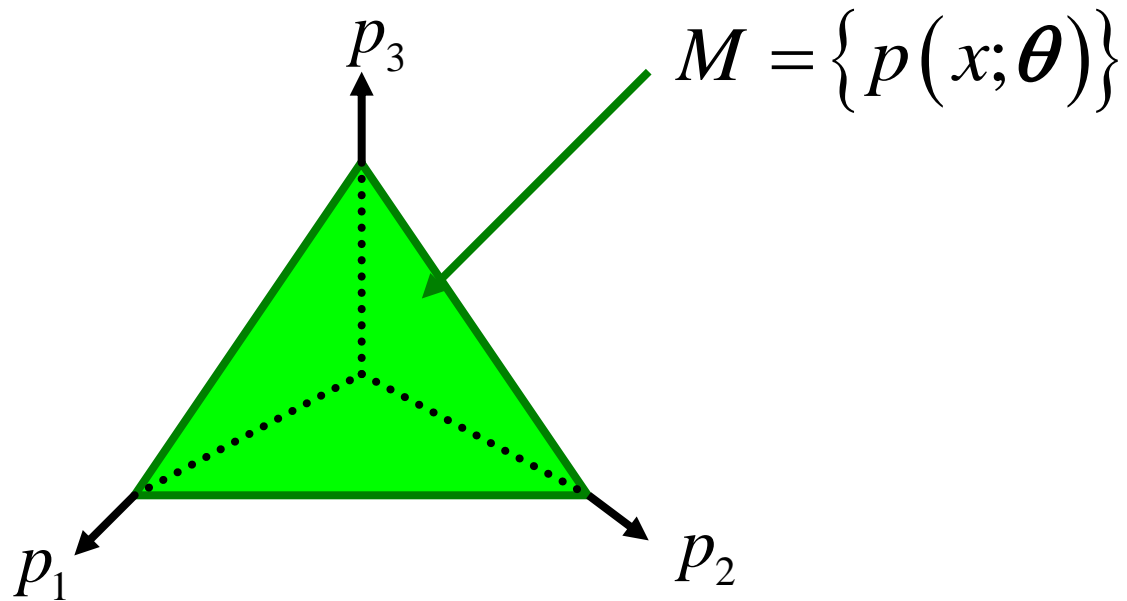
リーマン幾何  
双対アフィン接続



# 離散確率分布(三つ目さいころ)

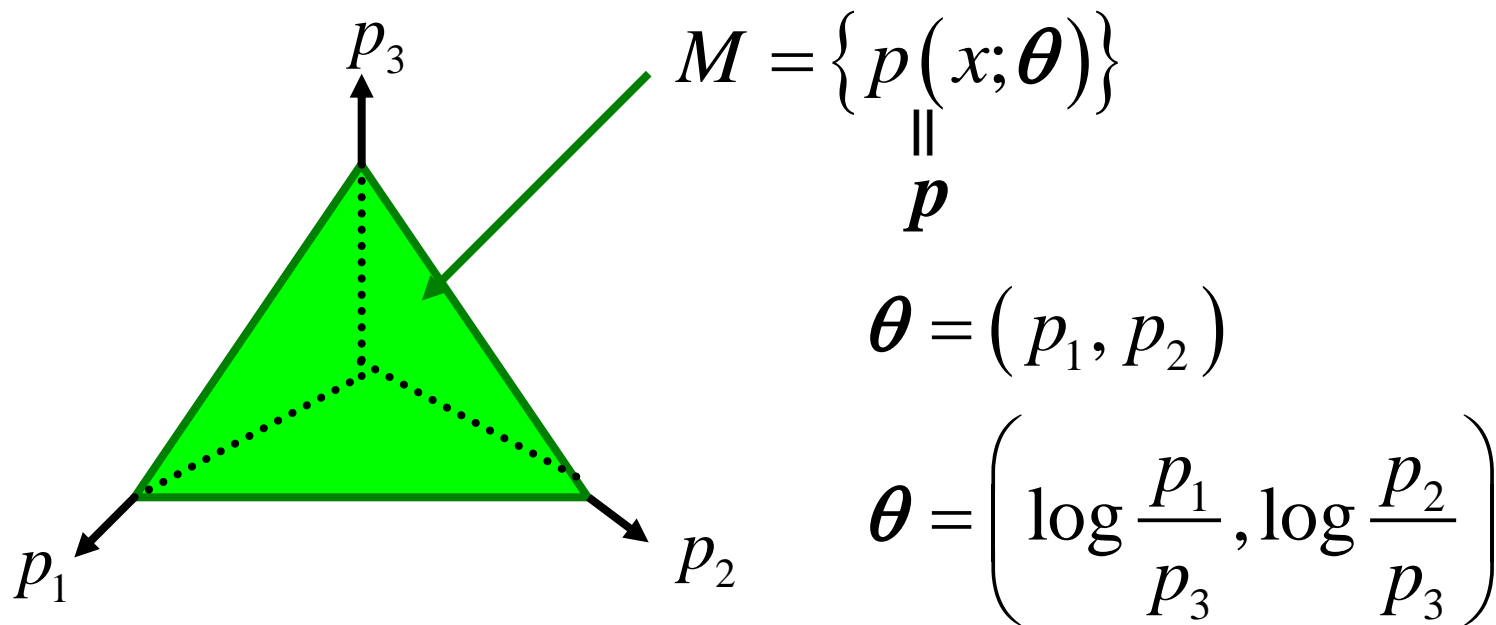
$$x = 1, 2, 3 \quad S_n = \{p(x)\} \quad n = 3$$

$$\mathbf{p} = (p_1, p_2, p_3), \quad p_1 + p_2 + p_3 = 1$$



# 確率分布族のつくる多様体(座標系)

$$\mathbf{p} = (p_1, p_2, p_3) \quad p_1 + p_2 + p_3 = 1$$



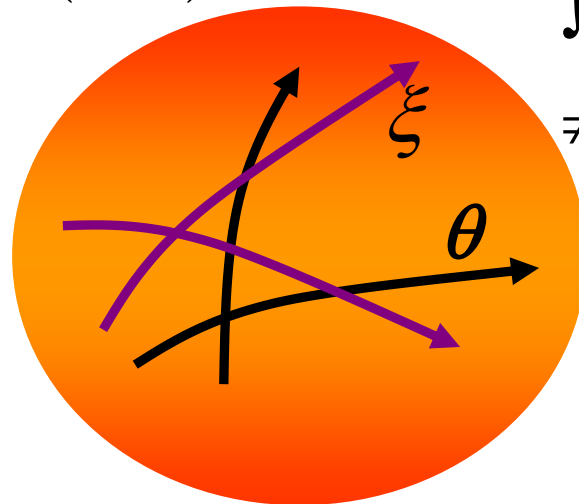
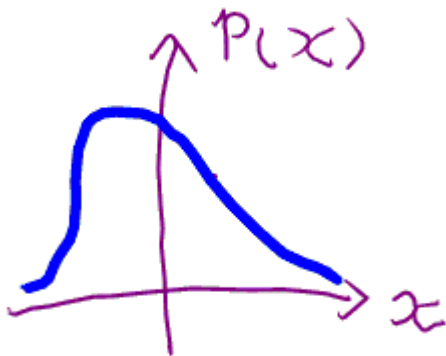
# 不変性の原理: $S = \{p(x, \theta)\}$

## 1. パラメータのとり方によらない

$$\xi = \xi(\theta), \quad \bar{p}(x, \xi) \quad D = \sum \theta_i^2 \neq \sum \xi_i^2$$

## 2. 確率変数の表示スケールによらない

$$y = y(x), \quad \bar{p}(y, \theta) \quad \int |p(x, \theta_1) - p(x, \theta_2)|^2 dx \\ \neq \int |\bar{p}(y, \theta_1) - \bar{p}(y, \theta_2)|^2 dy$$



# 確率分布空間の接空間

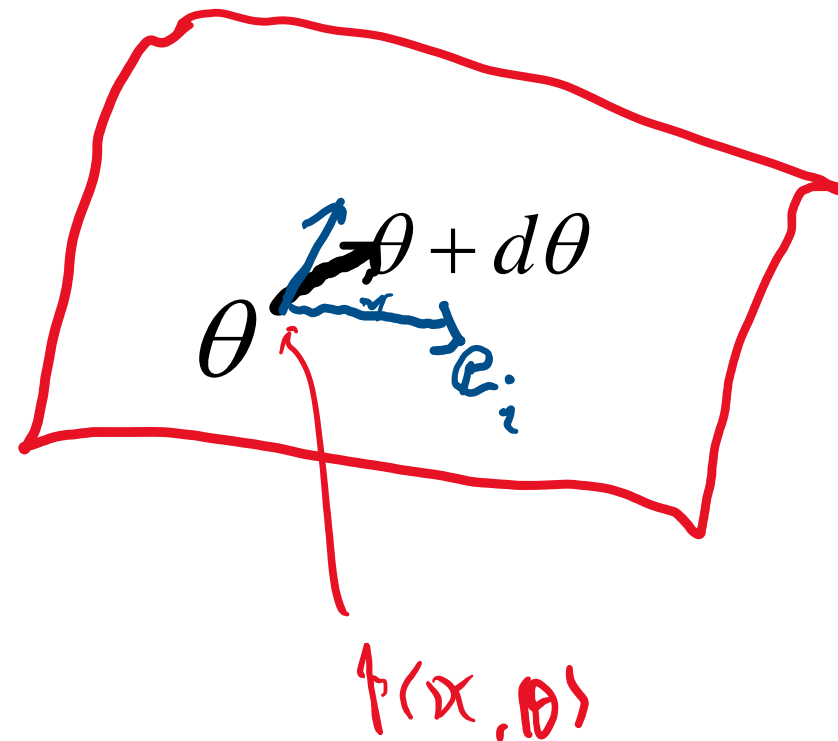
$$\mathcal{S} = \{p(x, \theta)\}$$

Spanned by scores

$$d\boldsymbol{\theta} = \sum d\theta^i \mathbf{e}_i$$

$$g_{ij}(\boldsymbol{\theta}) = \langle \mathbf{e}_i, \mathbf{e}_j \rangle$$

$$\mathbf{e}_i = \frac{\partial}{\partial \theta^i} \approx \frac{\partial}{\partial \theta^i} \log p(x, \theta)$$



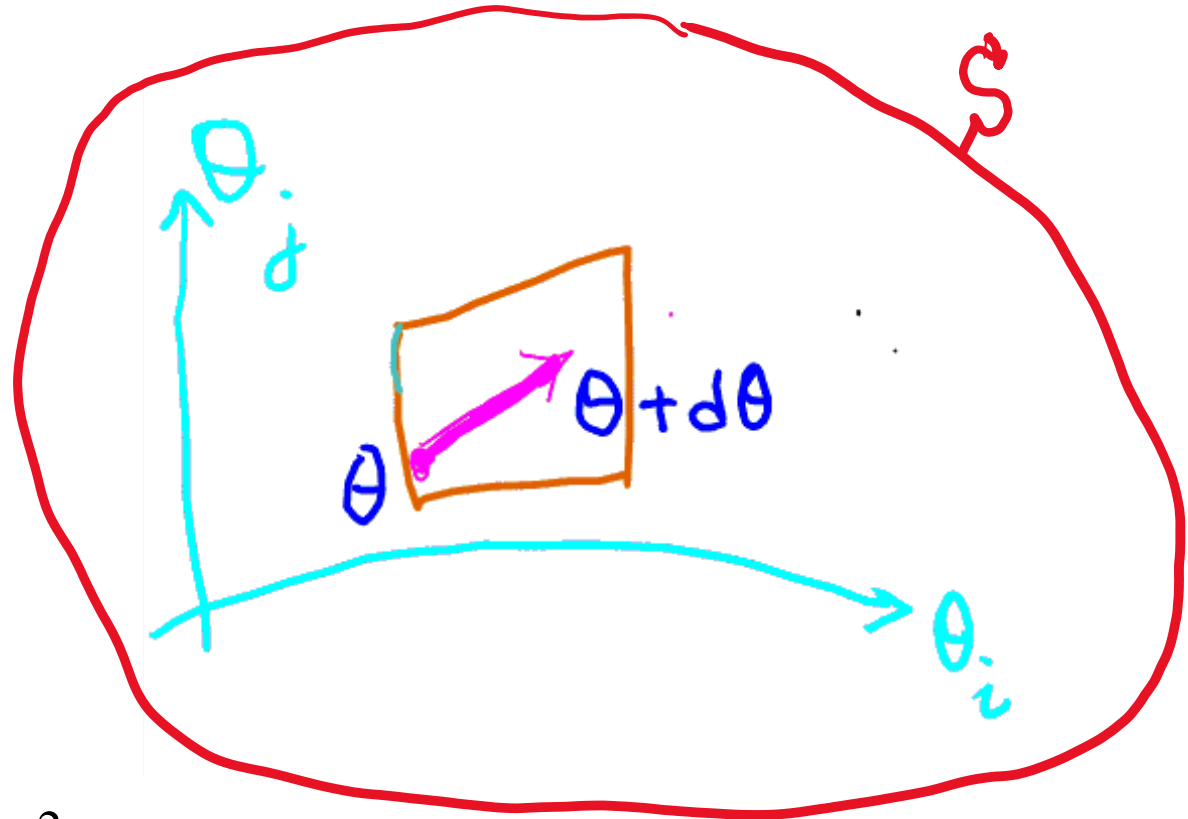


# リーマン構造

$$ds^2 = \langle d\boldsymbol{\theta}, d\boldsymbol{\theta} \rangle = \sum g_{ij}(\boldsymbol{\theta}) d\theta^i d\theta^j$$
$$= d\boldsymbol{\theta}^T G(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

$$G(\boldsymbol{\theta}) = (g_{ij}) = \langle d\mathbf{e}_i, d\mathbf{e}_j \rangle$$

Euclidean  $G = E$   $ds^2 = \sum (d\theta^i)^2$



# リーマン計量とアファイン接続

## 双対接続 $\{M, G, \nabla, \nabla^*\}$

Fisher情報行列

$$g = (g_{ij})$$

共変微分

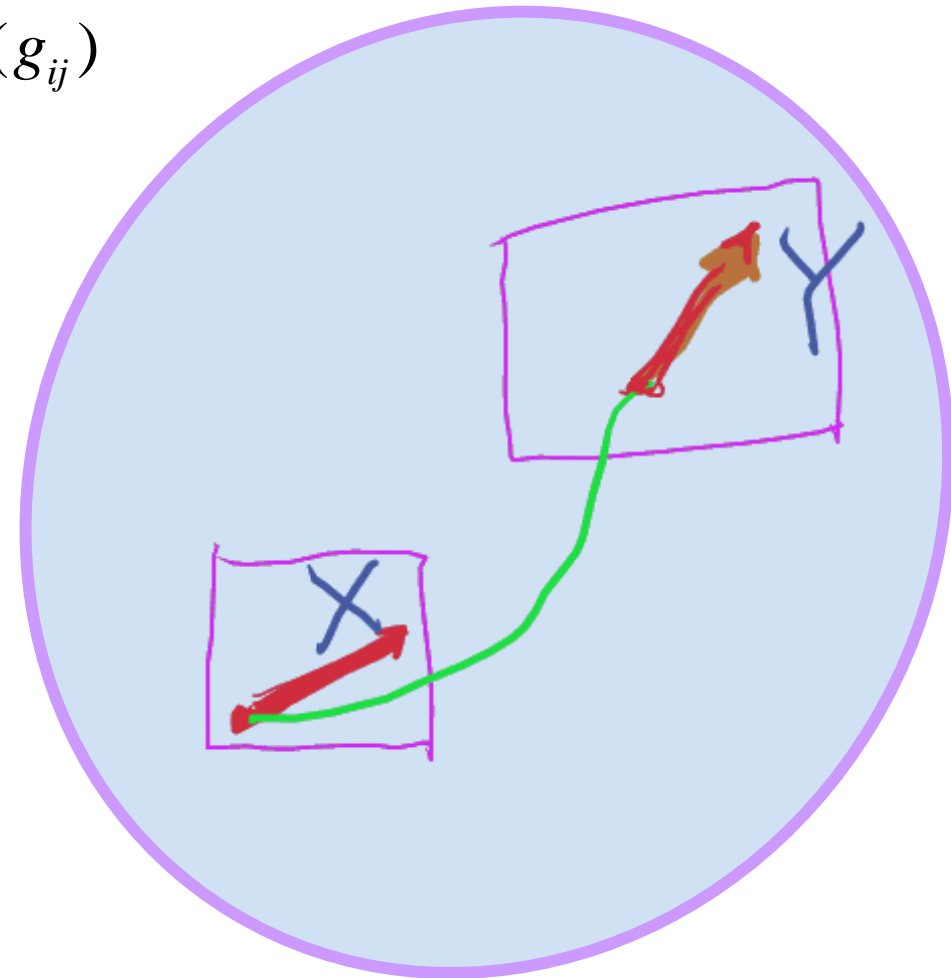
$$\nabla_X Y$$

$$\Pi_c X = Y$$

$$\text{測地線 } \Pi \dot{X} = \dot{X} \quad X = X(t)$$

$$s = \int \sqrt{\sum g_{ij}(\theta) d\theta^i d\theta^j}$$

最短距離: まっすぐ



# 二つのアファイン接続 $(\nabla, \nabla^*)$

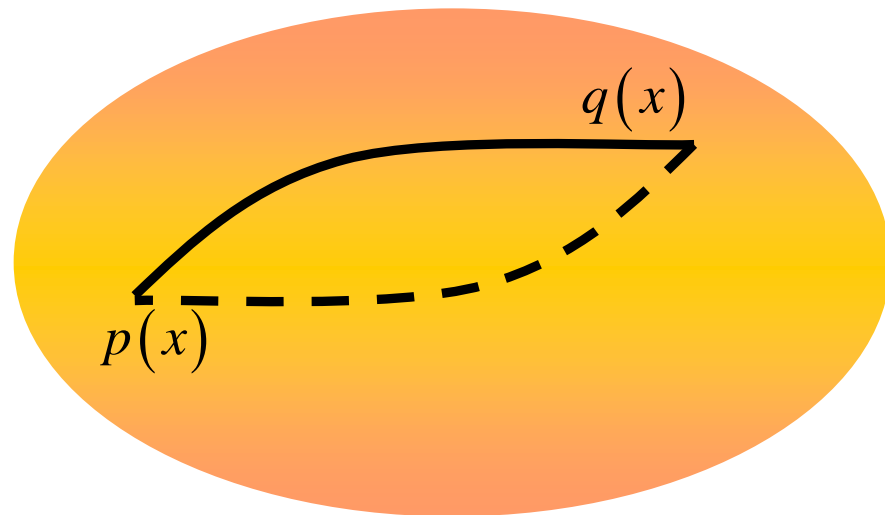
$$(\Pi, \Pi^*)$$

*e-geodesic*

$$\log r(x, t) = t \log p(x) + (1-t) \log q(x) + c(t)$$

*m-geodesic*

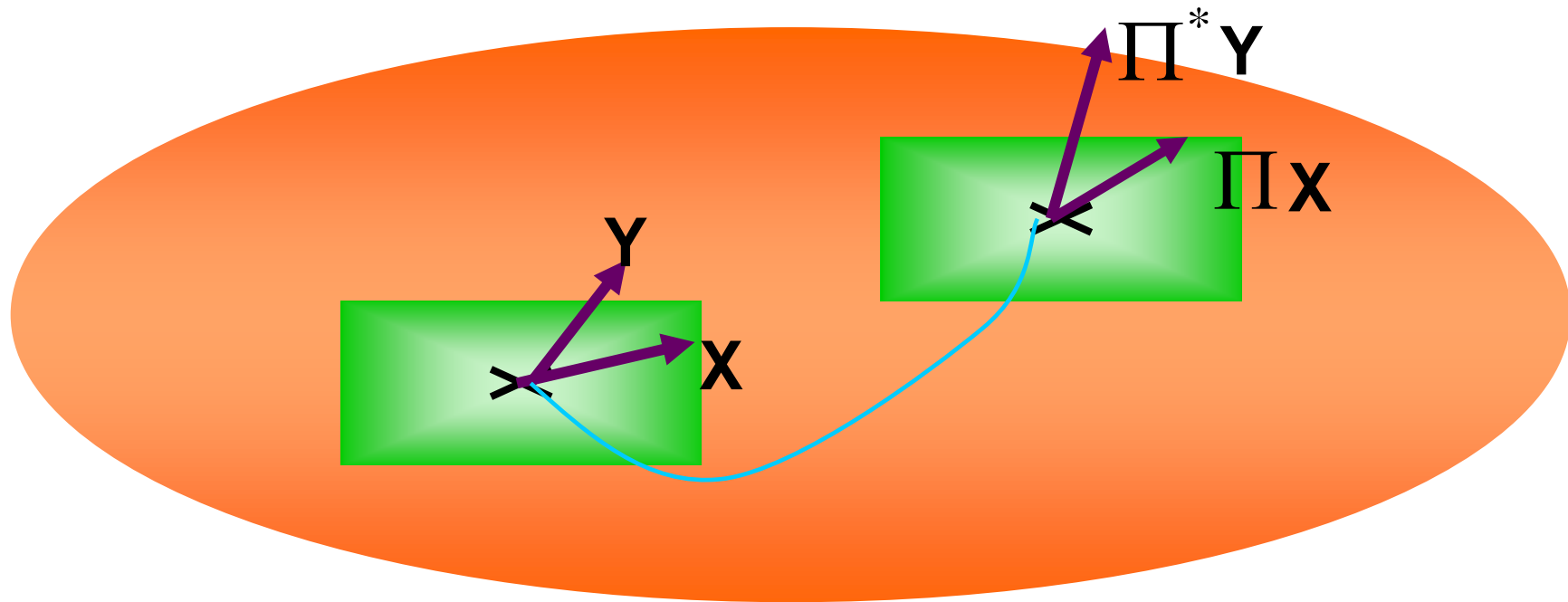
$$r(x, t) = tp(x) + (1-t)q(t)$$



# 双対接続:二つのアフィン接続

$$X\langle Y, Z \rangle = \langle \nabla_X Y, Z \rangle + \langle \nabla_X Z, Y \rangle$$

$$\langle X, Y \rangle = \langle \Pi X, \Pi^* Y \rangle \quad \langle X, Y \rangle = \sum g_{ij} X^i Y^j$$



Riemannian geometry:  $\Pi = \Pi^*$

# 指数型分布族： 双对平坦空間

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

$\psi(\theta)$ : convex function, free-energy

Gaussian:

$$e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$
$$\left[ \begin{array}{l} x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} x \\ x^2 \end{pmatrix}, \quad \theta = \begin{pmatrix} \theta^1 \\ \theta^2 \end{pmatrix} = \begin{pmatrix} -\frac{1}{2\sigma^2} \\ \frac{\mu^2}{\sigma^2} \end{pmatrix} \\ \theta \cdot x = -\frac{(x-\mu)^2}{2\sigma^2} + c \end{array} \right.$$

entropy .  $-\varphi(\eta) = -\int p(x, \theta) \log p(x, \theta) dx$

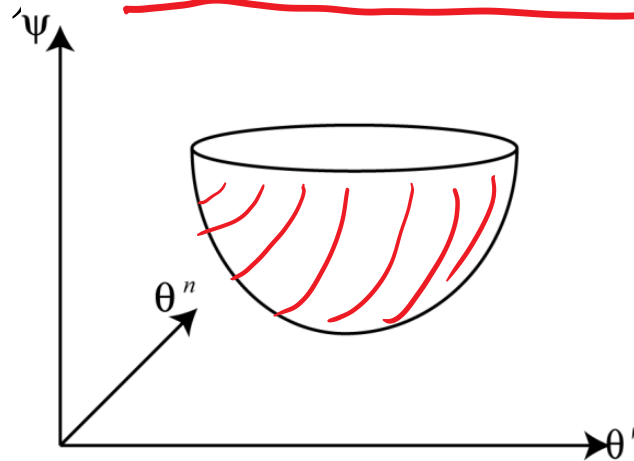
natural parameter :  $\theta = \frac{\partial}{\partial \eta} \varphi(\eta)$

expectation parameter :  $\eta = E[x] = \frac{\partial}{\partial \theta} \psi(\theta)$

# 凸関数、凸解析——双対平坦

$S$  : 座標系  $\theta = (\theta^1, \theta^2, \dots, \theta^n)$

$\psi(\theta)$  : 凸関数 function



$$\psi(\theta) = \frac{1}{2} \sum (\theta^i)^2$$

**negative entropy  
energy**

$$\varphi(p) = \int p(x) \log p(x) dx$$

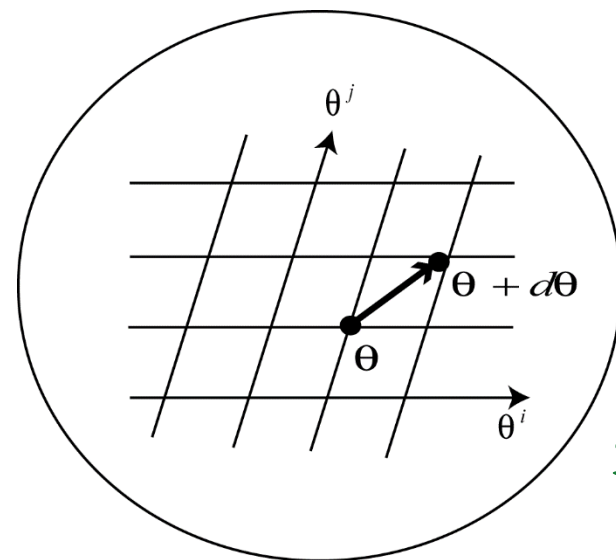
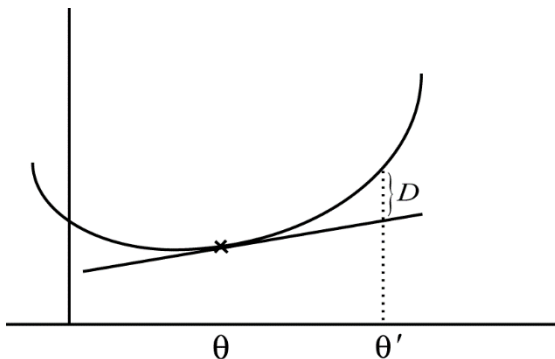
# リーマン計量と平坦性 $\{S, \psi(\theta), \theta\}$

convex:  $\tilde{\theta} = A\theta + c$

## Bregman divergence

$$D(\theta', \theta) = \psi(\theta') - \psi(\theta) - (\theta' - \theta) \cdot \text{grad } \psi(\theta)$$

affine structure



$$D(\theta, \theta + d\theta) = \frac{1}{2} \sum g_{ij}(\theta) d\theta^i d\theta^j$$

$$g_{ij} = \partial_i \partial_j \psi(\theta), \quad \partial_i = \frac{\partial}{\partial \theta^i}$$

straight line

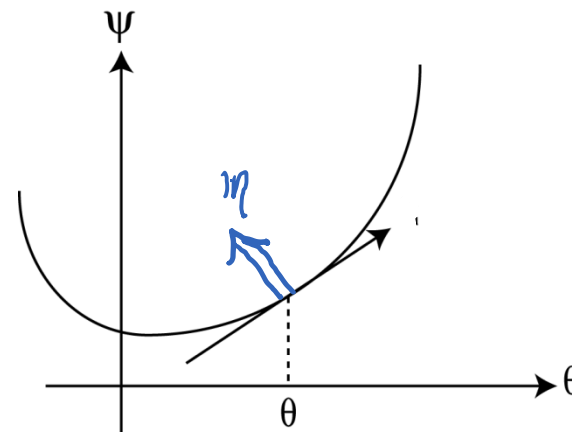
Flatness (affine)  $\theta$  : **geodesic** (not Levi-Civita)  $\leftarrow \theta(t) = t a + b$

# Legendre 变换

dual coordinates  $(\theta, \eta)$

$$\eta_i = \partial_i \psi(\theta), \quad \partial_i = \frac{\partial}{\partial \theta^i}$$

$\psi(\theta)$   $\theta \leftrightarrow \eta$   $\varphi(\eta)$   
one-to-one



$$\theta^i = \partial^i \varphi(\eta), \quad \partial_i = \frac{\partial}{\partial \eta_i}$$

$$\eta_i = \partial_i \psi(\theta), \quad \partial_i = \frac{\partial}{\partial \theta^i}$$

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

$$\varphi(\eta) + \psi(\theta) - \theta_i \eta^i = 0$$

: proof easy

$$D(\theta, \theta') = \psi(\theta) + \varphi(\eta') - \theta \cdot \eta'$$

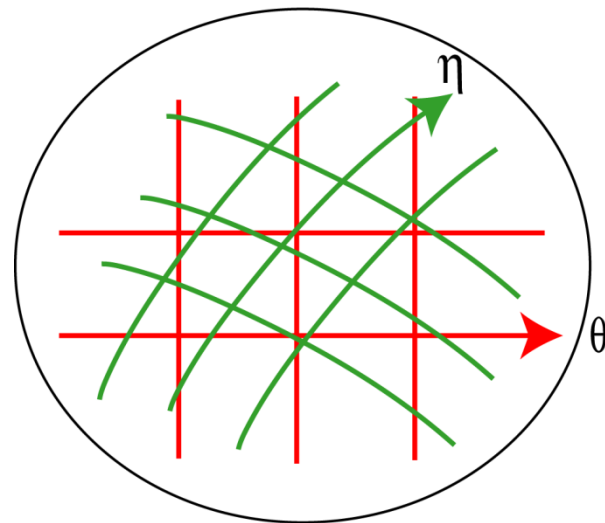


# 双対平坦空間の双対アフィン座標

$$\theta = (\theta_1, \dots, \theta_n)$$

$$\eta = (\eta_1, \dots, \eta_n)$$

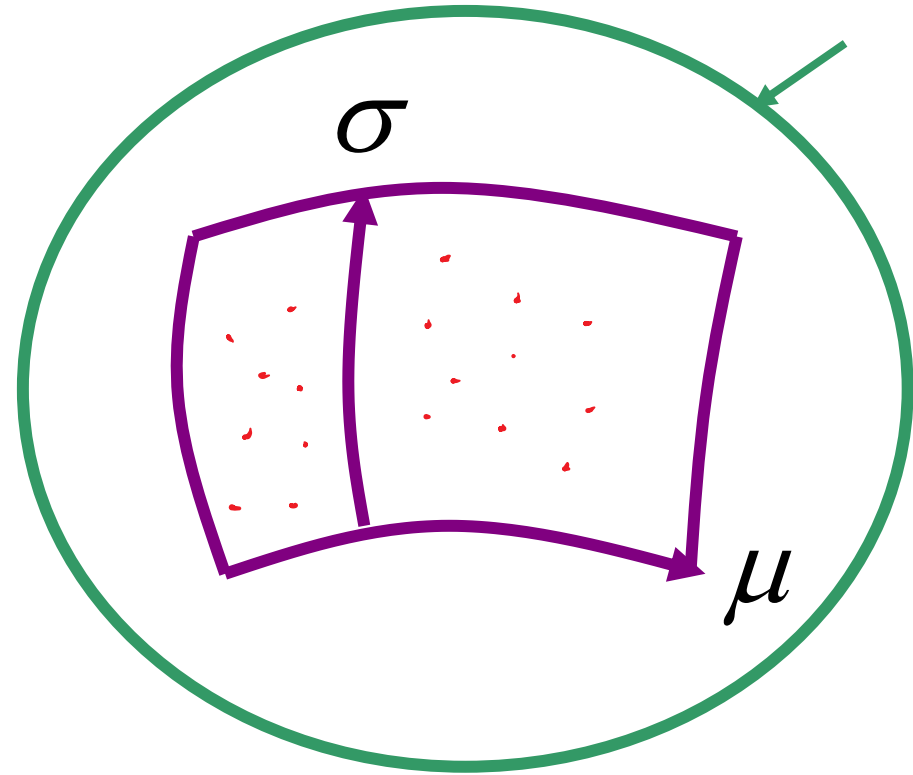
$$\eta = \eta(\theta) \longleftrightarrow \theta = \theta(\eta)$$



one-to-one  
differentiable

## Gaussian distributions

$$\begin{aligned}\mathcal{X} &= (\mu, \sigma^2), \\ \Theta &= \left(-\frac{1}{2\sigma^2}, \frac{\mu^2}{\sigma^2}\right), \\ \eta &= (\mu, \mu^2 + \sigma^2)\end{aligned}$$



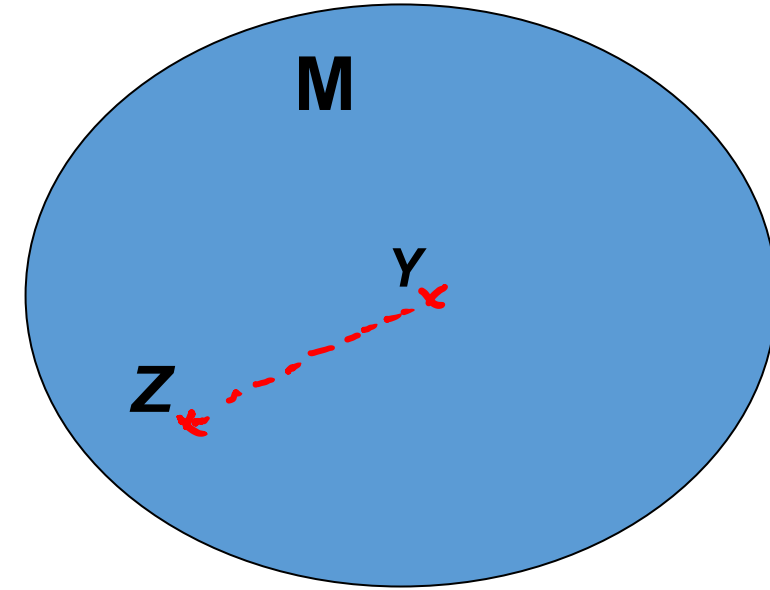
$$S = \{p(x; \mu, \sigma)\} \quad p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

# Divergence: $D[z : y]$

$$D[z : y] \geq 0$$

$$D[z : y] = 0, \quad \text{iff } z = y$$

$$D[z : z + dz] = \frac{1}{2} \sum g_{ij} dz_i dz_j$$



Not necessarily symmetric

$$D[z : y] = D[y : z]$$

positive-definite  $G = (g_{ij})$

Taylor expansion

$$D(z : z + dz) = \frac{1}{2} \sum g_{ij} dz_i dz_j + \frac{1}{6} \sum \kappa_{ijk} dz_i dz_j dz_k + \dots$$

# 双对平坦空間

$\theta$ -coordinates  $\leftrightarrow$   $\eta$ -coordinates

potential functions  $\psi(\theta), \varphi(\eta)$

$$g_{ij}(\theta) = \frac{\partial^2}{\partial \theta_i \partial \theta_j} \psi(\theta) \cdots g^{ij} = \frac{\partial^2}{\partial \eta_i \partial \eta_j} \varphi(\eta)$$

$$\psi(\theta) + \varphi(\eta) - \sum \theta_i \eta_i = 0$$

**exponential family:**  $p(x, \theta) = \exp\{\sum \theta_i x_i - \psi(\theta)\}$

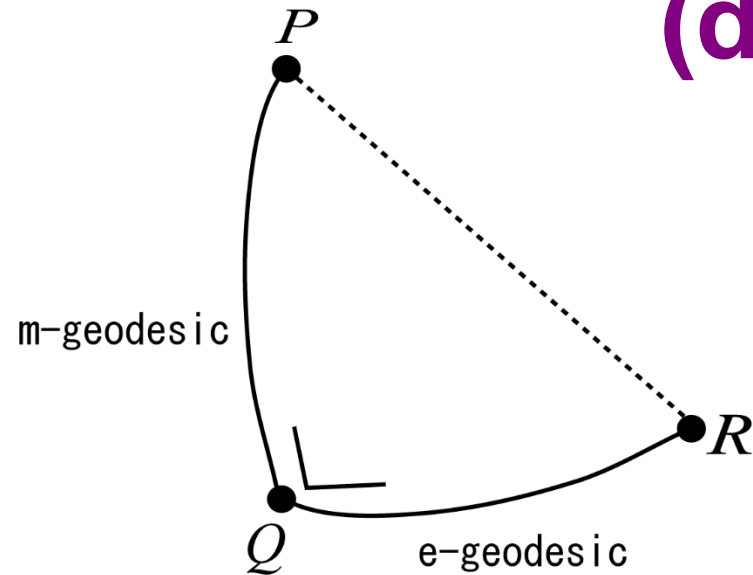
$\psi$ : cumulant generating function

$\varphi$ : negative entropy

canonical divergence  $D(P: P') = \psi(\theta) + \varphi(\eta') - \sum \theta_i \eta'_i$

# 拡張ピタゴラスの定理

(dually flat manifold)



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

proof

ユークリッド空間: 自己双対

$$\theta = \eta$$

$$\psi(\theta) = \frac{1}{2} \sum (\theta_i)^2$$

# 射影定理

$$\min_{Q \in M} D[P:Q]$$

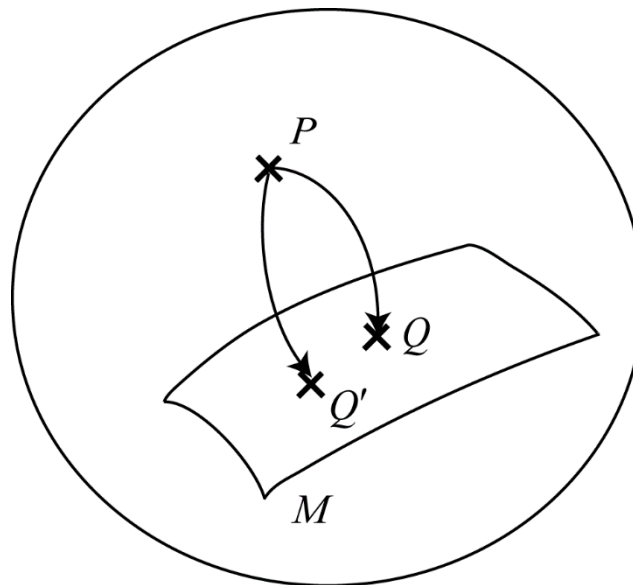
**Q = m-射影 P から M**

**unique when M is e-flat**

$$\min_{Q \in M} D[Q:P]$$

**Q' = e-射影 P から M**

**unique when M is m-flat**



## 双对平坦幾何

Convex function – Bregman divergence – exponential family

– Dually flat Riemannian divergence

$$\psi(\theta) \Rightarrow \mathcal{D}_\psi[\theta : \theta'] \Rightarrow \{\theta, \eta\}, \quad \mathcal{G} = \nabla \nabla \psi$$

Dually flat R-manifold – convex function – canonical divergence

KL-divergence

$$\{\theta, \eta\} \Rightarrow \psi(\theta) \Rightarrow \mathcal{D}_\psi[\theta : \theta']$$
$$\mathcal{G} = \frac{\partial^2 \psi}{\partial \theta^2} \quad \cdot \quad \mathcal{D}_{KL}[p(x) : q(x)]$$

# 一般にダイバージェンスは計量と接続を与える リーマン計量 (Eguchi)

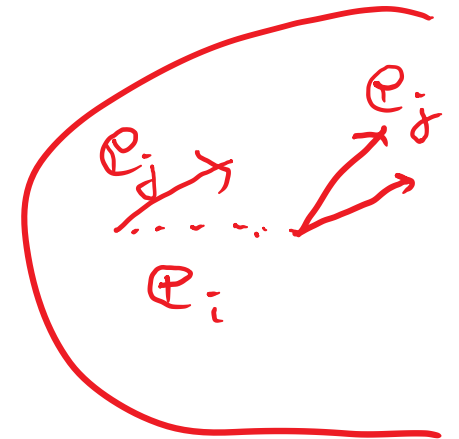
$$g_{ij}(\mathbf{z}) = \partial_i \partial_j D[\mathbf{z} : \mathbf{y}]_{|y=z} : D[\mathbf{z} : \mathbf{y}] = \frac{1}{2} g_{ij}(\mathbf{z}) (z_i - y_i)(z_j - y_j)$$

## 二つのアファイン接続 $\{\nabla, \nabla^*\}$

$$\nabla_{\mathbf{e}_i} \mathbf{e}_j = \Gamma_{ij}^k \mathbf{e}_k$$

$$\Gamma_{ijk}(\mathbf{z}) = -\partial_i \partial_j \partial'_k D[\mathbf{z} : \mathbf{y}]_{|y=z} \quad \partial_i = \frac{\partial}{\partial z_i}, \quad \partial'_i = \frac{\partial}{\partial y_i}$$

$$\Gamma_{ijk}^*(\mathbf{z}) = -\partial'_i \partial'_j \partial_k D[\mathbf{z} : \mathbf{y}]_{|y=z} \quad T = \Gamma^* - \Gamma$$





# 不変でないダイバージェンス

## Wasserstein距離

### q-ダイバージェンス $\leftrightarrow$ $\alpha$ -ダイバージェンス

$$D_{\alpha}[p : q] = \frac{4}{1 - \alpha^2} \sum (1 - p_i^{\frac{1-\alpha}{2}} q_i^{\frac{1+\alpha}{2}})$$

$$D_q[p : q] = \frac{4}{1 - q} \sum (1 - p_i^q q_i^{1-q}); \quad q = 2\alpha - 1$$

projectively-dually flat

# divergence

( $n > 1$ )

$S = \{\mathbf{p}\}$  : space of probability distributions

**invariance**

**dually flat space**

**invariant divergence**

**Flat divergence**

**F-divergence**  
**Fisher inf metric**  
**Alpha connection**

**KL-divergence**

**convex functions**  
**Bregman**

$$D[\mathbf{p} : \mathbf{q}] = \int \mathbf{p}(\mathbf{x}) \log\left\{\frac{\mathbf{p}(\mathbf{x})}{\mathbf{q}(\mathbf{x})}\right\} d\mathbf{x}$$



$q$  – exponential family

$$\log_q(u) = \frac{1}{1-q} (u^{1-q} - 1)$$

$$\log\{p(x, \theta)\} = \theta \cdot \mathbf{x} - \psi_q(\theta)$$

$\psi_q(\theta)$ : convex function

**双対平坦空間(非不変)  
Pythagorasの定理**

$$D_q^*(p : q) = \frac{1}{h_q(\theta)} D_q^*(p : q) : \quad h_q = \sum p_i^q$$

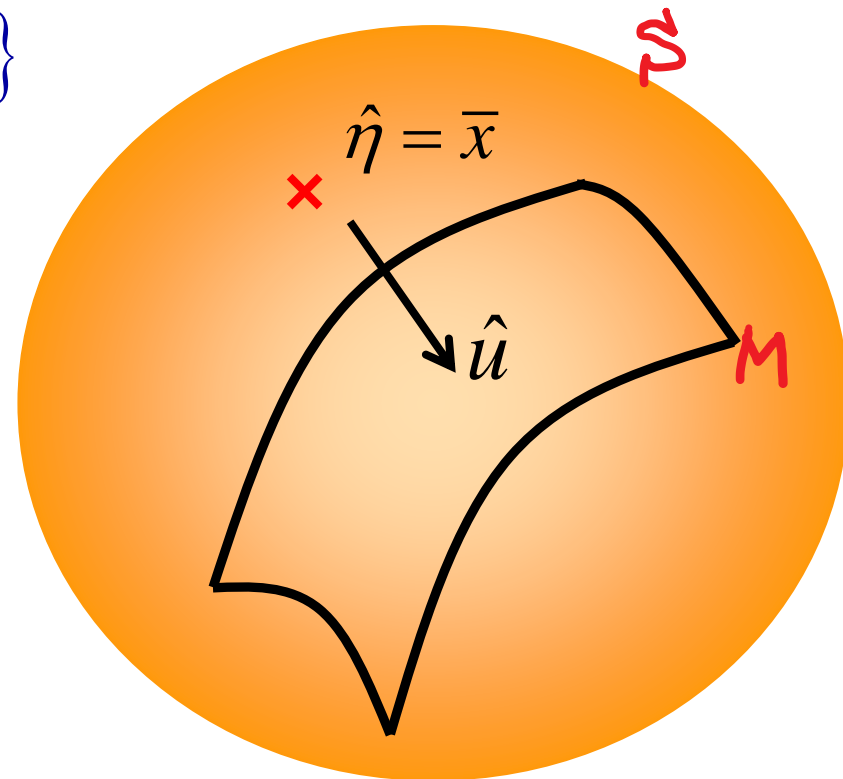
# 統計学への応用： 曲指数型分布族：

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

$$p(x, u) = p(x, \theta(u)) \square x_1, x_2, \dots, x_n$$

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

$\hat{u}(x_1, \dots, x_n)$  : estimator



# 統計学への応用

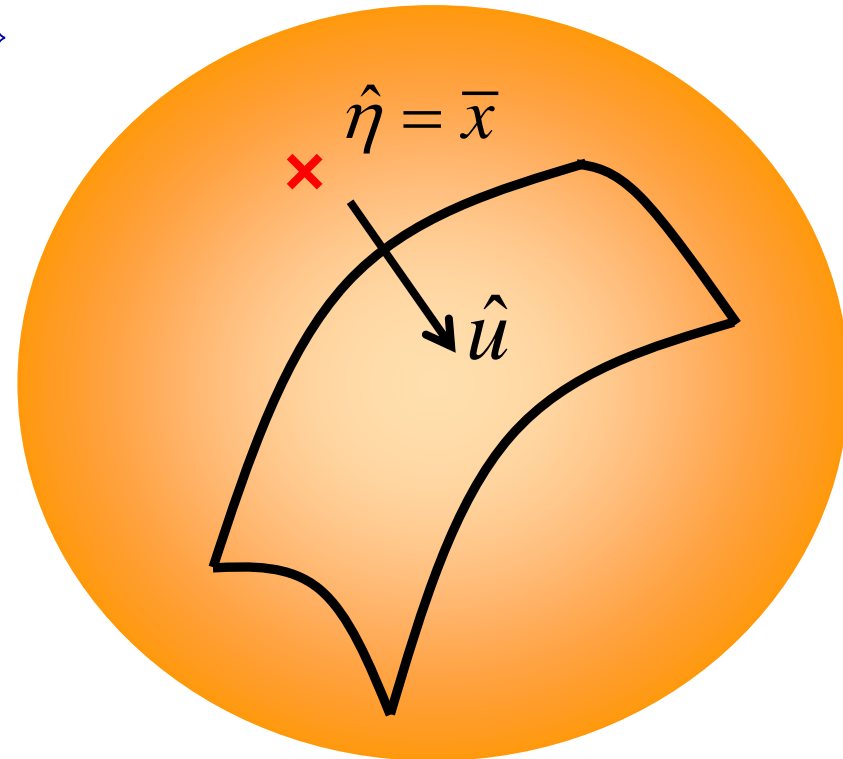
$$p(x, u) \square x_1, x_2, \dots, x_n \quad p(x, \theta) = \exp\{\theta \cdot x - \psi(\theta)\}$$

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

$\hat{u}(x_1, \dots, x_n)$  :推定

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x(k)$$

$H_0 : u = u_0$  :検定



# 推定誤差

$$e = \hat{\eta} - \eta = \frac{1}{N} [\alpha_i - E[\alpha]], \quad \tilde{e} = \sqrt{N} e$$

*(Red arrows point from  $\bar{\alpha}$  to  $\alpha_i$  and from  $\eta$  to  $E[\alpha]$ )*

Cramer-Rao bound

$$E[\tilde{e}_i] = 0$$

$$E[\tilde{e}_i \tilde{e}_j] = \delta_{ij}$$

$$E[\hat{e}_i \hat{e}_j \tilde{e}_k] = \frac{1}{\sqrt{N}} T_{ijk}$$

$$E[\hat{e}_i \hat{e}_j \hat{e}_k \hat{e}_l] = \frac{1}{N} S_{ijkl}$$

$$E[e_i e_j] \geq \frac{1}{N} g_{ij}$$

∇∇∇∇

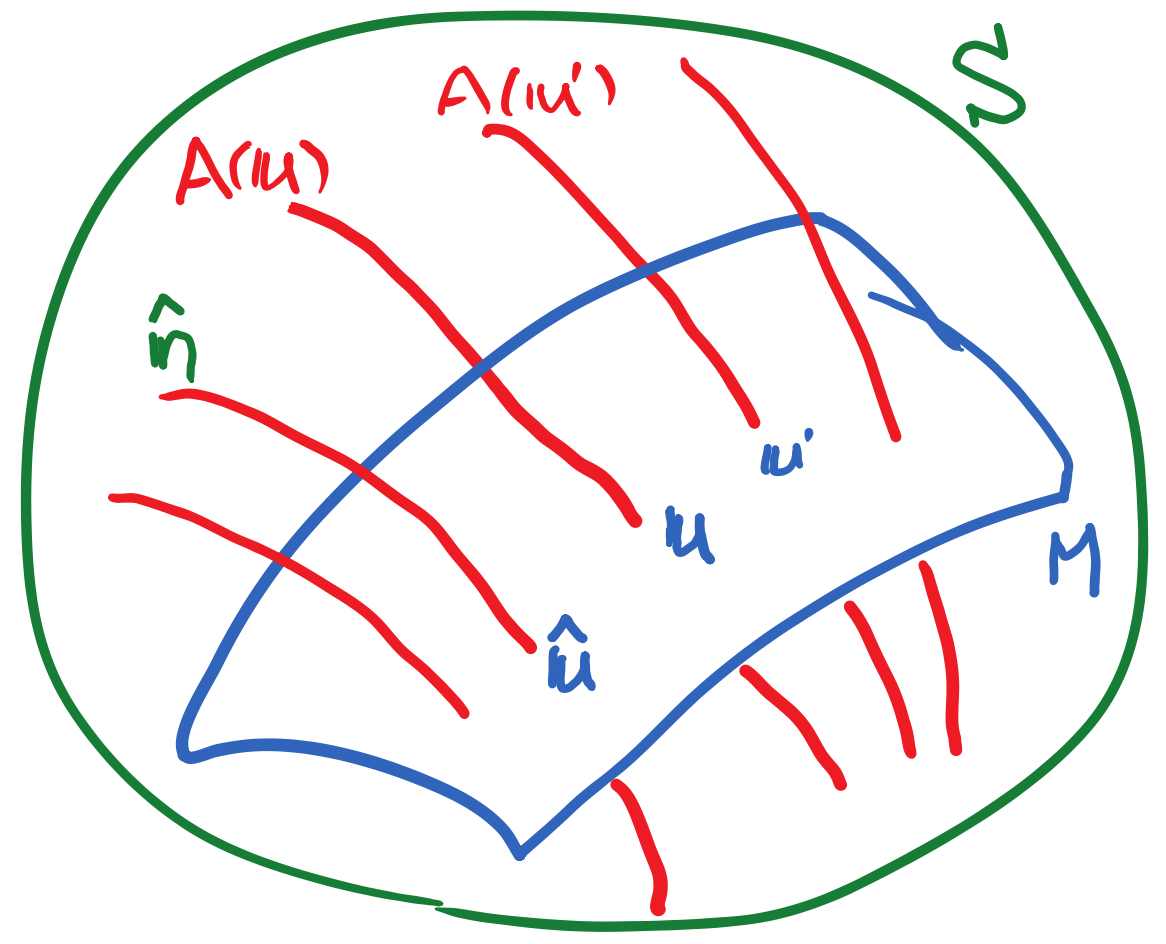
∇∇∇∇

∇∇∇∇

# 補助多様体族

推定量 ---  $\hat{u} = f(\hat{\eta})$

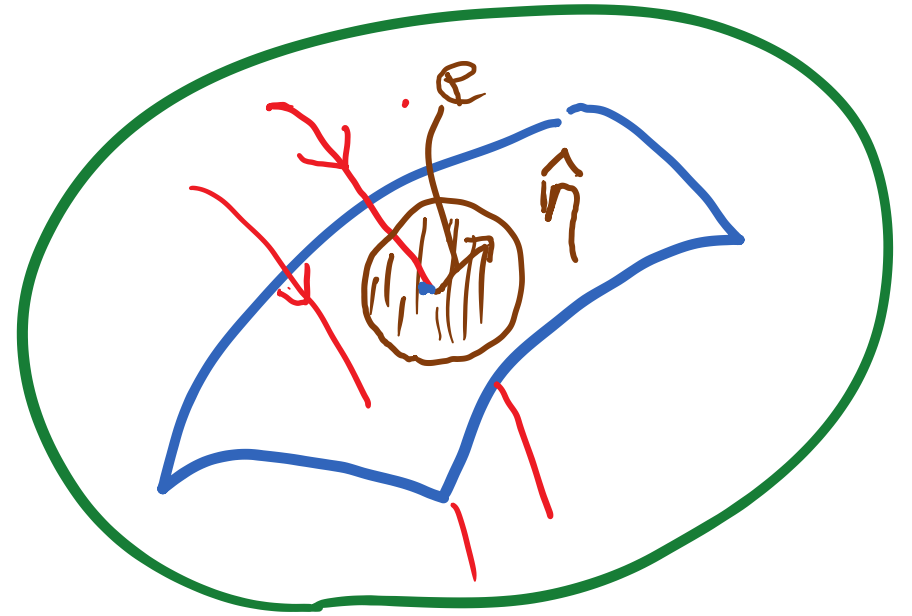
Ancillary family  $A(u)$



# 最尤推定は一致性を持ち有効

$$\hat{\mu}_{MLE} : \min_{\eta \in M} KL[\hat{\eta} : \eta(\mu)]$$

$M$ -projection of  $\hat{\eta}$  to  $M$



Efficient estimator --- orthogonal projection



# 誤差の高次漸近理論

$$p(x, \theta(u)) \quad : x_1, \dots, x_N$$

$$\hat{u} = u(x_1, \dots, x_n)$$

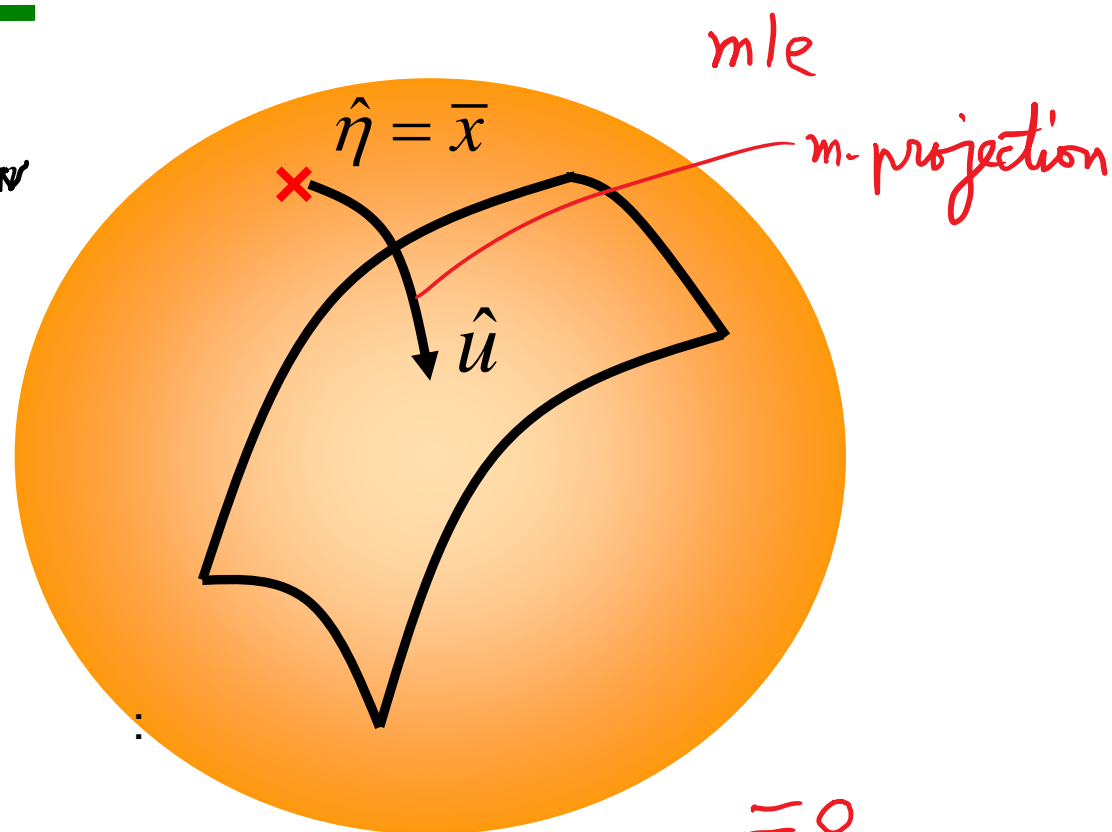
$$e = E \left[ (\hat{u} - u)(\hat{u} - u)^T \right]$$

$$e = \frac{1}{N} G_1 + \frac{1}{N^2} G_2$$

$$G_1 \geq G^{-1} \quad : \text{Cramér-Rao: linear theory}$$

$$G_2 = H_M^{(e)^2} + H_A^{(m)^2} + \Gamma^{(m)^2}$$

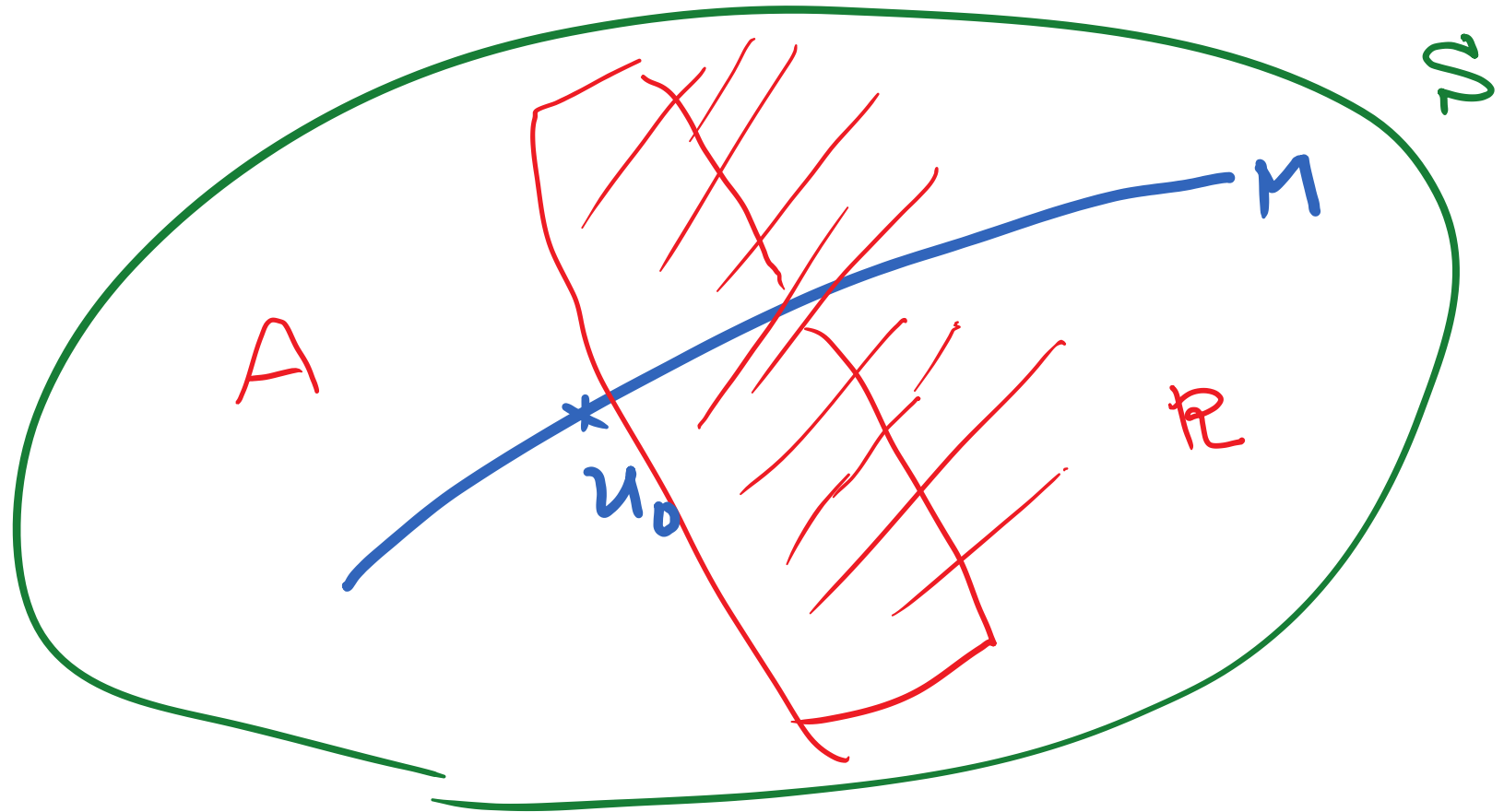
quadratic approximation



# 仮説検定

$$H_0 : \mu = \mu_0$$

$$H_1 : \mu = \mu_1$$





# Neyman-Scott問題： 無限個の局外母数

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

$$x_1 \square p(x, \theta, \xi_1)$$

$$x_2 \square p(x, \theta, \xi_2)$$

-----

$$x_N \square p(x, \theta, \xi_N)$$

$\theta$  : parameter of interest

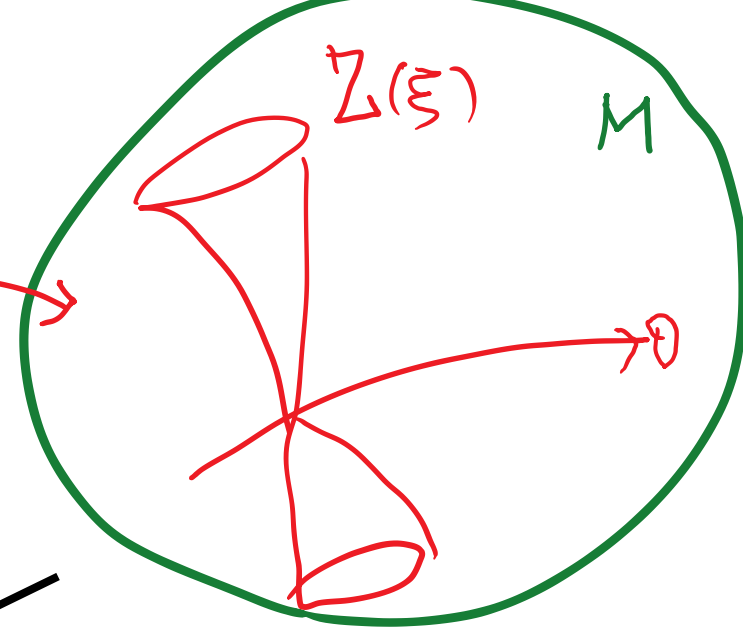
$\xi$  : nuisance parameter

# Semiparametric 統計モデル: 比例定数の推定

$$M = \{p(x, \theta, Z)\}$$

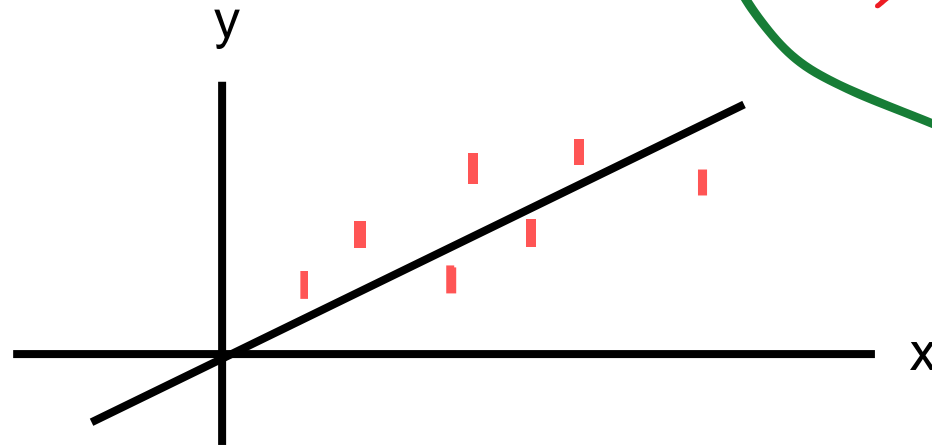
$$\xi \square Z(\xi)$$

関数自由度の未知母数



linear relation  $\mathbf{x} = (x, y)$

$$y = \theta x$$



$$\begin{cases} y_i = \theta \xi_i + \varepsilon_i \\ x_i = \xi_i + \varepsilon_i' \end{cases} \quad p(x, y; \theta, Z) = \int p(x, y; \xi, \theta) Z(\xi) d\xi$$

mle, least square, total least square

# 統計 Model

$$p(x, y|\theta, \xi) = c \exp \left\{ -\frac{1}{2}(x - \xi)^2 - \frac{1}{2}(y - \theta\xi)^2 \right\}$$

$$\prod p(x_i, y_i|\theta, \xi_i) : \theta, \xi_1, \dots, \xi_n$$

$$p(x, y|\theta, Z) = \int p(x, y|\theta, \xi)Z(\xi)d\xi$$

———— semiparametric

# 最小二乗法は良いか？

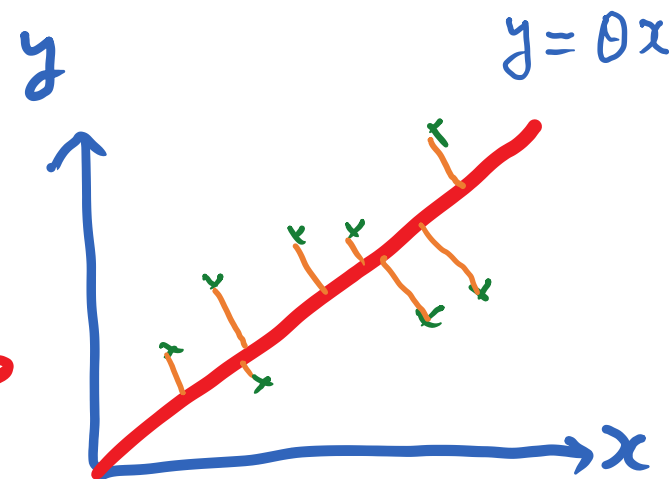
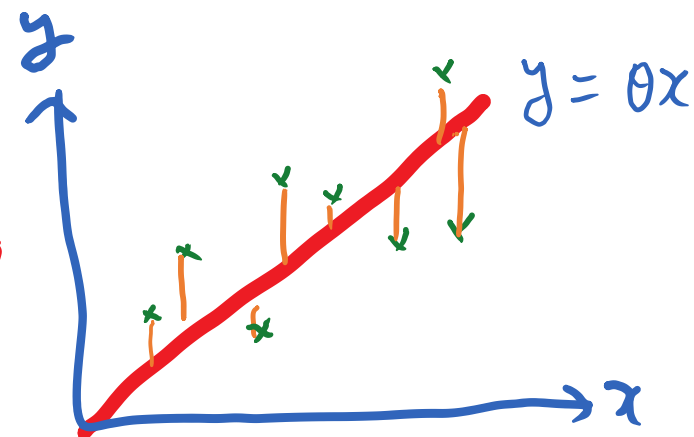
$$L(\theta) = \sum (y_i - \theta x_i)^2 \rightarrow \min \quad : \hat{\theta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

$$\frac{1}{n} \sum \frac{y_i}{x_i} : \text{averag} \quad \frac{\sum y_i}{\sum x_i} : \text{gross average}$$

mle, TLS

$$\sum (y_i - \theta x_i)(\theta y_i + x_i) = 0$$

Neyman-Scott



# セミパラ統計モデル

$$x_1, x_2, \dots \square p(x, \theta, Z)$$

推定関数

$$f(x, \theta)$$

$$E_{\theta, Z} [f(x, \theta)] = 0$$

$$E_{\theta', Z} [f(x, \theta)] \neq 0 \quad \theta' \neq \theta$$

推定方程式

$$\sum f(x_i, \theta) = 0 \quad \Rightarrow \hat{\theta}$$

$$\frac{1}{N} \sum f(x_i, \theta) \Rightarrow E_{\theta, Z} [f(x, \theta)]$$



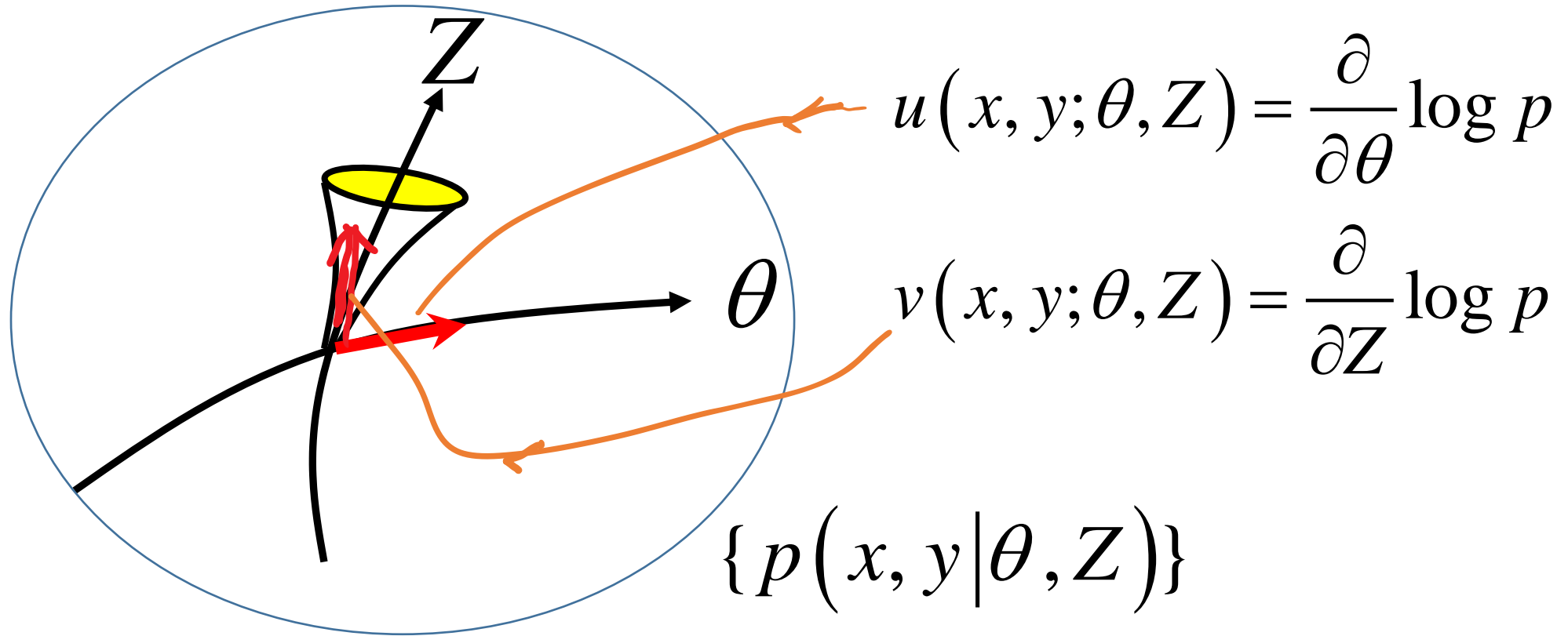
# 推定関数

$$E_{\theta, z} [f(x, \theta)] = 0: \text{ unbiased}$$

$$\sum_{i=1}^n f(x_i, \hat{\theta}) = 0: \hat{\theta} = \theta + e$$

$$E \left[ (\hat{\theta} - \theta)^2 \right] = \frac{1}{n} \frac{E[f^2]}{E[(\partial_{\theta} f)^2]}$$

# Fiber Bundle



# Estimating Function $f(x, \theta)$

$$e\text{-invariant} : E_{\theta, z} [f(x, \theta)] = 0$$

$$\prod_z^e f(x, \theta) = f$$

$$T_\theta = T_\theta^I \oplus T_\theta^N \oplus T_\theta^A$$

$$m\text{-orthogonality} : \langle v, f \rangle = 0$$

$$\left\langle \prod_z^m v, f \right\rangle = 0$$

$$\int p(x, \theta, \xi) z(\xi) f(x, \theta) dx d\xi$$

$$\langle \delta z, f \rangle = 0$$

$u^I(x, \theta, z)$ : optimal estimating function

## Efficient Score

$$\partial_{\theta} \ell = \partial_{\theta} \log p(x, \theta, \xi)$$

$$\bar{\partial}_{\theta} \ell = \partial_{\theta} \ell - \mathbb{E}_{\text{aux}} g^{xx} \partial_{\xi} \ell$$

$$\dot{\ell}^E(x, \theta, Z) = \int \bar{\partial}_{\theta} \ell(x, \theta, \xi) Z(\xi) d\xi$$

orthogonal

$$f(x, \theta) = \dot{\ell}^E(x, \theta, Z_0) + a(x)$$

$$\sum f(x_i, y_i; \theta) = 0$$

$$f(x, y; \theta) = (x + \theta y + c)(y - \theta x)$$

$$c = \frac{\bar{\xi} \sigma^2}{\bar{\xi}^2 - (\bar{\xi})^2} \quad \begin{cases} \bar{\xi} = 1 \\ \bar{\xi}^2 = 2 \end{cases}$$

$$c = 0: \quad V = \frac{1}{n} \frac{(2 + \sigma^2) \sigma^2}{4} \quad : \frac{3}{4}$$

$$c = 1: \quad V = \frac{1}{n} \left( 1 - \frac{1}{\sigma^2 + 2} \right) \sigma^2 \quad : \frac{2}{3}$$

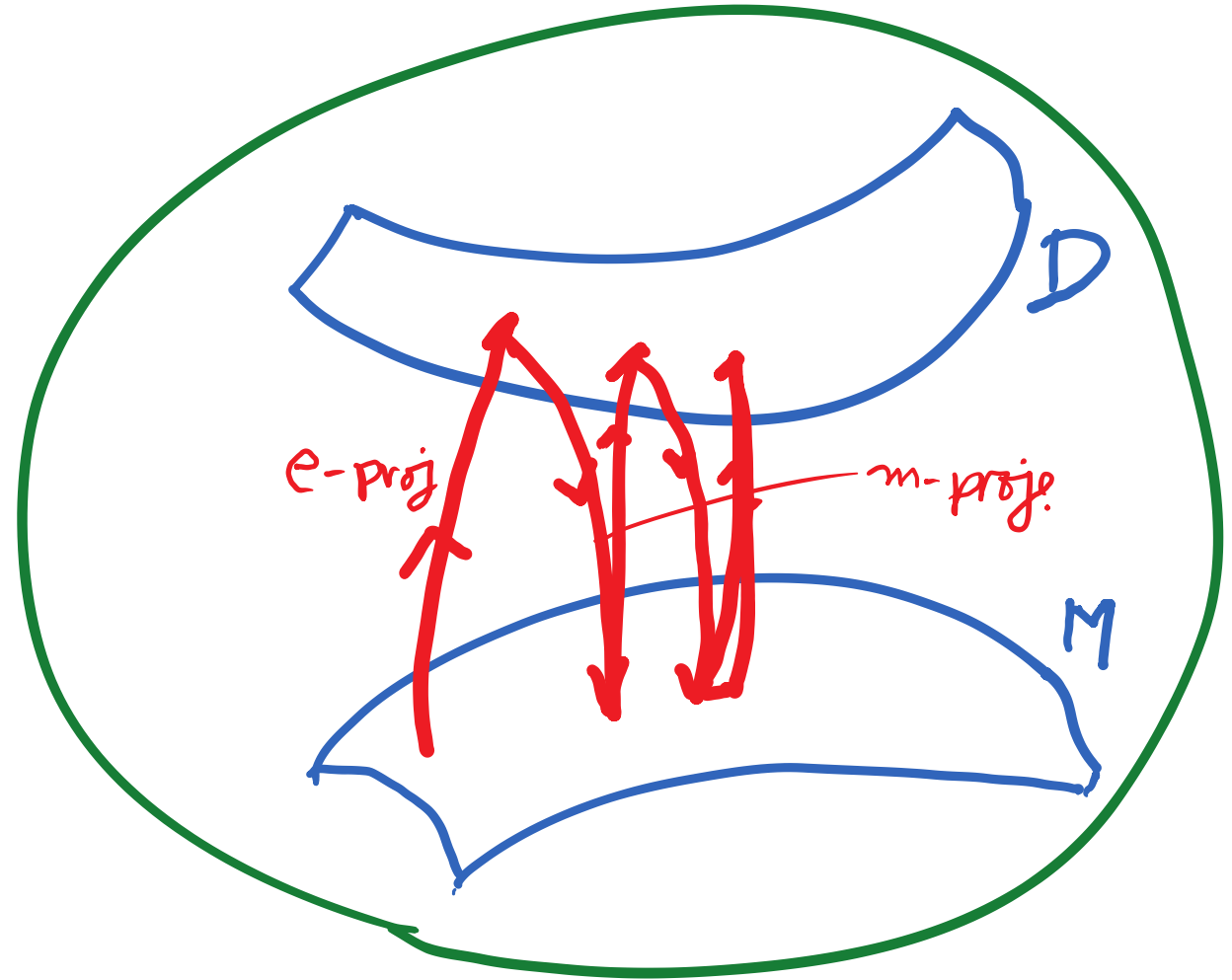
$$c = \infty: \quad v = \frac{1}{n} \sigma^2 \quad : 1$$

# em-algorithm EM-algorithm

## Variational Bayes

$$\min D_{KL} [q(x) : p(x)]$$

$$q(x) \in \mathcal{D}, \quad p(x) \in \mathcal{M}$$



# EM algorithm

*hidden*

hidden variables

*observe*

$$p(\mathbf{x}, \mathbf{y}; \mathbf{u})$$

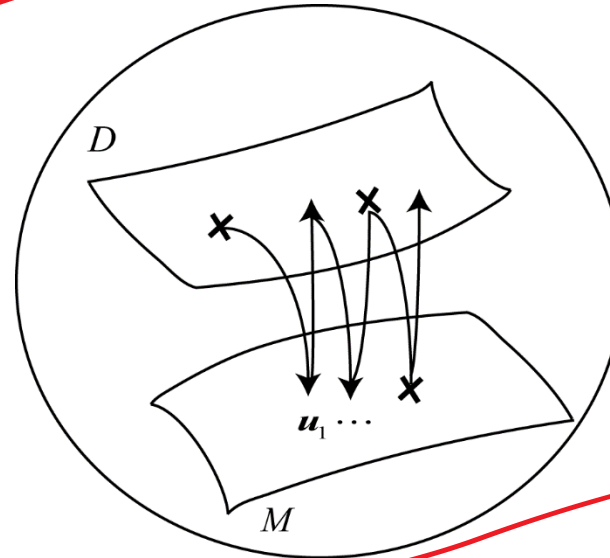
$$D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$

$$M = \{p(\mathbf{x}, \mathbf{y}; \mathbf{u})\}$$

$$D_M = \{\mathcal{F}(\mathbf{x}, \mathbf{y}) \mid \mathcal{F}(\mathbf{x}) = \mathcal{F}_D(\mathbf{x})\}$$

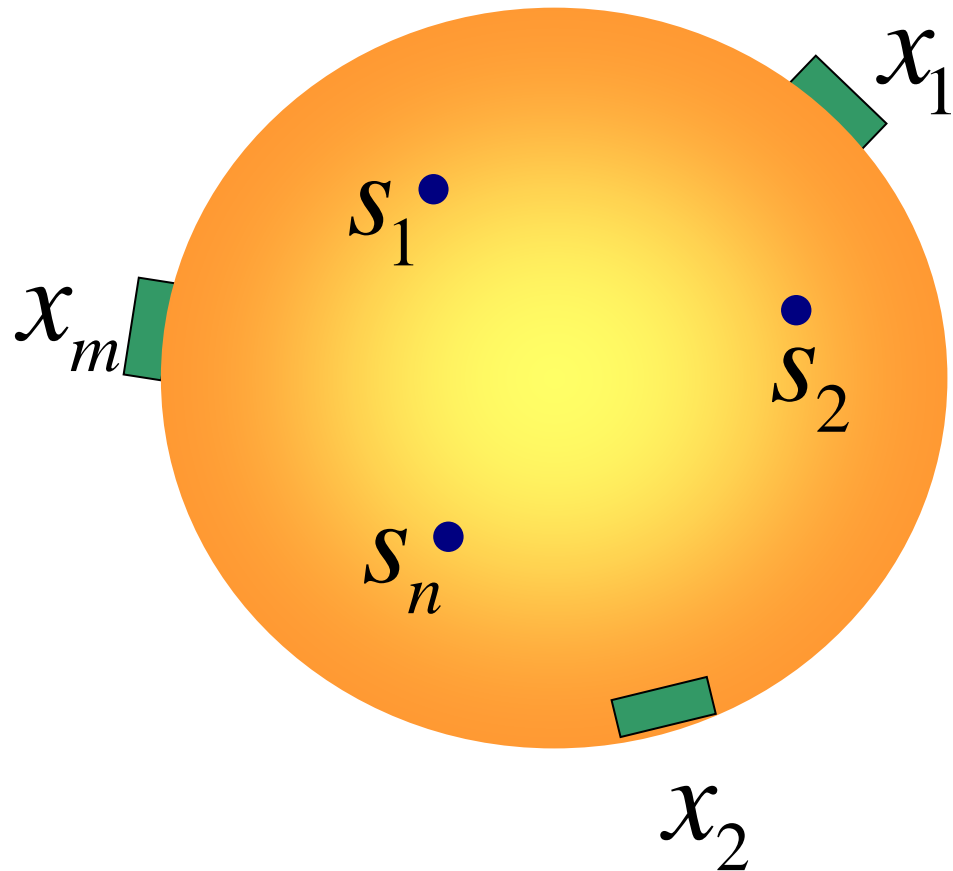
$$\min KL[\hat{p}(\mathbf{x}, \mathbf{y}) : p \in M] \quad \text{m-projection to } M$$

$$\min KL[p \in D : p(\mathbf{x}, \mathbf{y}; \hat{\mathbf{u}})] \quad \text{e-projection to } D$$



$$\left\{ \begin{aligned} \mathcal{F}_D &= \frac{1}{N} \sum \delta(\mathbf{x} - \mathbf{x}_i) \\ \mathcal{F}(\mathbf{x}, \mathbf{y}) &= \mathcal{F}_D(\mathbf{x}) r(\mathbf{y} | \mathbf{x}) \end{aligned} \right.$$

# 信号の混合と分解



生体情報解析

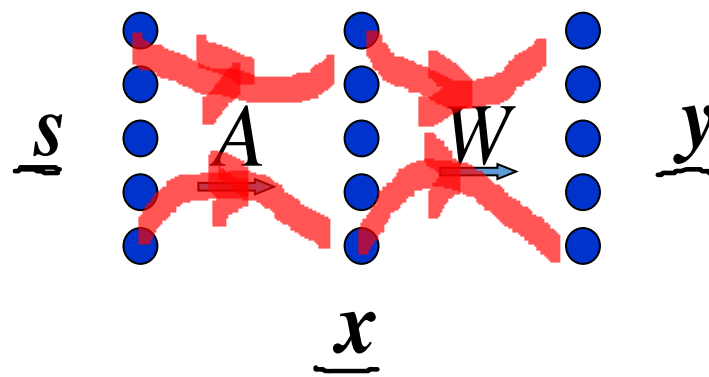
カクテルパーティ効果

移動体通信

画像解析



# 独立成分分析



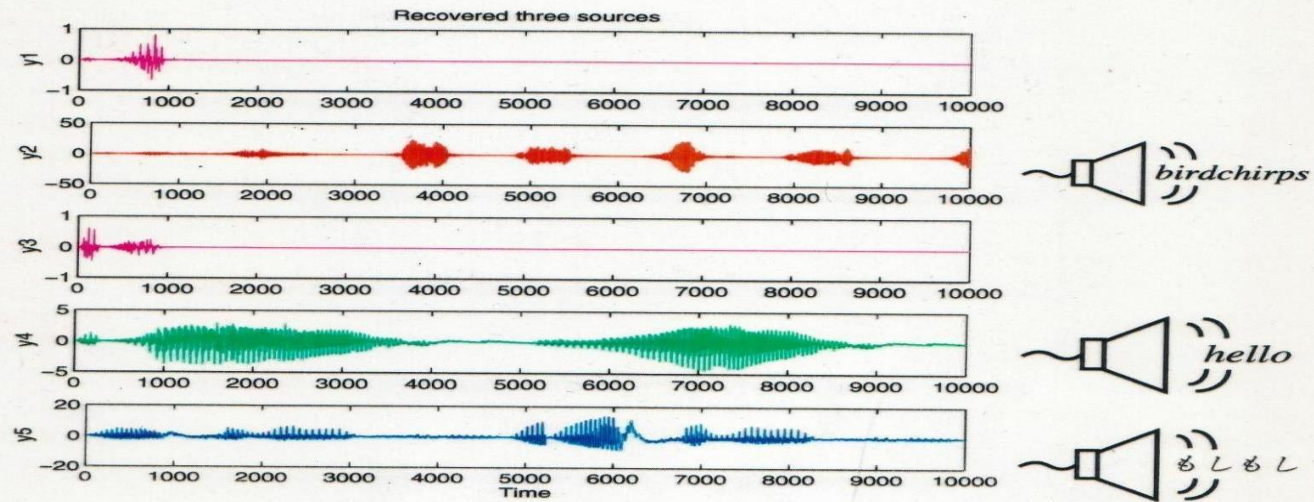
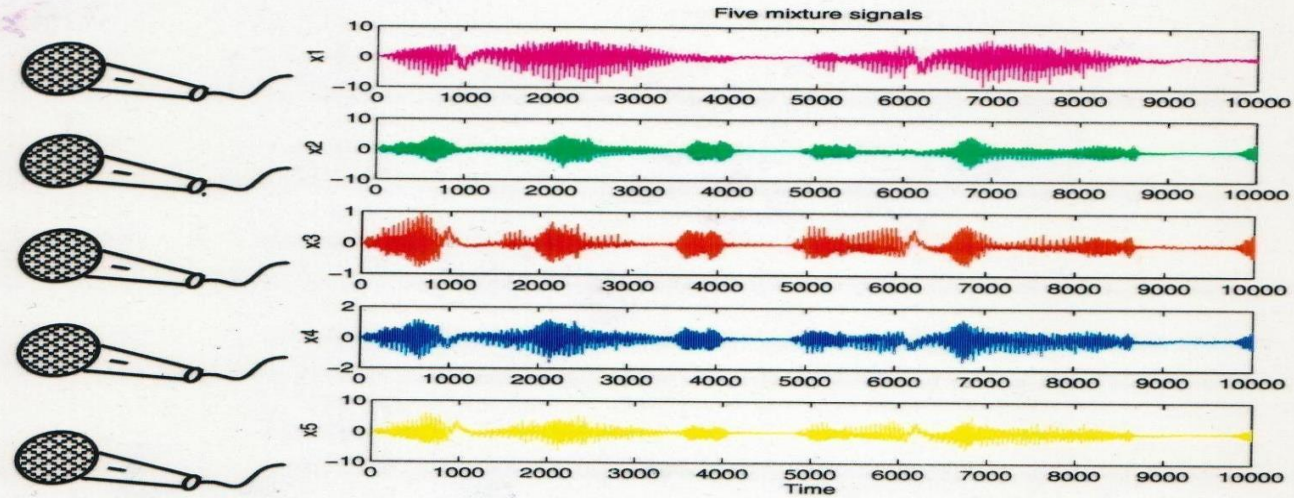
$$\underline{x} = A\underline{s} \quad x_i = \sum A_{ij} s_j$$

$$\underline{y} = W\underline{x} \quad W = A^{-1}$$

觀測信号:  $x(1), x(2), \dots, x(t)$

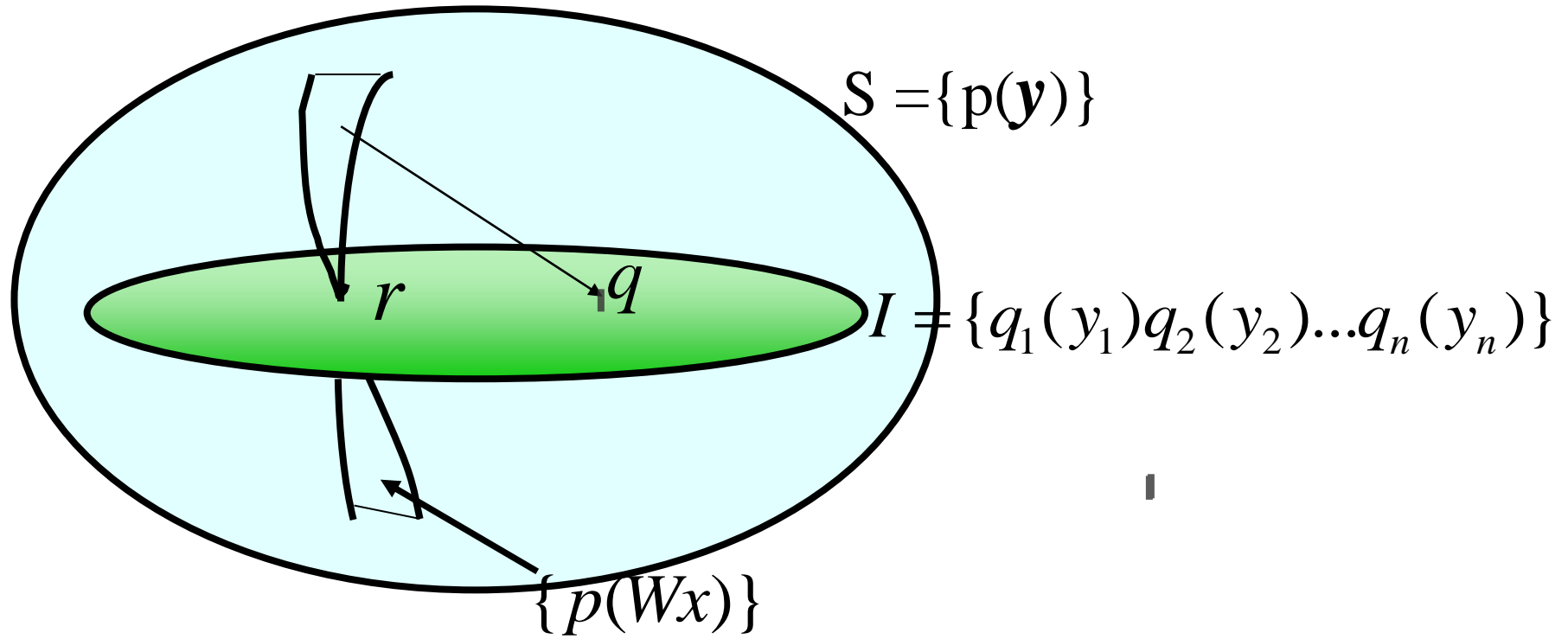
復元信号:  $s(1), s(2), \dots, s(t)$

# Cocktail party experiment



- 5 microphones (sensors) and only 3 speakers

# 情報幾何による評価関数



$$l(\mathbf{W}) = KL[p(\mathbf{y}; \mathbf{W}) : q(\mathbf{y})]$$

$$r(\mathbf{y})$$

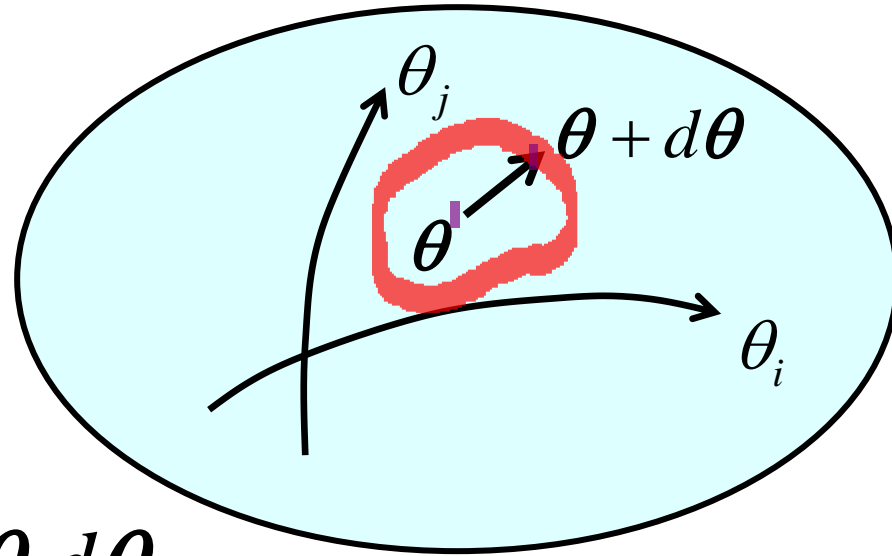
# 行列Wの空間: $GL(n)$ リーマン空間

$$\theta \rightarrow W$$

$$ds^2 = |d\theta|^2$$

$$= \sum g_{ij}(\theta) d\theta_i d\theta_j$$

$$= d\theta^T G(\theta) d\theta$$



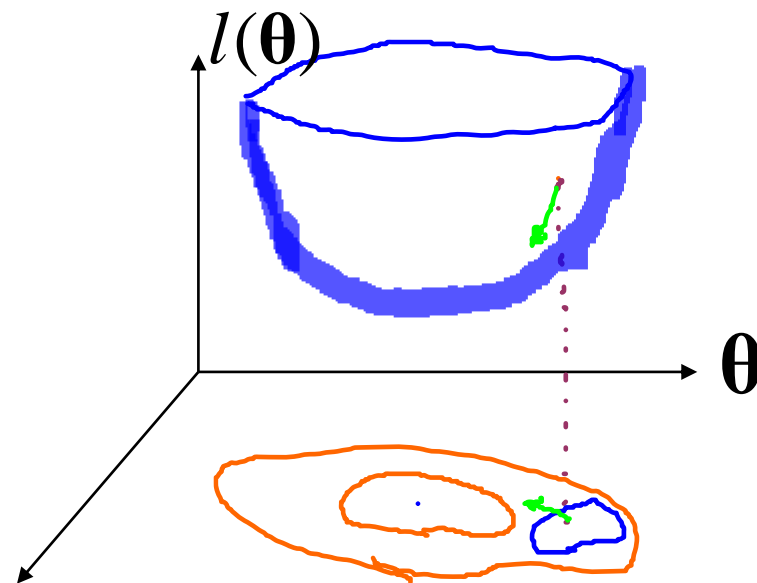
Euclid:  $G = I$

# 自然勾配 (Natural Gradient)

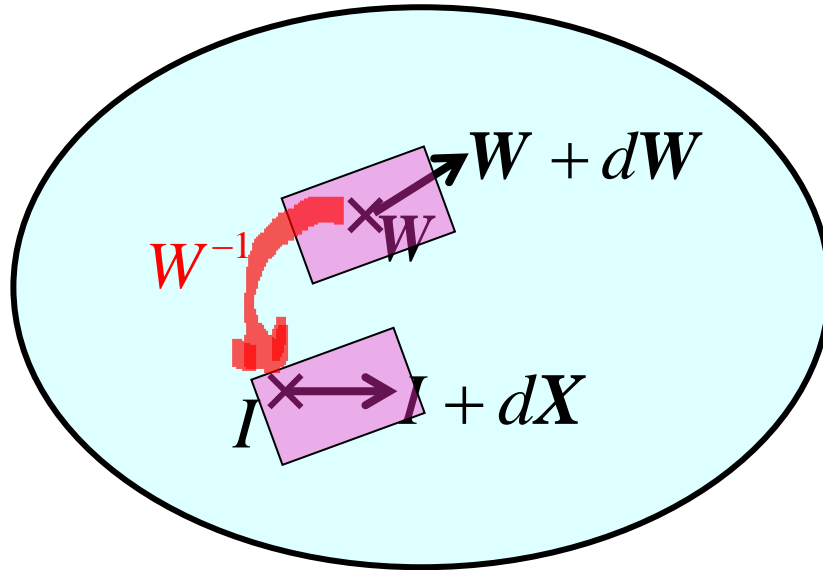
$$\max \quad dl = l(\boldsymbol{\theta} + d\boldsymbol{\theta}) - l(\boldsymbol{\theta})$$

$$|d\boldsymbol{\theta}|^2 = \varepsilon$$

$$\nabla l = G^{-1}(\boldsymbol{\theta}) \nabla l$$



# 行列の空間: Lie群



$$dX = dW W^{-1}$$

$$|dW|^2 = \text{tr}(dX dX^T) = \text{tr}(dW W^{-1} W^{-T} dW^T)$$

$$\nabla l = \frac{\partial l}{\partial W} W^T W$$

$dX$  : non-holonomic basis

# 自然勾配

$$\begin{aligned}\Delta \mathbf{W} &= -\eta \mathbf{G}^{-1} \frac{\partial l(\mathbf{y}, \mathbf{W})}{\partial \mathbf{W}} \\ &= -\eta \frac{\partial l(\mathbf{y}, \mathbf{W})}{\partial \mathbf{W}} \mathbf{W}^T \mathbf{W}\end{aligned}$$

**Example of color image separation :**

Five original images (but unknown to the neural net)



Five mixed images for separation



Final (stable states) of five separated images



# ICAから派生したもの

$$x = As$$

非負行列分解  
スパース信号解析



(a) Three binary edge images (reverse images are used in the experiment)



(b) Two edge image mixtures



(c) Reconstructed binary edge images (after reversion)

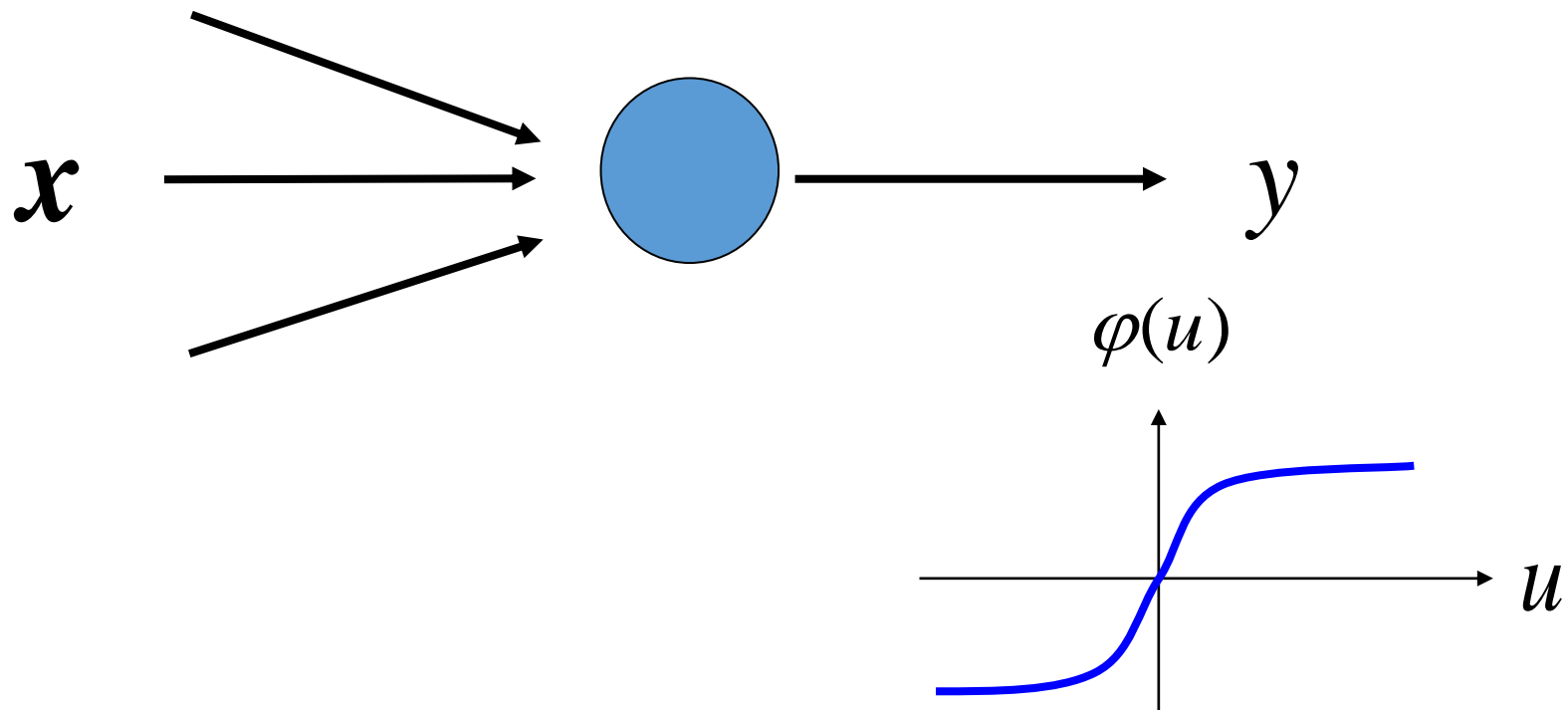
Fig. 5: Example of edge image image reconstruction: (a) the three binary edge images (reverse image copies are supplied for processing) , (b) their two mixtures, (c) the three extracted edge images (after reversion).

# 多層パーセプトロンの情報幾何

Natural Gradient and Singularities

# 数理ニューロン

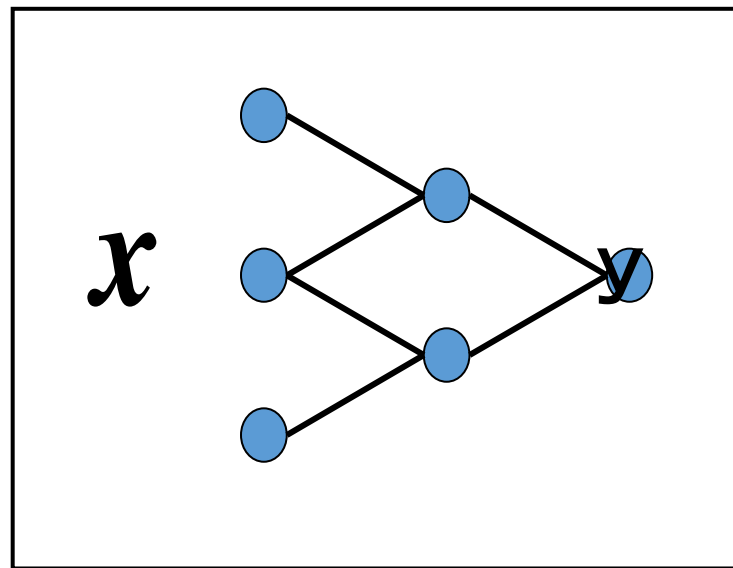
$$y = \varphi\left(\sum w_i x_i - h\right) = \varphi(\mathbf{w} \cdot \mathbf{x})$$



# 多層パーセプトロン

$$y = \sum v_i \varphi(\mathbf{w}_i \cdot \mathbf{x}) + n$$

$$\mathbf{x} = (x_1, x_2, \dots, x_n)$$

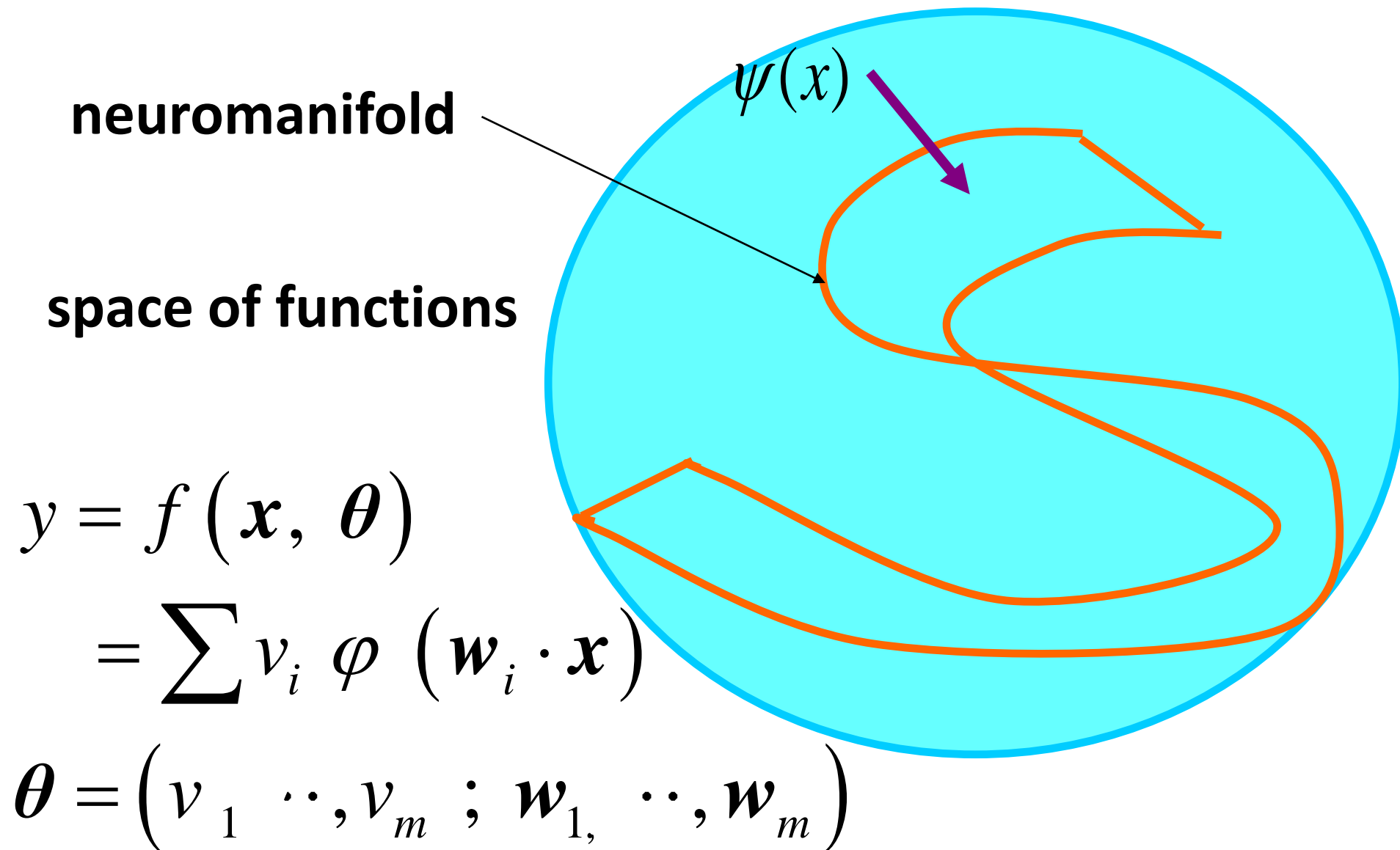


$$p(y|\mathbf{x};\boldsymbol{\theta}) = c \exp\left\{-\frac{1}{2}(y - f(\mathbf{x},\boldsymbol{\theta}))^2\right\}$$

$$f(\mathbf{x},\boldsymbol{\theta}) = \sum v_i \varphi(\mathbf{w}_i \cdot \mathbf{x})$$

$$\boldsymbol{\theta} = (\mathbf{w}_1, \dots, \mathbf{w}_m; v_1, \dots, v_m)$$

# 多層パーセプトロンと神経多様体



## 例題からの学習

$$\psi(\mathbf{x}) \approx f(\mathbf{x}, \hat{\theta})$$

多数の例題  $\dots D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

learning ; estimation

# Backpropagation --- 確率降下学習

examples :  $(y_1, \mathbf{x}_1), \dots, (y_t, \mathbf{x}_t)$  -- training set

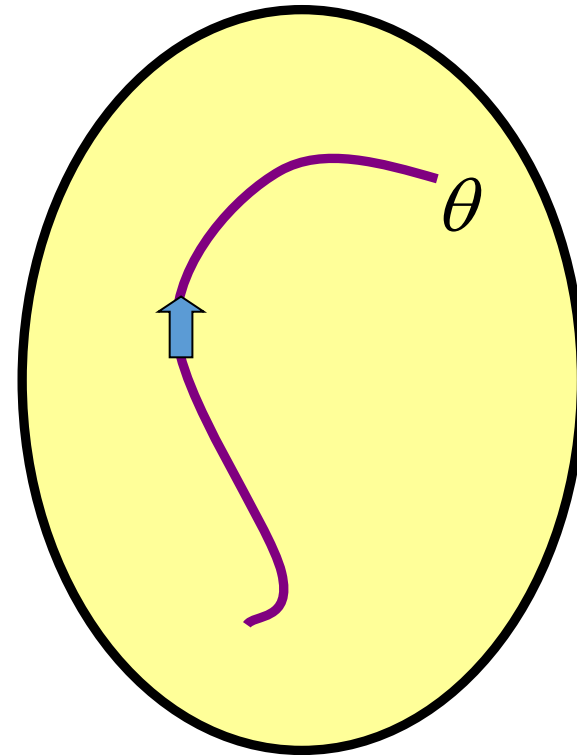
$$y = f(x, \theta) + n$$

$$E(y, x; \theta) = \frac{1}{2} |y - f(x, \theta)|^2$$

$$= -\log p(y, x; \theta)$$

$$\Delta \theta_t = -\eta_t \frac{\partial E}{\partial \theta}$$

$$f(x, \theta) = \sum v_i \varphi(\mathbf{w}_i \cdot \mathbf{x})$$

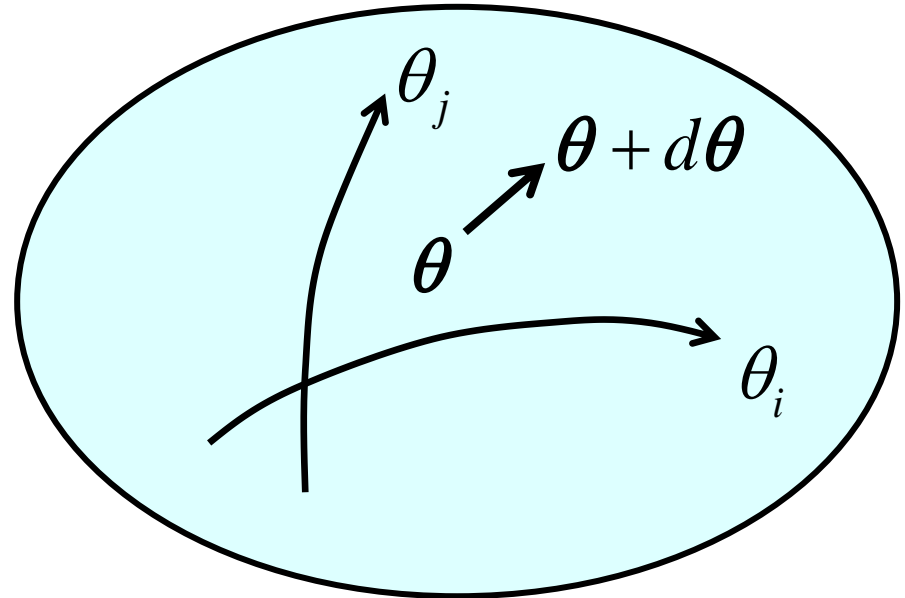




# 計量: 実はリーマン空間であった

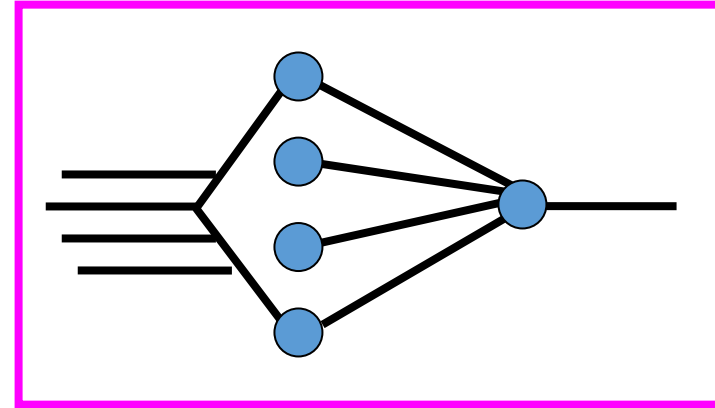
$$g_{ij}(\boldsymbol{\theta}) = E\left[\frac{\partial \log p(y | x; \boldsymbol{\theta}) \partial \log p(y | x; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j}\right]$$

$$\begin{aligned} ds^2 &= |d\boldsymbol{\theta}|^2 \\ &= \sum g_{ij}(\boldsymbol{\theta}) d\theta_i d\theta_j \\ &= d\boldsymbol{\theta}^T G(\boldsymbol{\theta}) d\boldsymbol{\theta} \end{aligned}$$

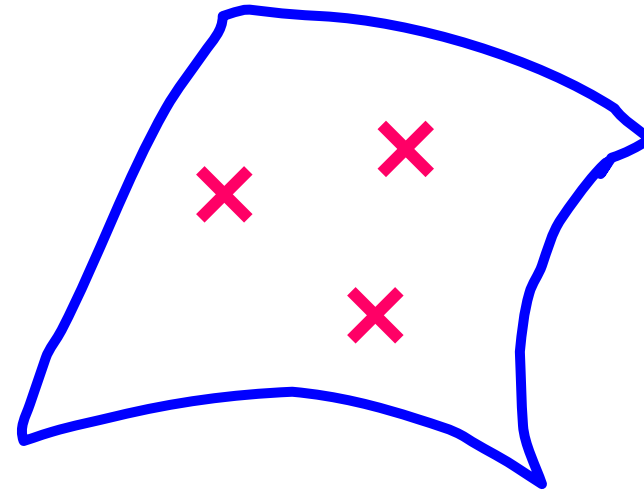
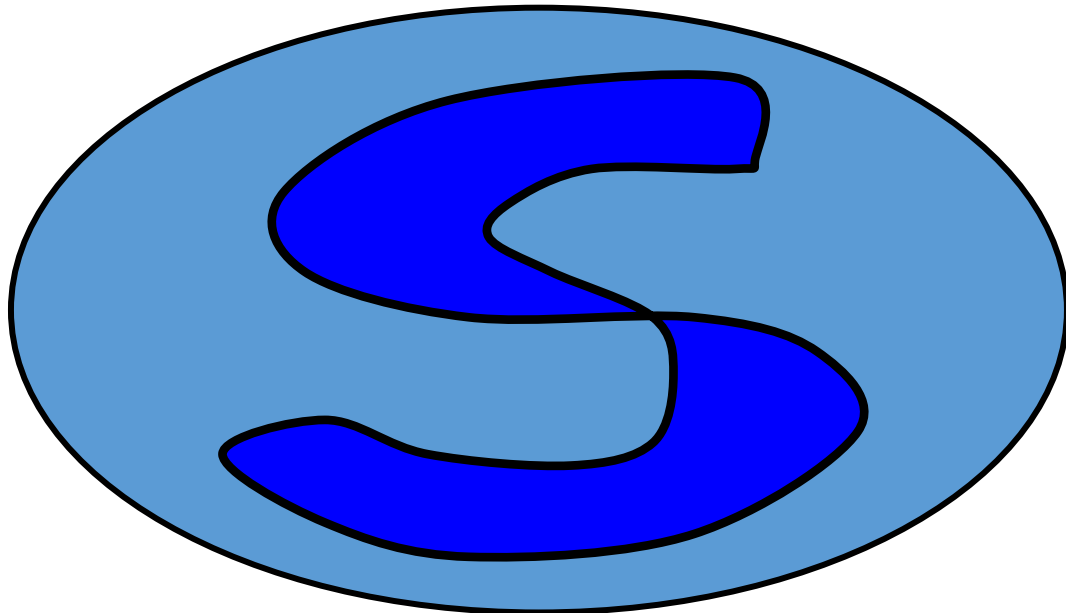


# Topology: Neuromanifold

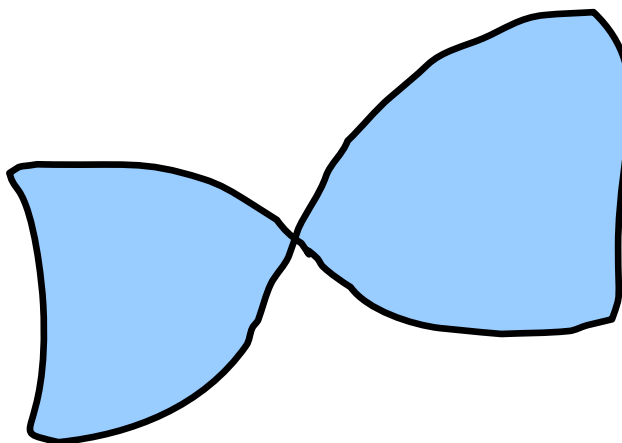
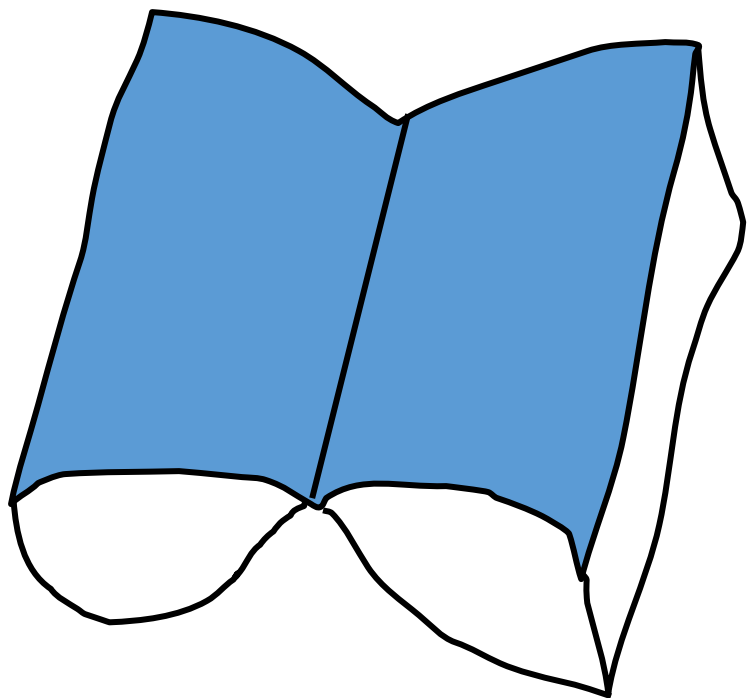
- Metrical structure
- Topological structure



$\theta$



# singularities



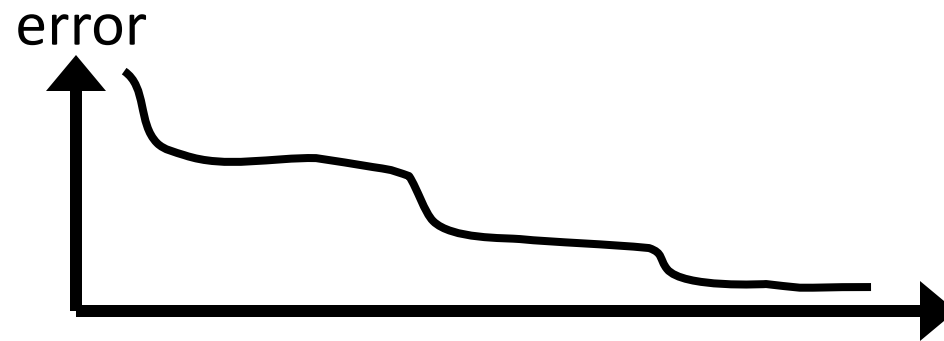
# Backprop の問題点

$$\Delta\theta_t = -\eta_t \nabla l(x_t, y_t; \theta_t)$$

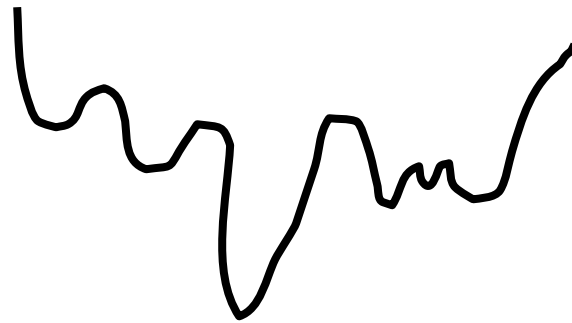
- **slow convergence---plateau---saddle**
- **local minima**

# MLP学習の欠陥

slow convergence : plateau

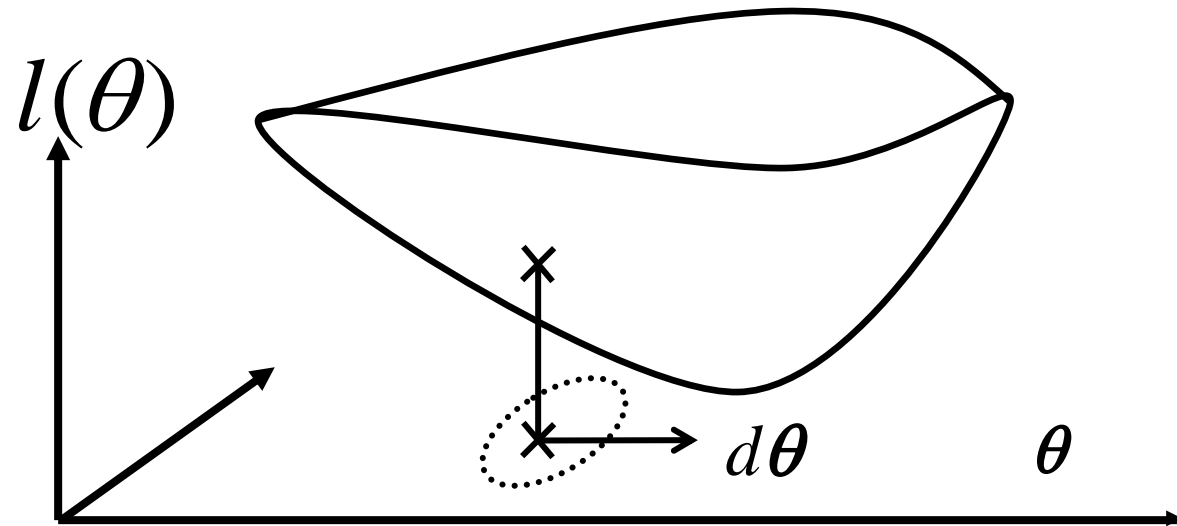


local minima



➡ Boosting, Bagging, SVM

# 最急降下方向--- **Natural Gradient**



$$\nabla l = \left( \frac{\partial l}{\partial \theta_1}, \dots, \frac{\partial l}{\partial \theta_n} \right)$$

$$\Delta \theta_t = -\eta_t \nabla l(x_t, y_t; \theta_t)$$

$$\nabla l = G^{-1}(\theta) \nabla l$$

$$|d\theta|^2 = d\theta^T G d\theta = \sum G_{ij} d\theta^i d\theta^j$$

# 自然勾配學習 Natural Gradient

$$\max \quad dl = l(\boldsymbol{\theta} + d\boldsymbol{\theta}) - l(\boldsymbol{\theta}) = \nabla l \cdot d\boldsymbol{\theta}$$

$$\text{under } |d\boldsymbol{\theta}|^2 = \sum g_{ij} d\theta_i d\theta_j = \varepsilon^2$$

$$d\boldsymbol{\theta} \approx \nabla l = G^{-1}(\boldsymbol{\theta}) \nabla l$$

$$\Delta \boldsymbol{\theta}_t = -\eta_t \tilde{\nabla} l(x_t, y_t; \boldsymbol{\theta}_t)$$

# MLPの情報幾何

Natural Gradient Learning :  
S. Amari ; H.Y. Park

$$\Delta \boldsymbol{\theta} = -\eta G^{-1}(\boldsymbol{\theta}) \frac{\partial l}{\partial \boldsymbol{\theta}}$$

Adaptive natural gradient learning

$$G_{t+1}^{-1} = (1 + \varepsilon) G_t^{-1} - \varepsilon G_t^{-1} \nabla f \nabla f^T G_t^{-1}$$



# Landscape of error at singularity

Milner attractor

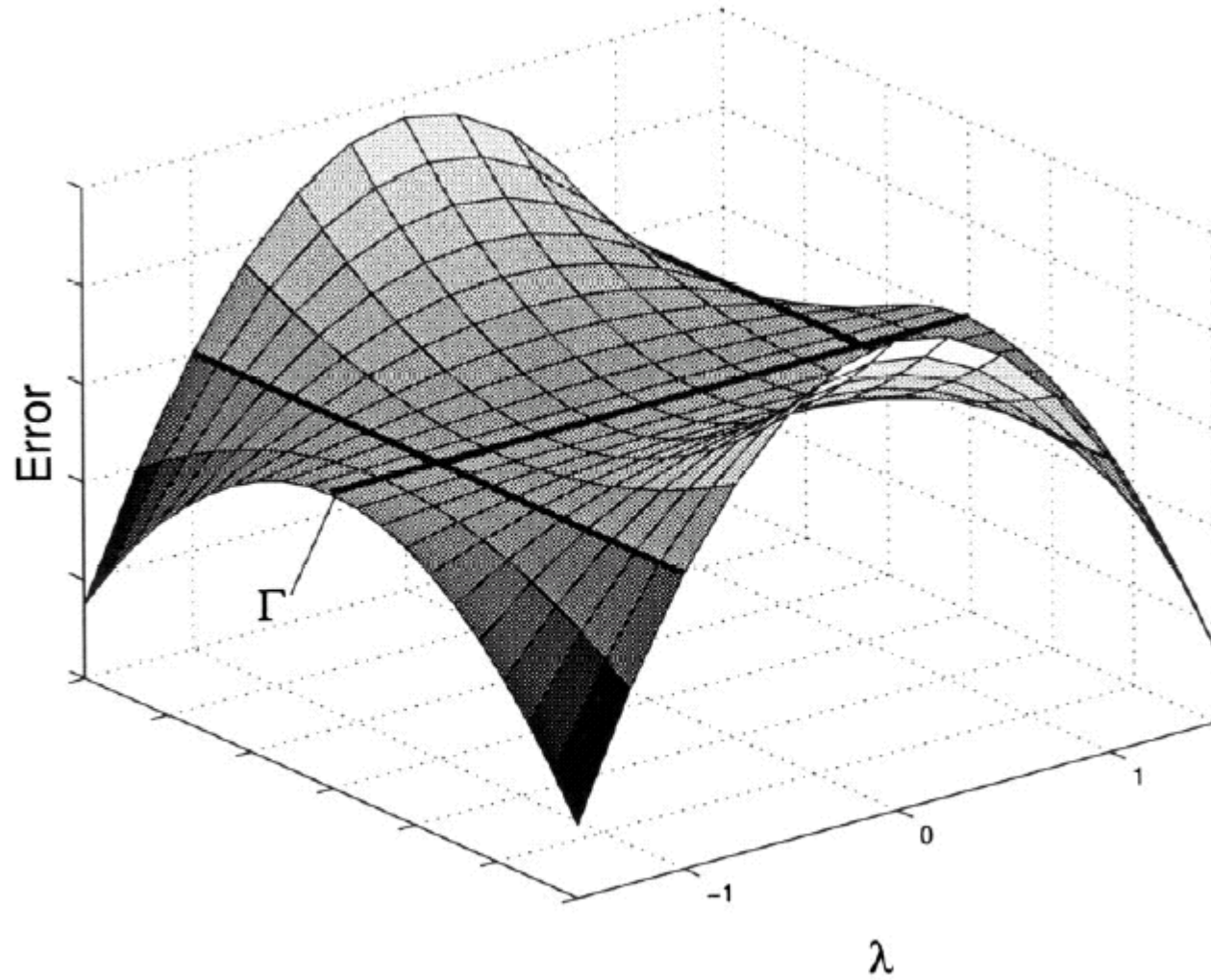


Fig. 5. Critical set with local minima and plateaus.

# 統計神経力学

Rozonoer (1969)

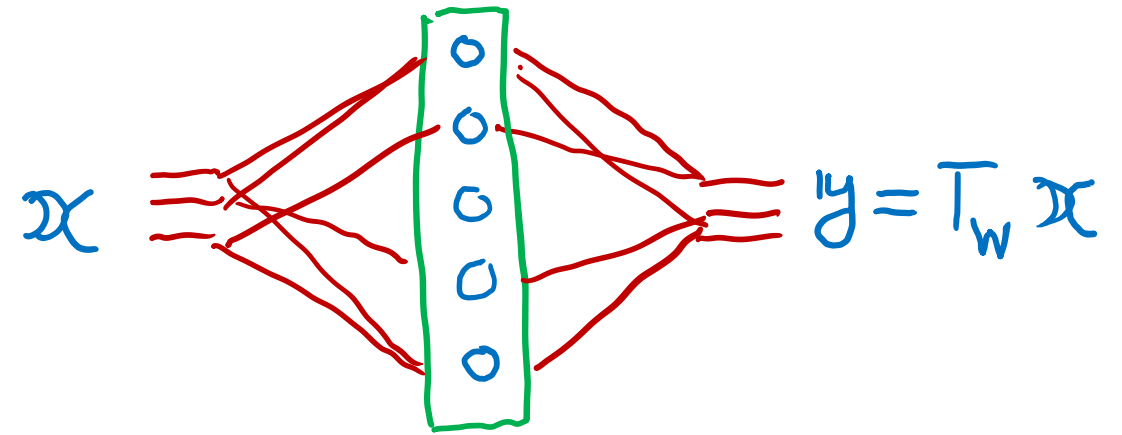
Amari (1971; 1974)

Amari et al (2013)

Toyoizumi et al (2015)

Poole, ..., Ganguli (2016)

Schoenholz et al (2017)



$$w_{ij} \sim N(0, 1)$$

## 巨視的振舞い

ほとんどすべての(典型的)回路に共通

# 巨視變數

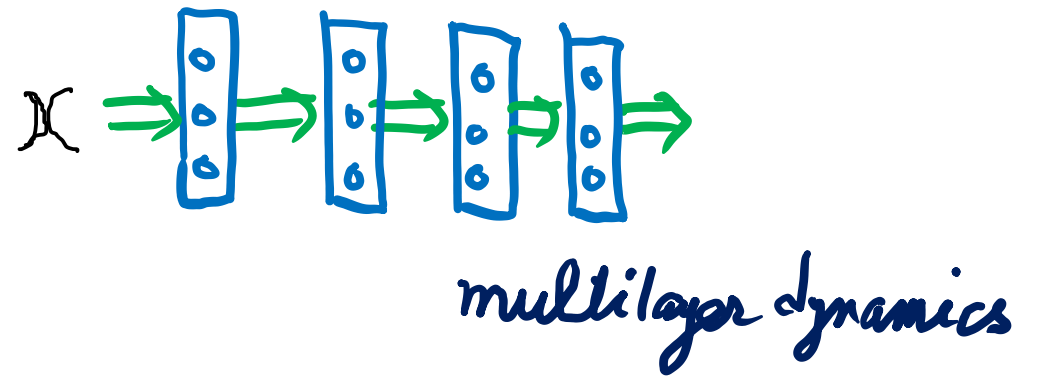
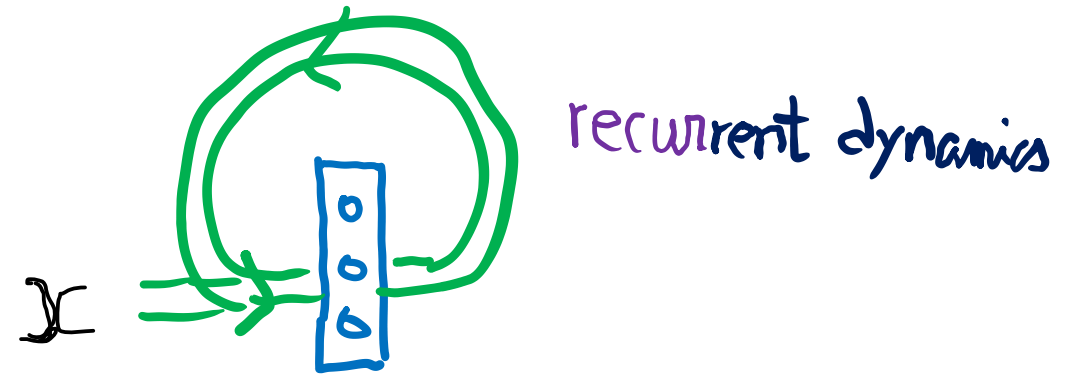
活動度:  $A = \frac{1}{n} \sum x_i^2$

距離・計量:  $D = D[\mathbf{x} : \mathbf{x}']$

曲率:

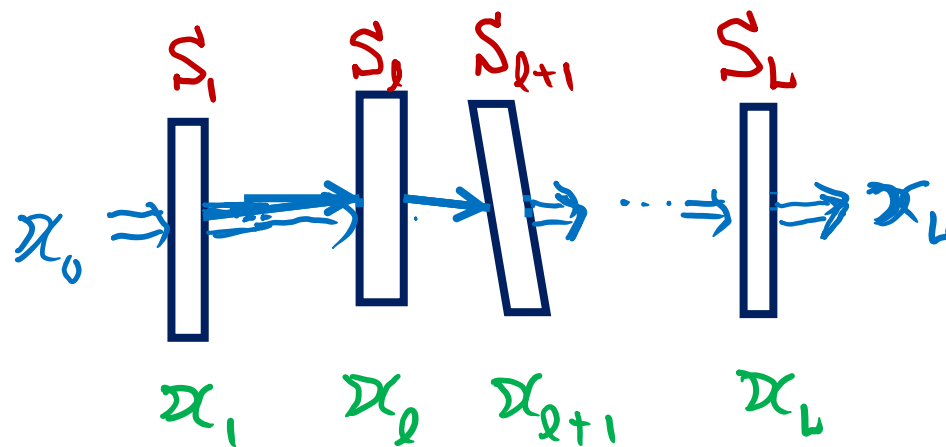
$$A_{l+1} = F(A_l)$$

$$D_{l+1} = K(D_l)$$



# 深層回路

$$x_{l+1} = \varphi\left(\sum w_{ij} x_l + w_{0i}\right)$$



$$A_l = \frac{1}{n_l} \sum x_l^2$$

$$w_{ij} \sim N(0, \sigma^2 / \sqrt{n})$$

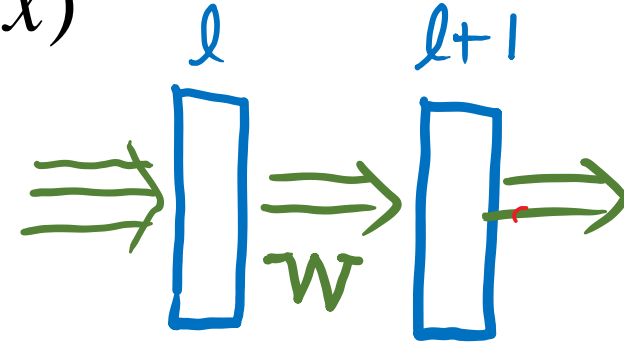
$$A_{l+1} = F(A_l)$$

$$w_{0i} = b \sim N(0, \sigma_b^2)$$

# Dynamics of Activity: law of large numbers

$$\tilde{x}_\alpha = \varphi\left(\sum w_{\alpha k} x_k + b_\alpha\right) = \varphi(u_\alpha) : \tilde{x} = \phi(Wx)$$

$$u_\alpha \sim N(0, A)$$



$$x \rightarrow \tilde{x} : A \rightarrow \tilde{A}$$

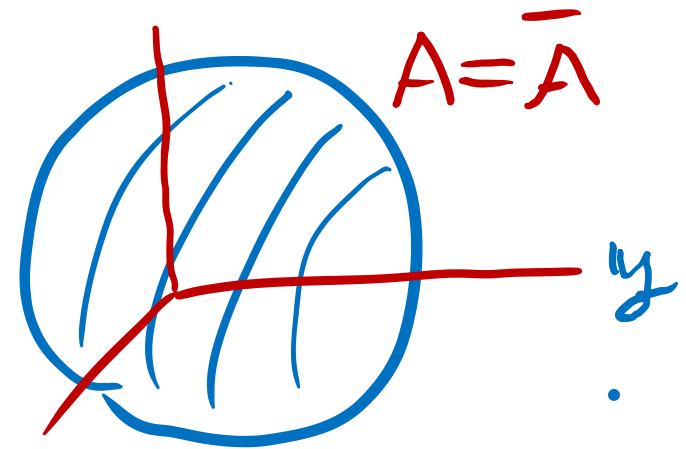
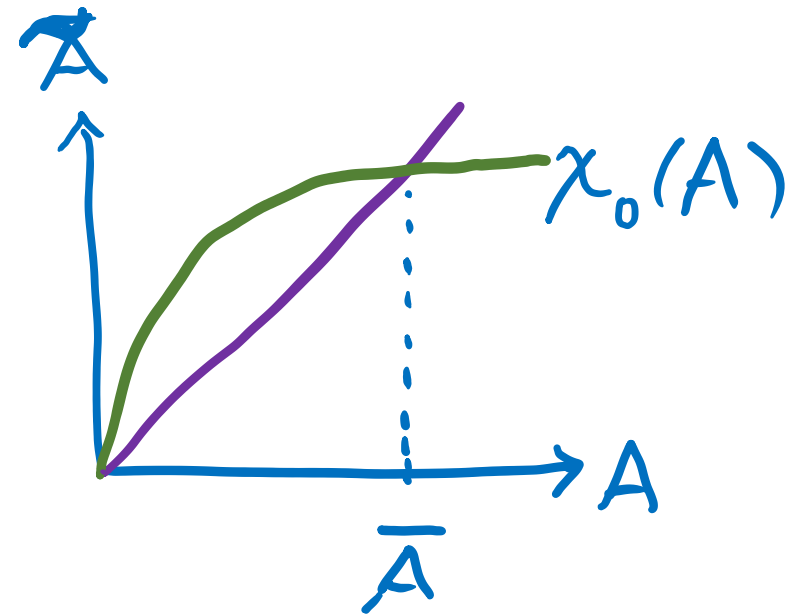
$$\tilde{A} = \frac{1}{n_{l+1}} \sum (\tilde{x}_\alpha)^2 = E[\varphi(u_\alpha)^2] = \chi_0(A)$$

$$\chi_0(A) = \int \varphi^2(\sqrt{A}v) Dv \quad v \sim N(0,1)$$

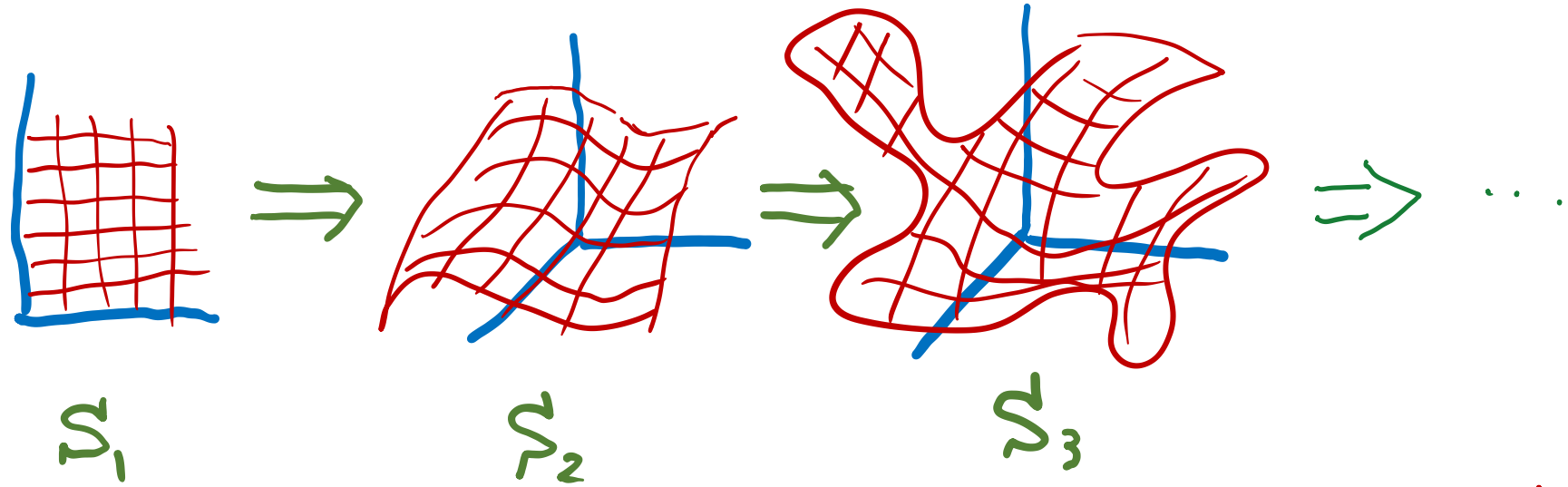
$$\chi_0'(0) > 1$$

$$\bar{A} = \chi_0(\bar{A})$$

$$\sum x_i^2 \rightarrow \text{converge}$$

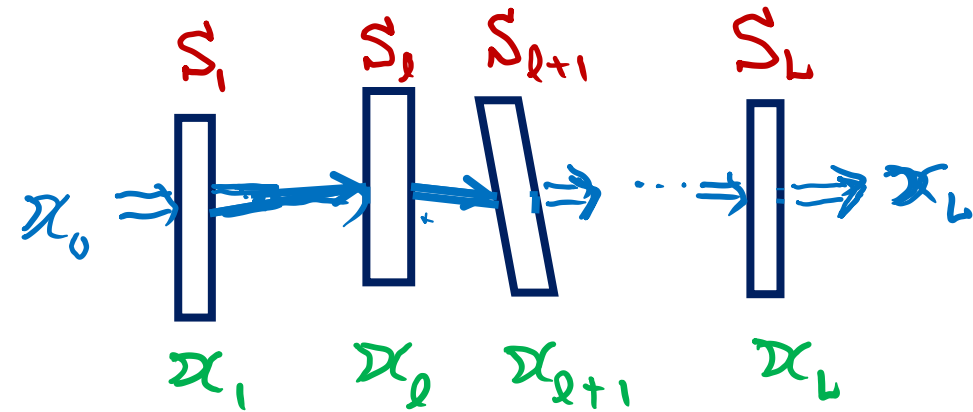


# 引き戻し計量 (リーマン計量・距離・曲率)

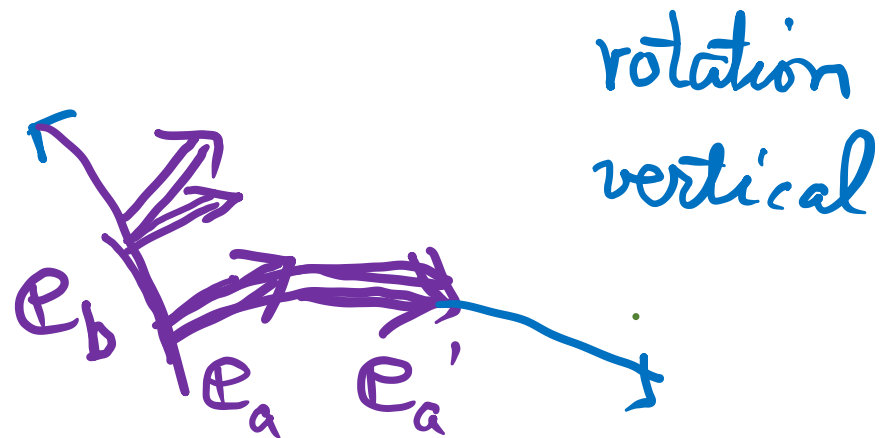
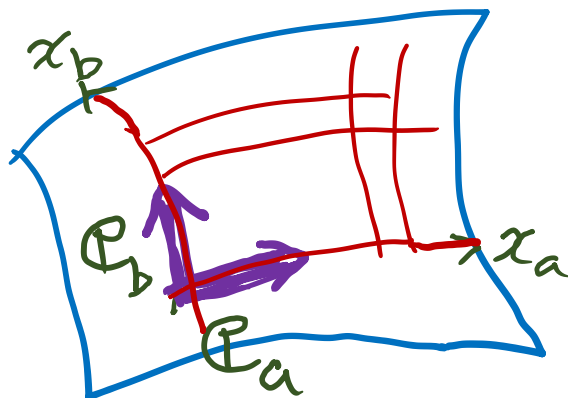


$$ds^2 = \sum g^l_{ab} dx^a dx^b = \frac{1}{n_l} d\mathbf{x}^l \cdot d\mathbf{x}^l$$

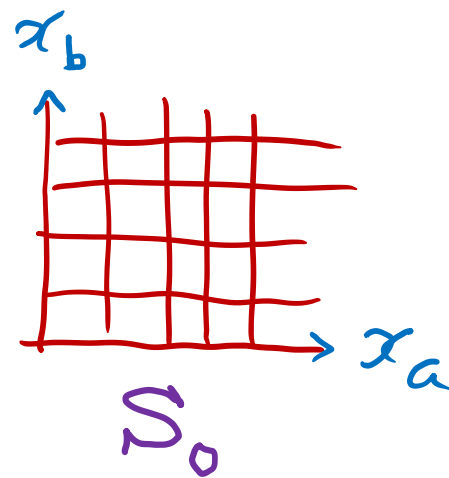
$$g^l_{ab} = \mathbf{e}^l_a \cdot \mathbf{e}^l_b$$



# 曲率



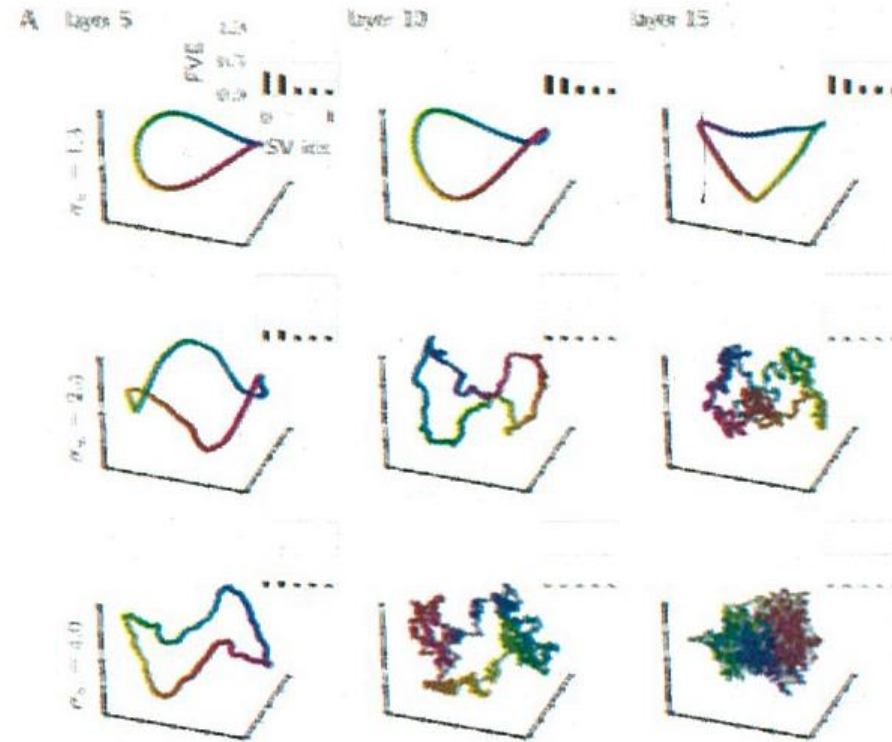
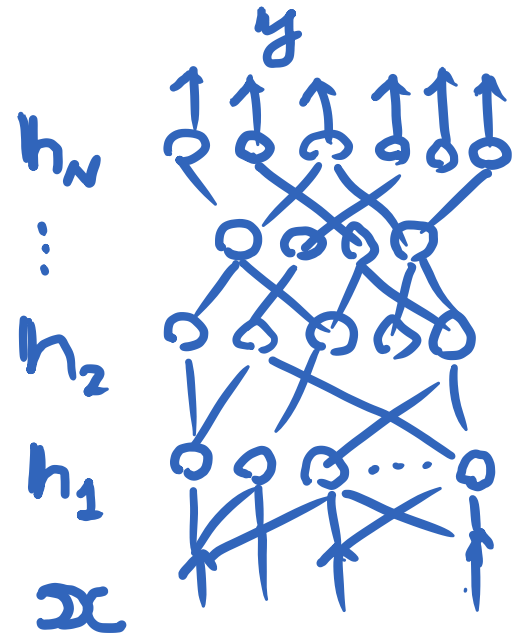
$$H_{abi}^{\ell} = \nabla_a e_b^{\ell}$$





# Poole et al (2016)

## Random deep neural networks



# Basis vectors

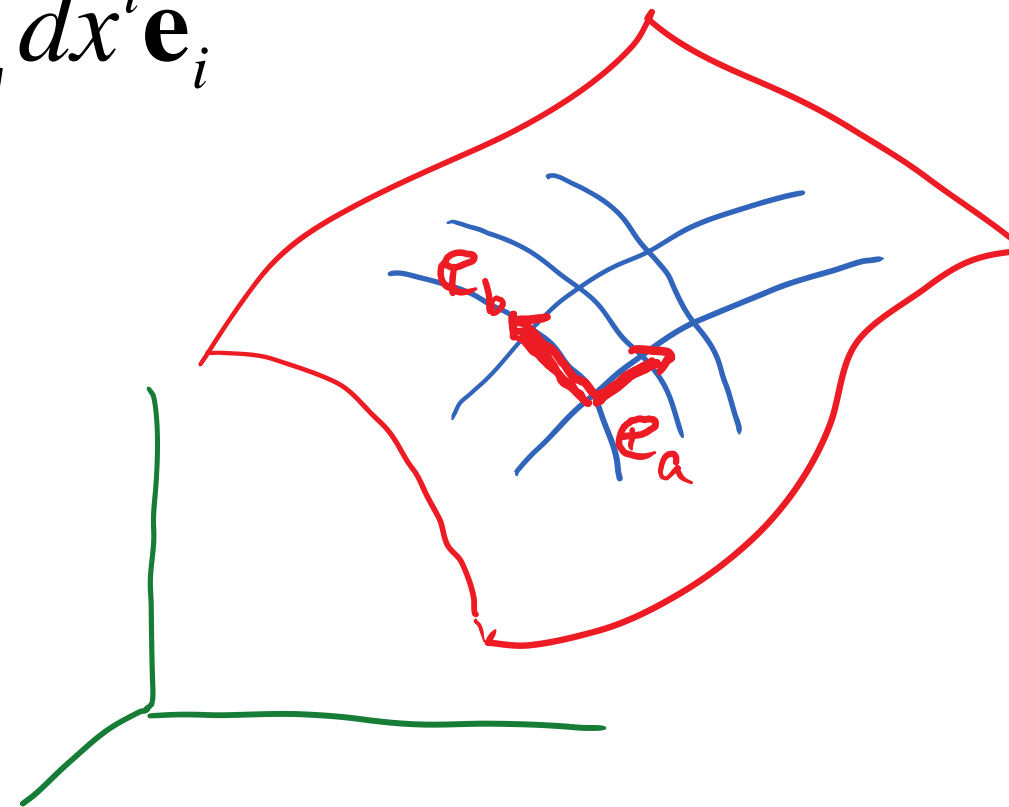
$$dx_{i_l} = \sum \varphi'(u_{i_l}) W_{i_{l-1}}^{i_l} dx_{i_{l-1}} = \sum B_{i_{l-1}}^{i_l} dx_{i_{l-1}}$$

$$d\mathbf{x} = B d\mathbf{x} = B \dots B d\mathbf{x}$$

$$d\mathbf{x} = \sum dx^i \mathbf{e}_i$$

$$B_{i_{l-1}}^{i_l} = \varphi'(u_{i_l}) W_{i_{l-1}}^{i_l}$$

$$\mathbf{e}_a = B \mathbf{e}_a = B \dots B \mathbf{e}_a$$



# リーマン計量の力学

$$\tilde{y}_\alpha = \varphi\left(\sum w_{\alpha k} y_k + b_\alpha\right) = \varphi(u_\alpha)$$

$$d\tilde{y}_\alpha = \sum B_k^\alpha dy_k \quad \tilde{\mathbf{e}}_a = B\mathbf{e}_a$$

$$ds^2 = \sum g_{ij} dy^i dy^j = \langle d\mathbf{y}, d\mathbf{y} \rangle$$

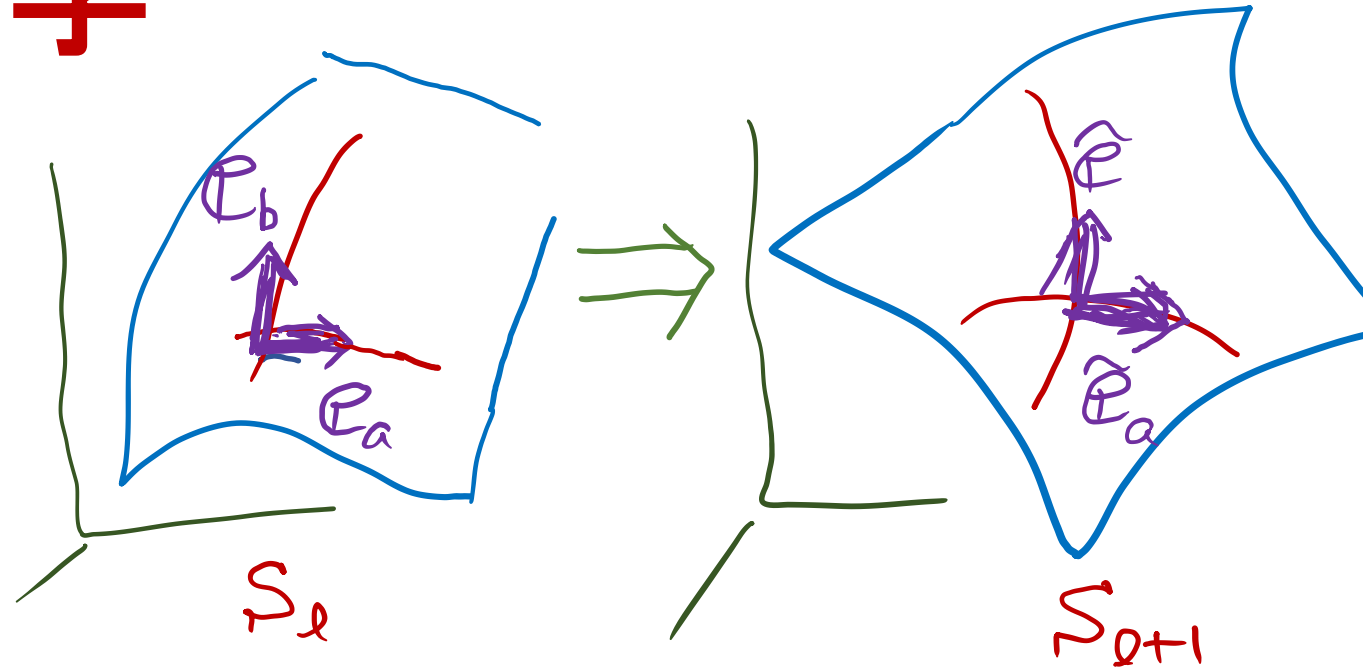
$$B = (B_k^\alpha) = (\varphi'(u_\alpha) w_k^\alpha)$$

$$\langle \tilde{\mathbf{e}}_\alpha, \tilde{\mathbf{e}}_\beta \rangle = \sum B_k^\alpha B_j^\alpha \langle \mathbf{e}_k, \mathbf{e}_j \rangle = \chi_1 \delta_\alpha^k \delta_\beta^j g_{jk}$$

$$E[(\varphi'(u_\alpha))^2 w_k^\alpha w_j^\alpha] = E[(\varphi'(u_\alpha))^2] E[w_k^\alpha w_j^\alpha]$$

平均場近似

$$\chi_1(A) = \int \sigma^2 \{ \varphi'(\sqrt{A}v) \}^2 Dv = \frac{1}{2\pi} \frac{\sigma^2 A + \sigma_b^2}{\sqrt{1 + 2(\sigma^2 A + \sigma_b^2)}}$$



# Metric

Law of large numbers

$${}^l g_{ab} = \left\langle {}^l \mathbf{e}_a, {}^l \mathbf{e}_b \right\rangle = BB \quad {}^{l-1} g_{ab}$$

$$ds^2 = \sum {}^l g_{ab} d {}^l x_a d {}^l x_b$$

$$BB = \sum_{i_l} W_{i_{l-1}}^{i_l} W_{i'_{l-1}}^{i_l} \varphi'(u_{i_l})^2 \approx \sigma_l^2 E[\varphi'^2] \delta_{i_{l-1} i'_{l-1}}$$

$$\chi_1 = \sigma_l^2 E\left[\varphi'(u_{i_l})^2\right]$$

# Metric

Law of large numbers

$${}^l g_{ab} = \left\langle {}^l \mathbf{e}_a, {}^l \mathbf{e}_b \right\rangle = BB {}^{l-1} g_{ab}$$

$$ds^2 = \sum {}^l g_{ab} d {}^l x_a d {}^l x_b$$

$$BB = \sum_{i_l} W_{i_{l-1}}^{i_l} W_{i'_{l-1}}^{i_l} \varphi'(u_{i_l})^2 \approx \sigma_l^2 E[\varphi'^2] \delta_{i_{l-1} i'_{l-1}}$$

$$\chi_1 = \sigma_l^2 E\left[\varphi'(u_{i_l})^2\right]$$

$$\tilde{g}_{ab} = \chi_1(A) g_{ab}$$

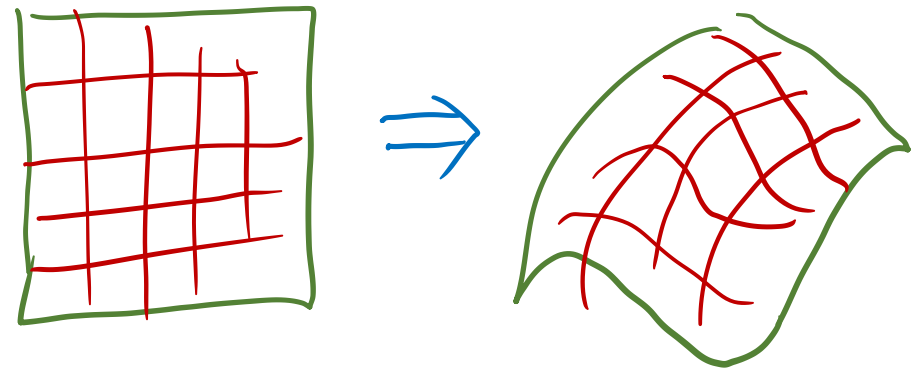
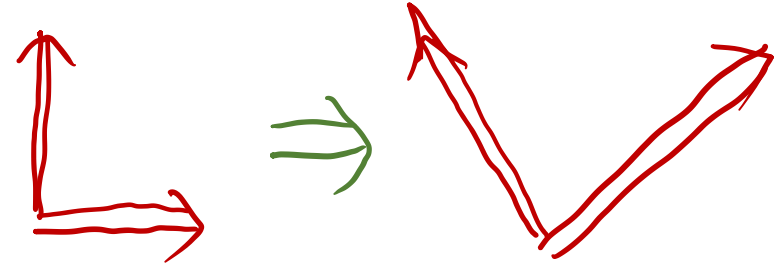
conformal transformation!

$$\bar{\chi}_1 = \bar{\chi}_1(\bar{A}) > 1:$$

拡大 (カオス、Lyapunov指数)

$$\Rightarrow g^l_{ab} = \prod \chi_1(A^s) \delta_{ab}$$

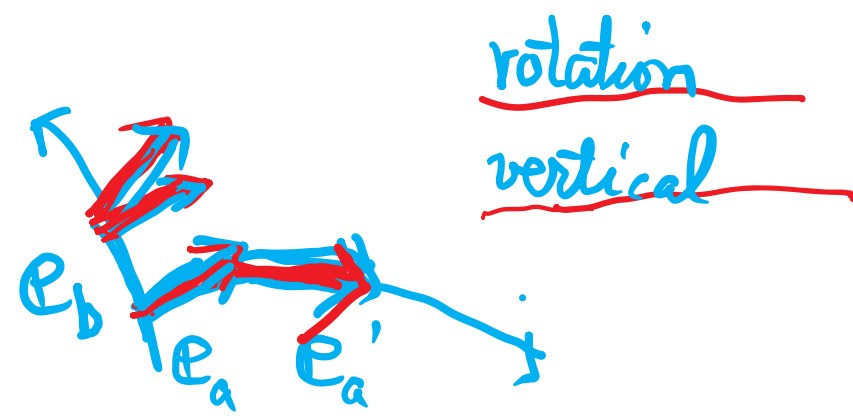
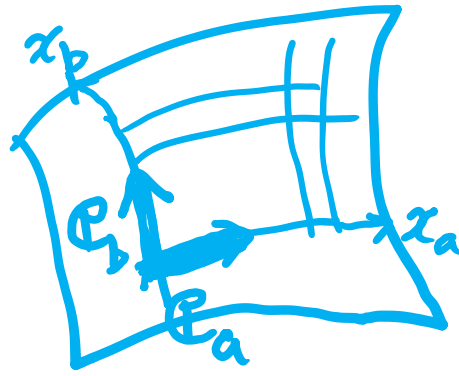
回転, 拡大・縮小



$${}^l g_{ab}(x) = \left( \prod \lambda_1(x) \right) g_{ab}(x)$$

**conformal geometry**

# 曲率の力学



$$\tilde{H}_{ab}^{\alpha} = \nabla_a \tilde{\mathbf{e}}_b^{\alpha} = \partial_a \partial_b \tilde{y}^{\alpha}$$

$$= \varphi''(u_{\alpha})(\mathbf{w} \cdot \mathbf{e}_a)(\mathbf{w} \cdot \mathbf{e}_b) + \varphi'(\mathbf{w} \cdot \partial_a \mathbf{e}_b)$$

$$\tilde{\mathbf{H}}_{ab} = \mathbf{H}_{ab}^{\perp} + \mathbf{H}_{ab}^{\square}$$

Euler-Schouten曲率  
Affine connection

$$\tilde{H}_{ab}^2 = |\tilde{\mathbf{H}}_{ab}|^2$$



# curvature & distortion

$$\mathbf{H}_{ab} = \nabla_a^l \mathbf{e}_b = \nabla_a \left( \mathbf{B}^{\ l-1} \mathbf{e}_b \right) = \mathbf{B} \nabla_a^{\ l-1} \mathbf{e}_b + (\nabla_a \mathbf{B})^{\ l-1} \mathbf{e}_b$$

$$\left| \mathbf{H}_{ab}^{\ l} \right|^2 = \chi_1 \left| \mathbf{H}_{ab}^{\ l-1} \right|^2 + \frac{1}{n \chi_1^2} (1 + 2\delta_{ab}) \chi_2$$

$$\chi_2 = \sigma^2 E \left[ \varphi''(u)^2 \right]$$

$$\chi_2(A) = \int \varphi''(\sqrt{A}v)^2 Dv$$

$$H_{ab}^{l+1} = \frac{1}{n\chi_1^2} \chi_2(A)(2\delta_{ab} + 1) + \chi_1(A) H_{ab}^{l^2}$$

$$\chi_1 > 1$$

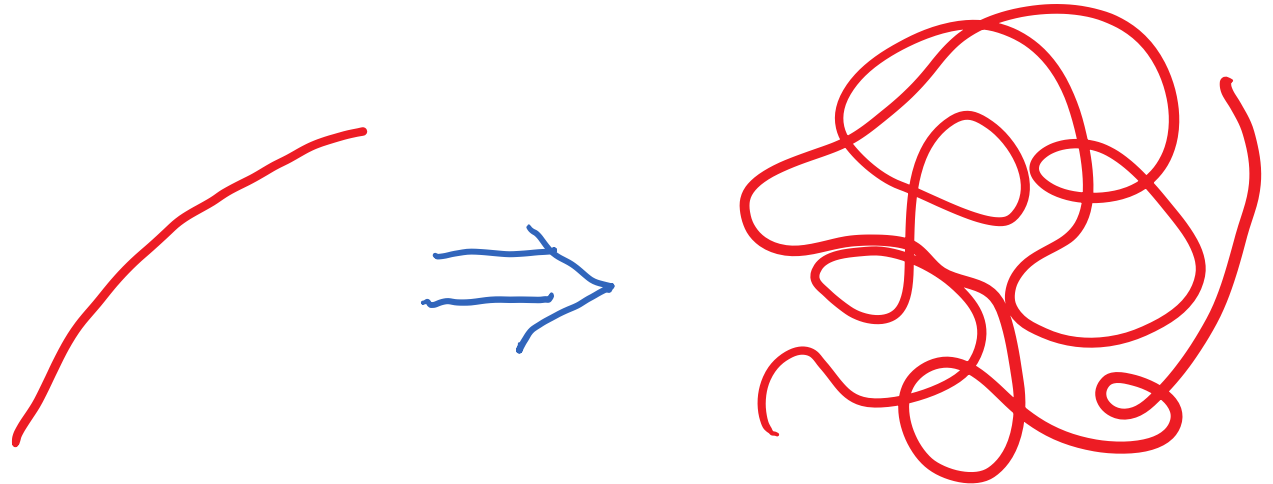
exponential expansion! creation is small!

# scalar curvature & distortion

$$\gamma^l = \frac{1}{\chi_1} \gamma^{l-1} + \frac{3}{n} \frac{\chi_2^2}{\chi_1^2}$$

$$\gamma^2 = H_{ab}^i H_{cd}^j g^{ac} g^{bd} \delta_{ij}$$

$$\gamma^l \rightarrow \frac{3\chi_2}{n\chi_1(\chi_1 - 1)} : \rightarrow \infty, \chi_1 \leq 1$$

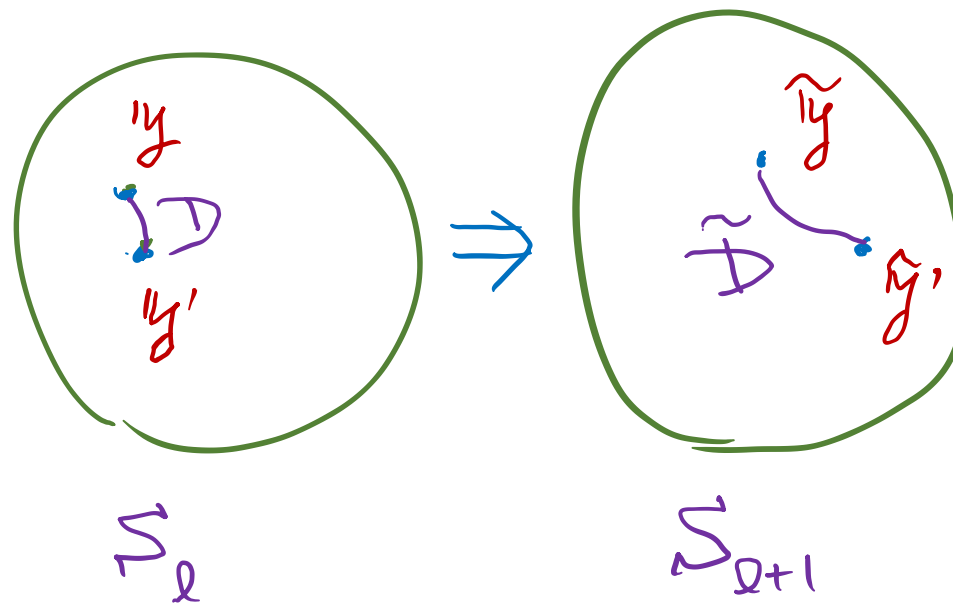


# 距離法則 (Amari, 1974)

$$D(x, x') = \frac{1}{n} \sum (x_i - x_i')^2$$

$$C(x, x') = \frac{1}{n} x \cdot x' = \sum x_i x_i'$$

$$D = A + A' - 2C$$



# Dynamics of Distance (Amari, 1974)

$$D(x, x') = \frac{1}{n} \sum (x_i - x_i')^2$$

$$C(x, x') = \frac{1}{n} x \cdot x' = \sum x_i x_i'$$

$$D = A + A' - 2C$$

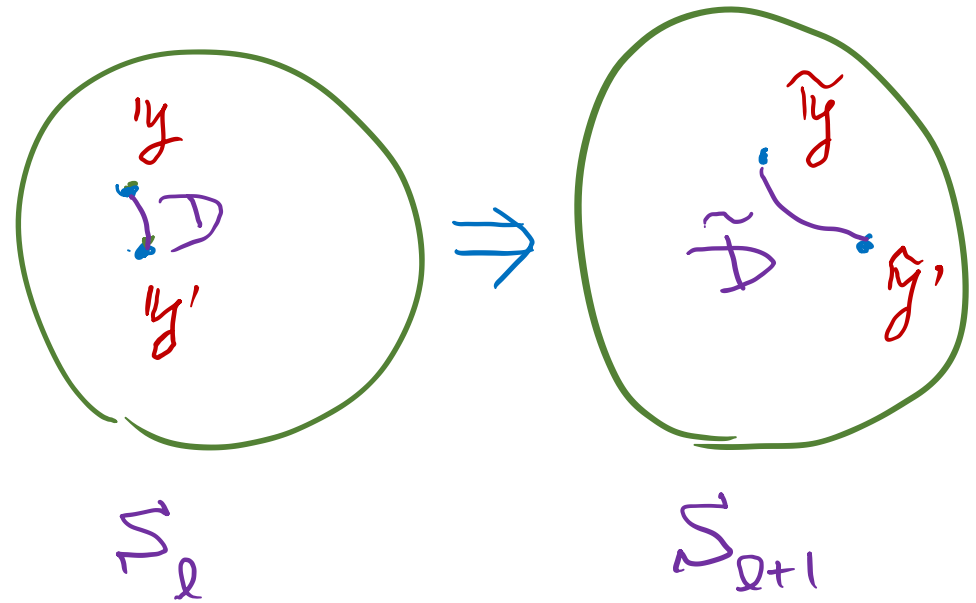
$$u_\alpha = \sum w_{\alpha k} y_k$$

$$\sim N(0, V)$$

$$u'_\alpha = \sum w_{\alpha k} y'_k$$

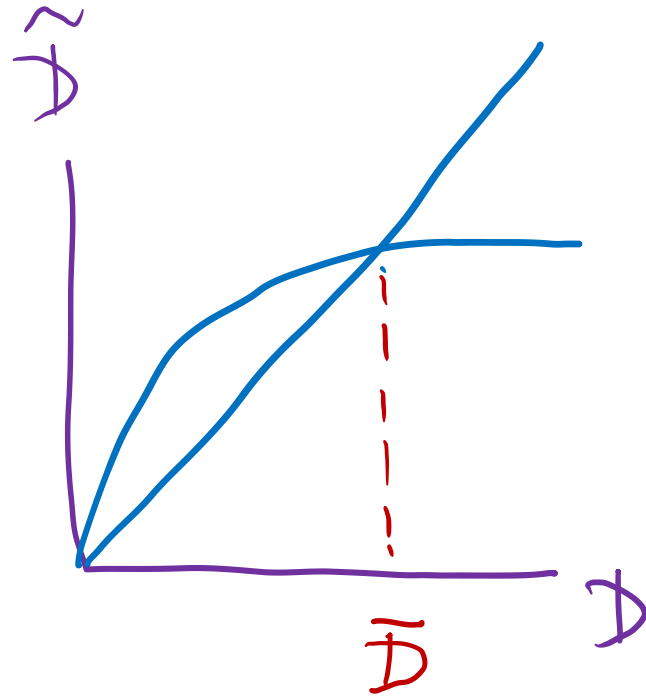
$$V = \begin{bmatrix} A & C \\ C & A' \end{bmatrix}$$

$$\tilde{C} = E[\varphi(\sqrt{A-C}\varepsilon + \sqrt{C}\nu)\varphi(\sqrt{A'-C}\varepsilon + \sqrt{C}\nu)]$$



$$D_{l+1} = K(D_l)$$

$$\frac{d\tilde{D}}{dD} \Big|_{D=0} = \chi_1 > 1$$



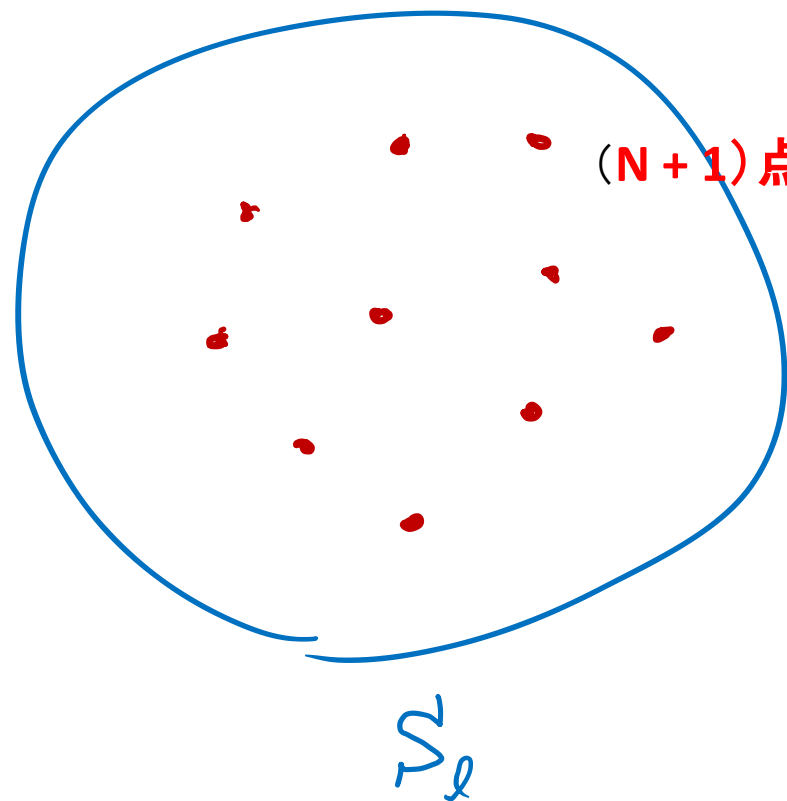
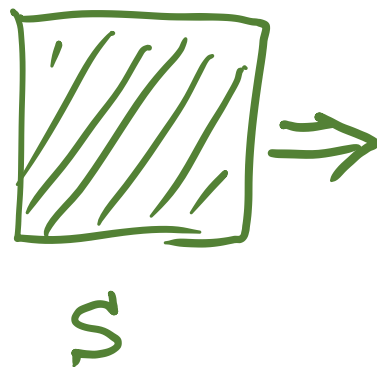
**本当か!**  $n \rightarrow \infty; l \rightarrow \infty$

*equidistance*

$$D(\mathbf{x}_l, \mathbf{x}'_l) \rightarrow \bar{D}$$

$$\bar{D} = \xi(\bar{D})$$

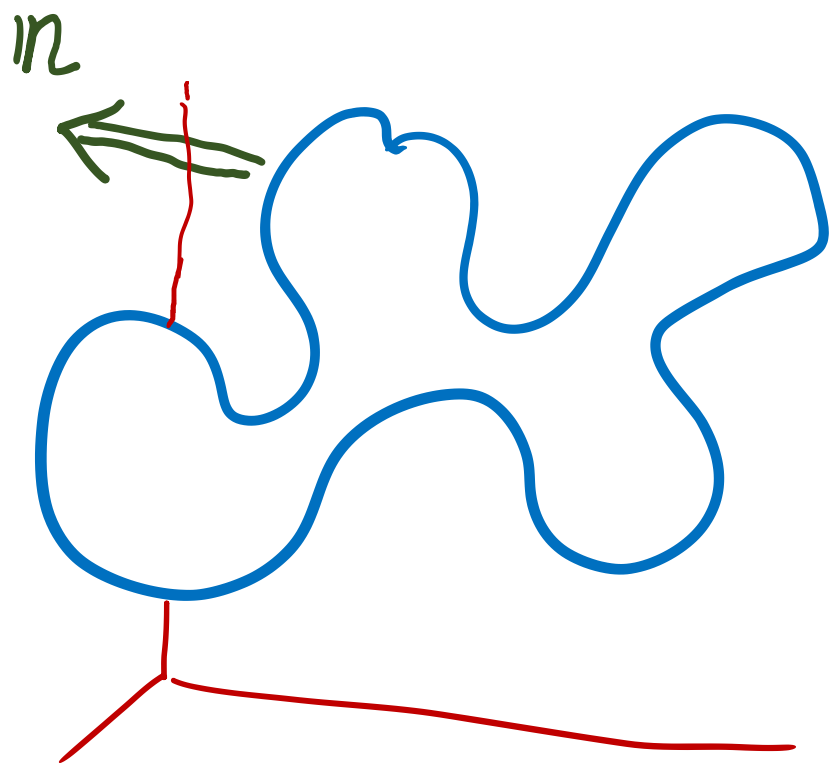
等距離



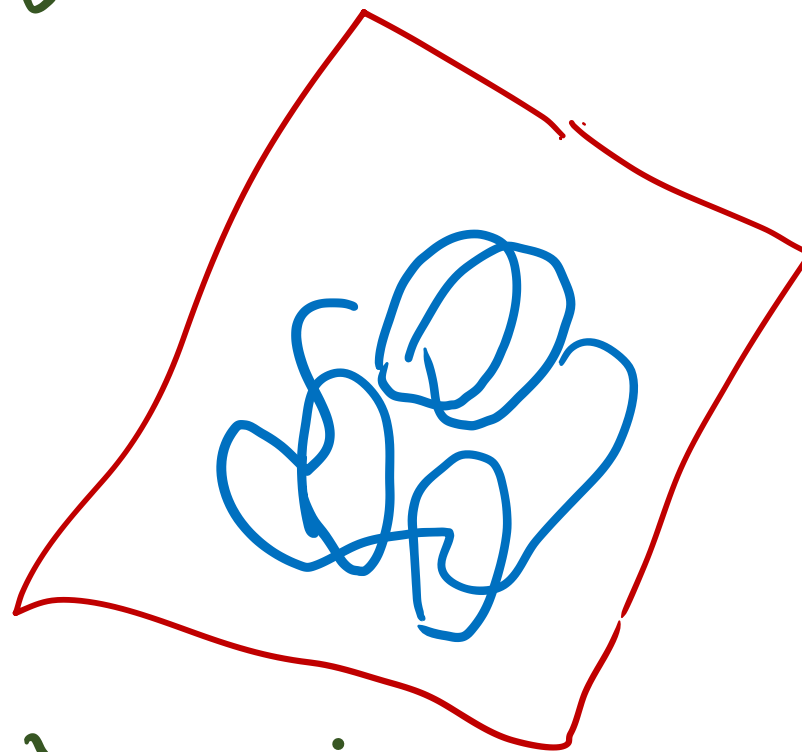
フラストレーション  
フラクタル

$$n_{l+1} < n_l$$

# 次元の縮小



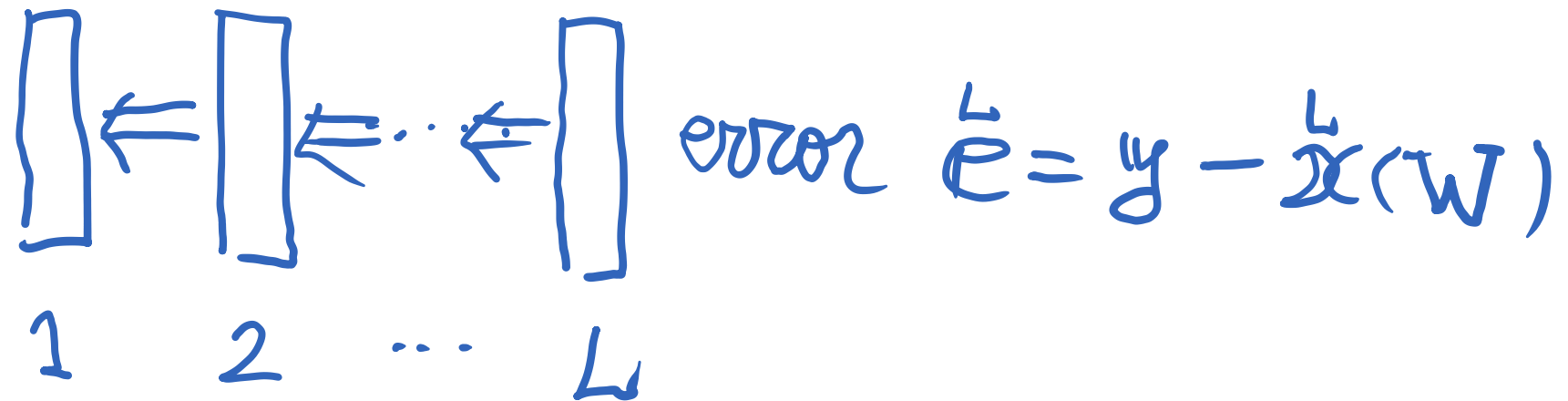
$$Wn=0$$



$$\bar{y} = \varphi(Wx)$$



# Fisher 情報行列と逆向き情報伝播



$$l(x, W) = \frac{1}{2} |y - \varphi(x; W)|^2 = |e(x, y)|^2$$

# 確率 model : 深層回路の多様体

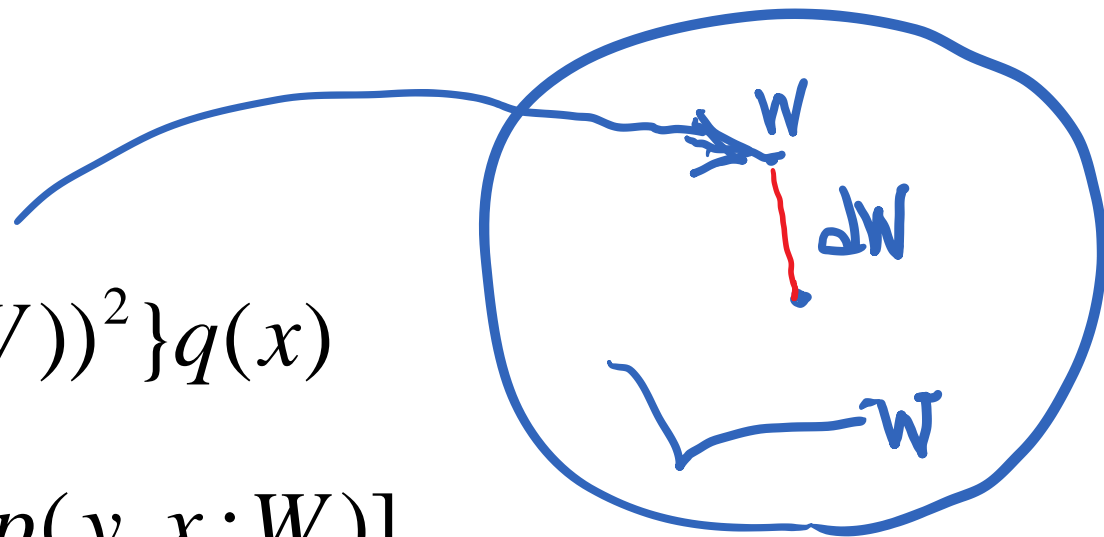
$$y = \varphi(u) + \varepsilon; \quad \varepsilon \sim N(0,1)$$

$$p(y, x : W) = c \exp\left\{-\frac{1}{2}(y - \varphi(x; W))^2\right\} q(x)$$

$$G = E_x[\nabla_W \log p(y, x : W) \nabla_W \log p(y, x : W)]$$

$$\underline{ds^2 = dW G dW}$$

↑ Fisher information



Riemannian

# Natural Gradient

$$\max \quad dl = l(\boldsymbol{\theta} + d\boldsymbol{\theta}) - l(\boldsymbol{\theta})$$

$$|d\boldsymbol{\theta}|^2 = \varepsilon \quad \text{KL}[p(\mathbf{x}, \boldsymbol{\theta}) : p(\mathbf{x}, \boldsymbol{\theta} + d\boldsymbol{\theta})] = \varepsilon$$

$$\nabla l = G^{-1}(\boldsymbol{\theta}) \nabla l$$

$$\Delta \boldsymbol{\theta}_t = -\eta_t \hat{\nabla} l(x_t, y_t; \boldsymbol{\theta}_t)$$

# Fisher information

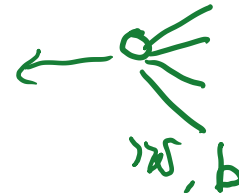
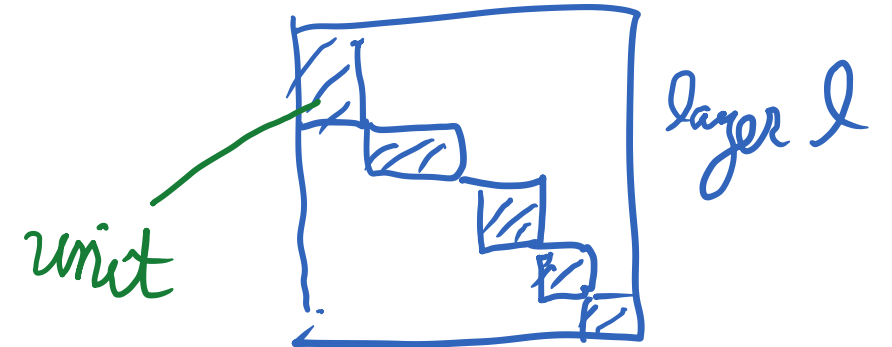
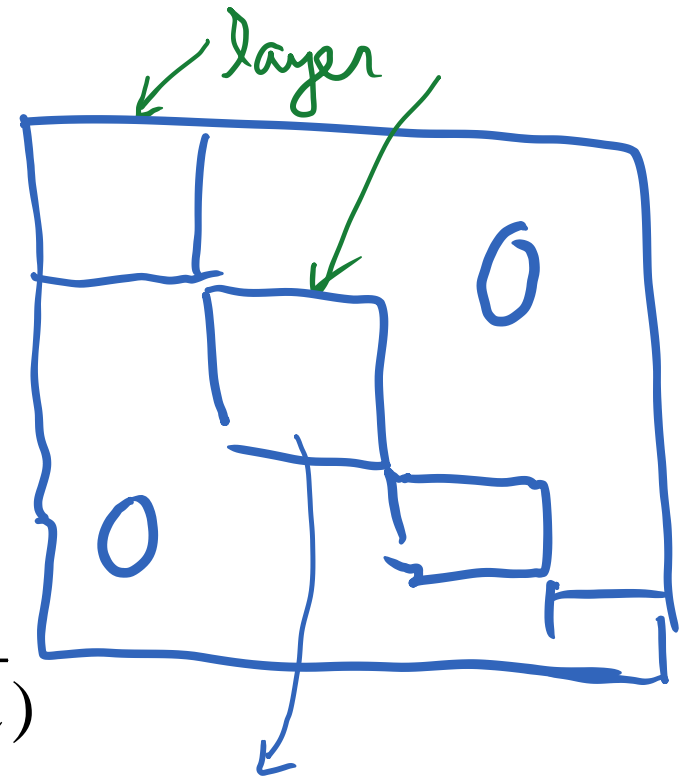
$$G = E_x \left[ \begin{array}{cc} \frac{\partial \phi}{\partial W_m} & \frac{\partial \phi}{\partial W_l} \end{array} \right]$$

$$\frac{\partial \phi^l}{\partial W_m} = \underbrace{\phi' W}_{\text{B}} \frac{\partial \phi^{l-1}}{\partial W_m} = B \frac{\partial \phi^{l-1}}{\partial W_m} = \underbrace{BB \dots B}_{\text{B}} \frac{\partial \phi^{m+1}}{\partial W_m}$$

$$G(W_l, W_m) = \prod \chi_1 E_x \left[ \begin{array}{cc} \phi' \left( \begin{array}{c} l \\ \mathbf{w}_i \end{array} \right)^2 & \begin{array}{cc} l-1 & l-1 \\ \mathbf{x} & \mathbf{x} \end{array} \end{array} \right] + O_p(1/\sqrt{n})$$

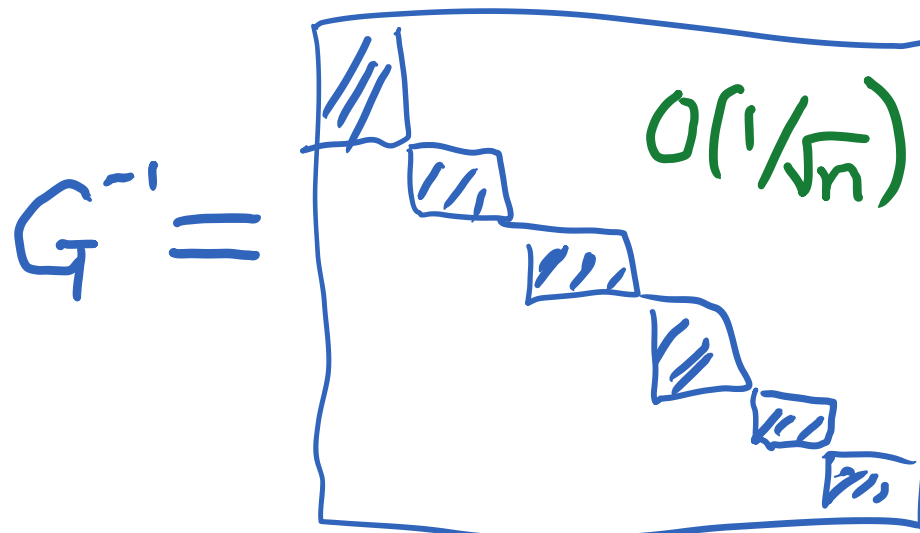
$$G(W_l, W_m) = 0 \sim O_p(1/\sqrt{n}), \quad l \neq m$$

$$G \left( \begin{array}{c} l \\ \mathbf{w}_i \end{array}, \begin{array}{c} l \\ \mathbf{w}_j \end{array} \right) = 0 \sim O_p(1/\sqrt{n}), \quad i \neq j$$

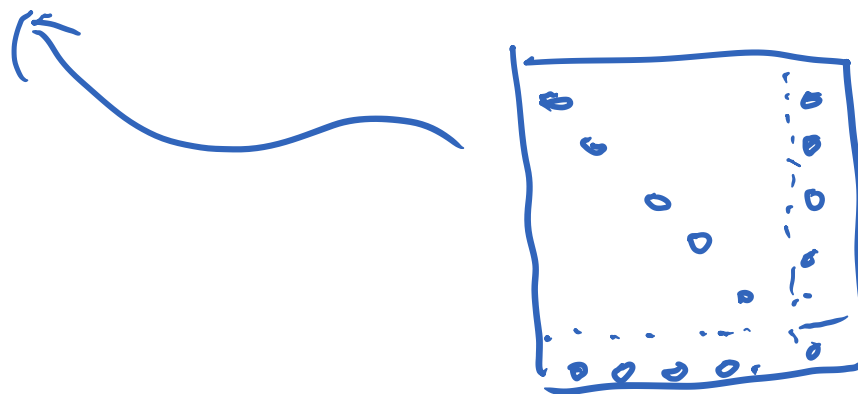


# Unitwise natural gradient

$$\Delta W = -\eta G^{-1} \nabla_W l$$

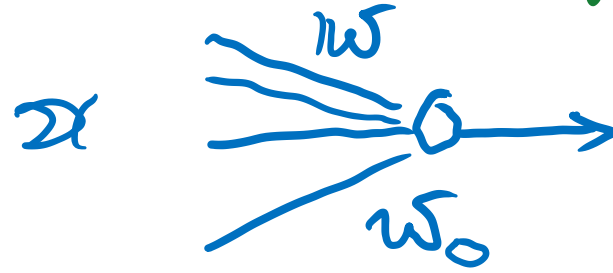


Y. Ollivier; Marceau-Caron



# 一個のニューロンの情報行列

$$\tilde{\mathbf{w}} = (w, w_0)$$



$$y = \varphi(\mathbf{w} \cdot \mathbf{x} + w_0)$$

$$G = \mathbb{E}_x [\partial_{\mathbf{w}} \varphi \partial_{\mathbf{w}} \varphi] = \mathbb{E}_x [(\varphi')^2 \mathbf{x}\mathbf{x}]$$

$$G(\tilde{\mathbf{w}}, \tilde{\mathbf{w}})$$

$\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\} \Rightarrow \{\mathbf{e}_1^*, \mathbf{e}_2^*, \dots, \mathbf{e}_n^*\}$  : ortho-normal basis

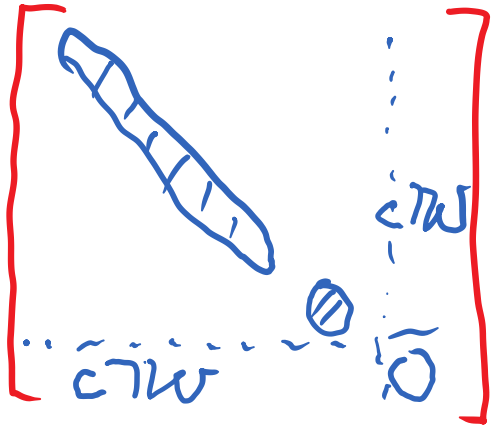
$$\mathbf{e}_n^* = \frac{\mathbf{w}}{w}$$

Input  $x$ : independent and identically distributed, 0-mean

$$G = A\mathbf{I} + \frac{B}{w^2} \mathbf{w}\mathbf{w} + \frac{C}{w} (\mathbf{w}\mathbf{b} + \mathbf{b}\mathbf{w})$$

$$\mathbf{b} = [0 \ 0 \ \dots \ 0 \ 1]'$$

$G^{-1}$  : similar form

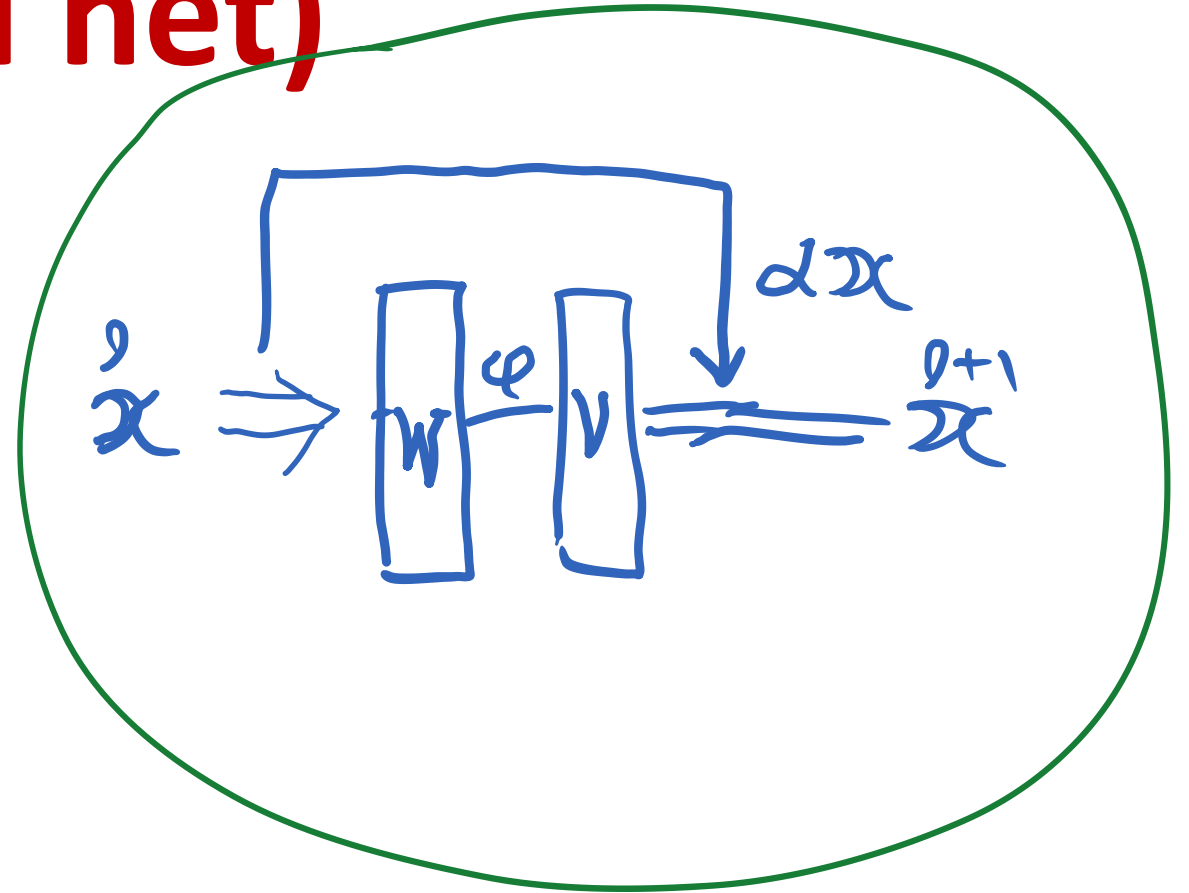


$A, B, C$   
 $(w, b)$

# Resnet (residual net)

$$\mathbf{x}^l = V \varphi \left( W \mathbf{x}^{l-1} \right) + \alpha \mathbf{x}^{l-1}$$

$$\chi_1 \rightarrow \sigma_v^2 \chi_1 + \alpha^2$$



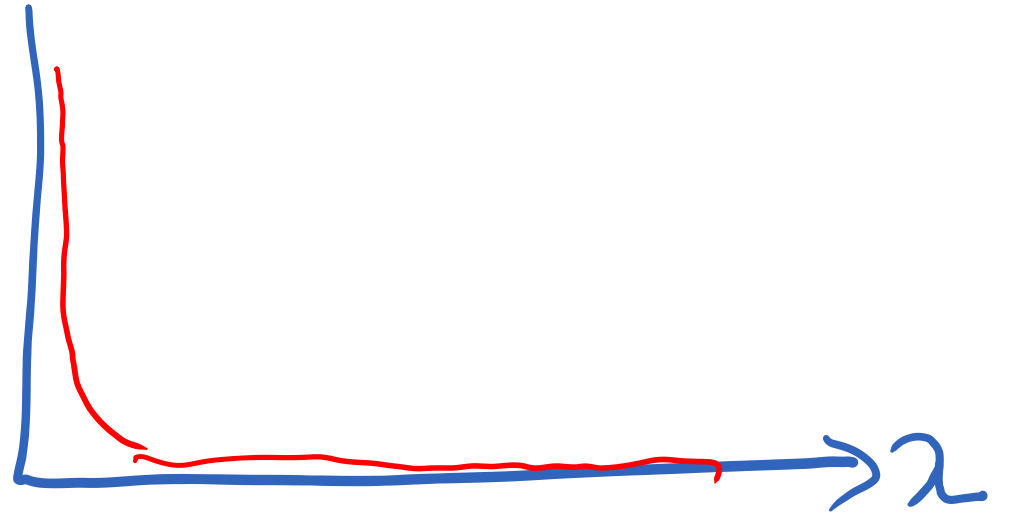


# Karakida theory

eigenvalues of  $G$

$$\frac{1}{P} \sum \lambda_i = \frac{1}{n}, \quad \frac{1}{P} \sum \lambda_i^2 = O(1)$$

distorted Riemannian metric



# Wasserstein Distance の情報幾何

**Shun-ichi Amari**

**RIKEN Brain Science Institute**

**R. Karakida. M. Oizumi**

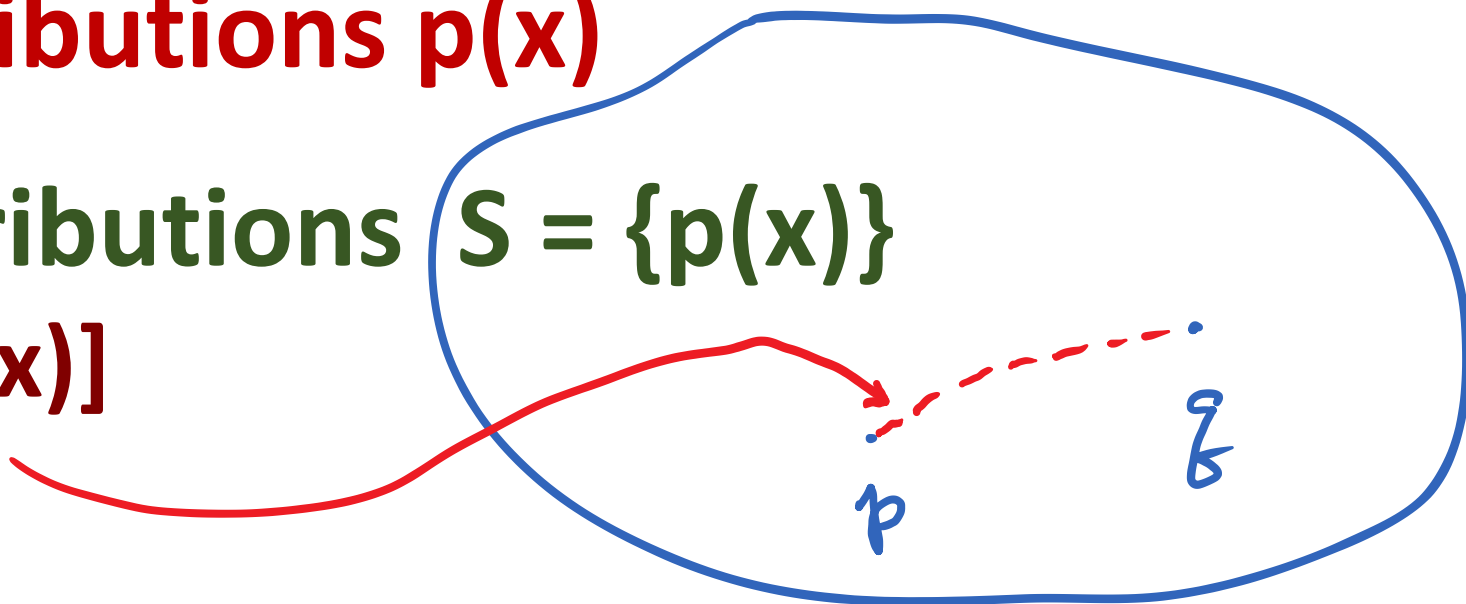
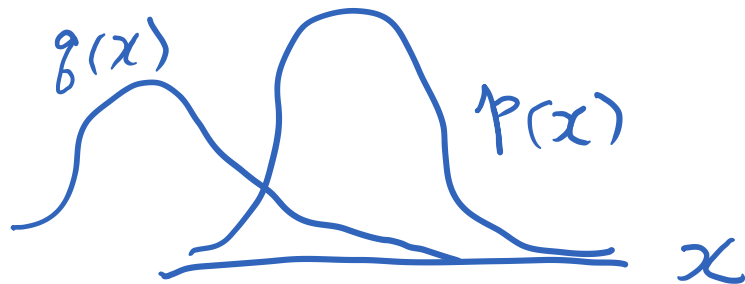
**Base space  $X$ ;  $X$ 上のパターン:  $p(x)$**

**picture  $X = (x, y)$ ; Boltzmann machine  $X = \{0, 1\}^n$   $\mathcal{X}$**

**probability distributions  $p(x)$**

**Geometry of Distributions  $S = \{p(x)\}$**

**distance  $D[p(x): q(x)]$**



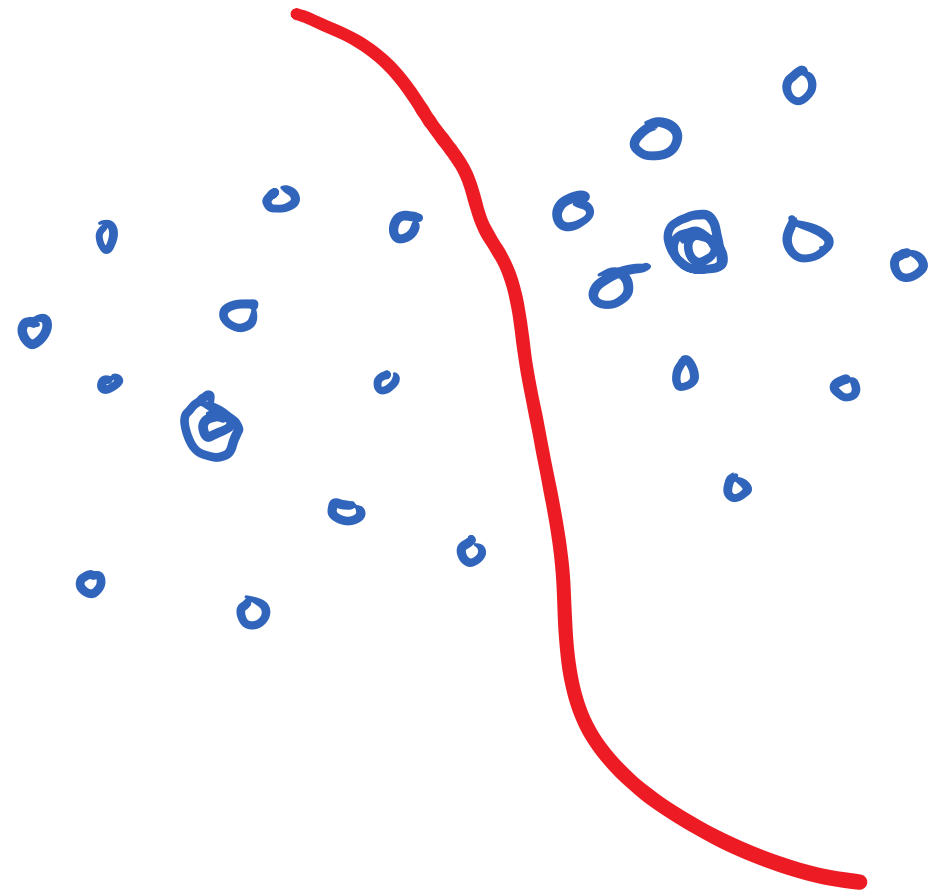
# ダイバージェンス

クラスタリング

Pat/パターン認識

機械学習

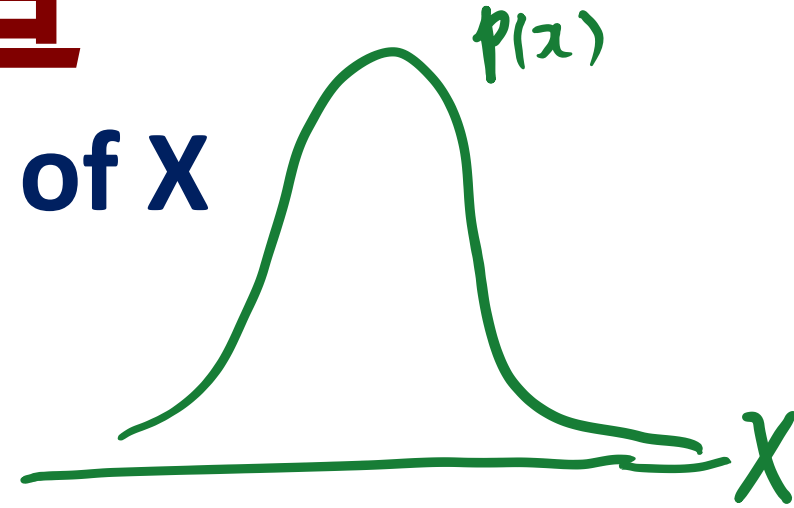
統計的推論



# 情報幾何と不変構造

invariant under transformations of  $X$

$$p(x) \sim p(y)$$



Fisher Information:

Affine connections:  $\alpha$ -connections

Duality: Dually coupled Riemannian manifold

# Wasserstein 距離

— Monge, Kantorovich

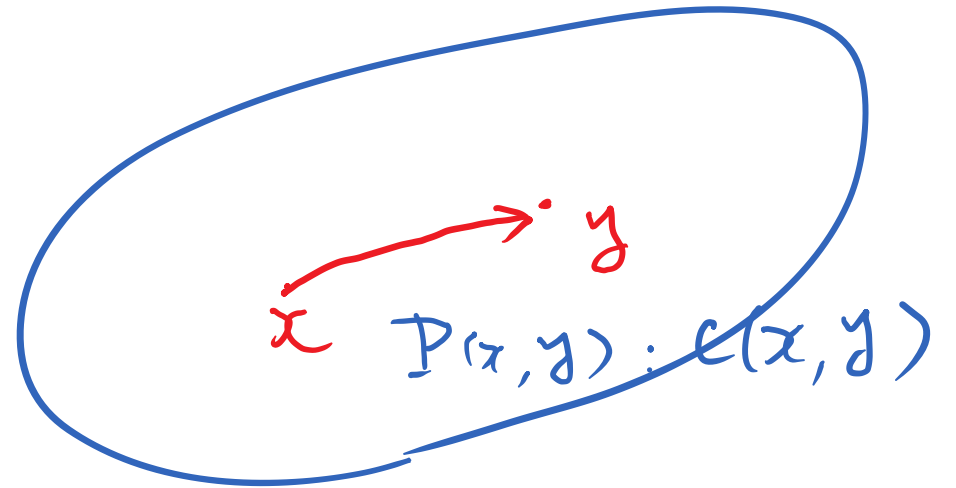
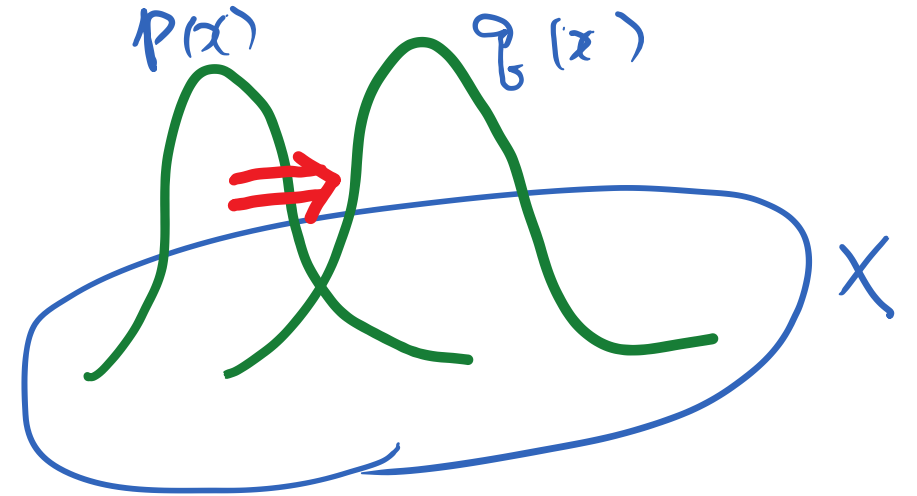
輸送問題  $p(x) \rightarrow q(x)$

cost  $c(x, y) = \text{metric over } X$

輸送計畫  $P(x, y)$

minimize  $\langle c, P \rangle =$

$$\int c(x, y) P(x, y) dx dy$$



# 線形計画問題

$$\text{minimize } \langle c, P \rangle = \int c(x, y) P(x, y) dx dy$$

under constraints

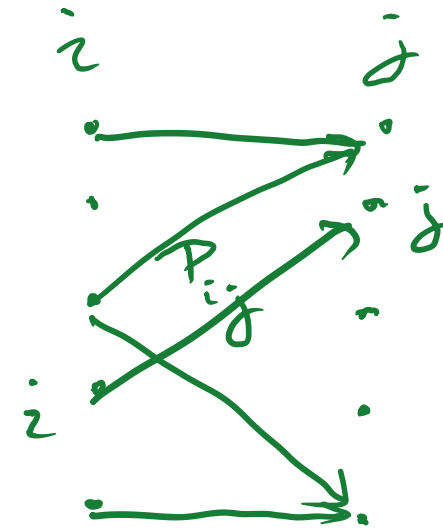
$$\int P(x, y) dy = p(x)$$

$$\int P(x, y) dx = q(y)$$

Discrete case :  $i \rightarrow j$

$$\text{minimize } \langle c, P \rangle = \sum c_{ij} P_{ij}$$

$$\text{constraints } \sum_j P_{ij} = p_i \quad \sum_i P_{ij} = q_j$$



# Entropy-正則化輸送問題

Marco

Cuturi

$$\min_{\mathbb{P}} F = \langle c, P \rangle - \lambda H[P(x, y)]$$

$\lambda \rightarrow 0$       **Wasserstein**

$\lambda \rightarrow \infty$       entropy term     $H[P(x, y)]$

$$P(x, y) = p(x)q(y) \quad \text{---} \quad \text{KL}[p(x): q(x)]$$

**information geometry**



# パターン $p(x)$ と $q(x)$ の距離?

KL-divergence

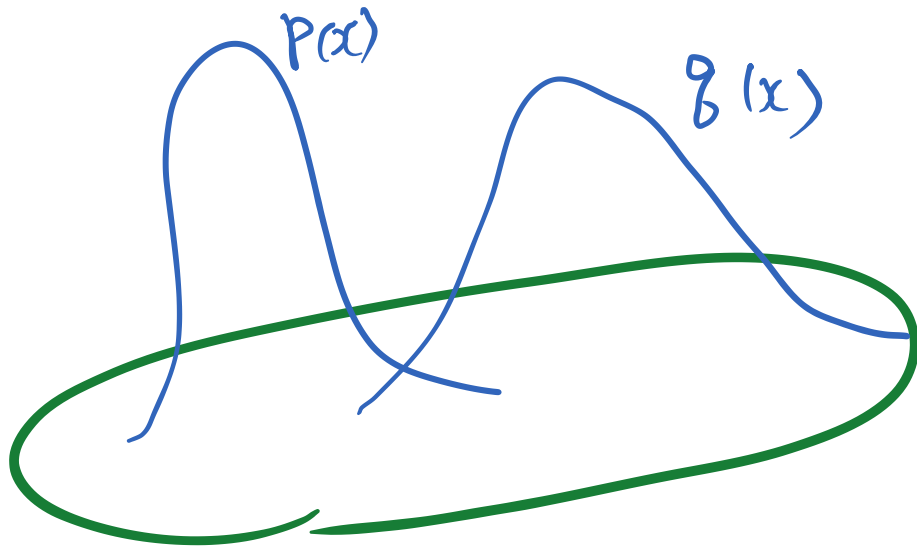
Hellinger distance

Wasserstein distance

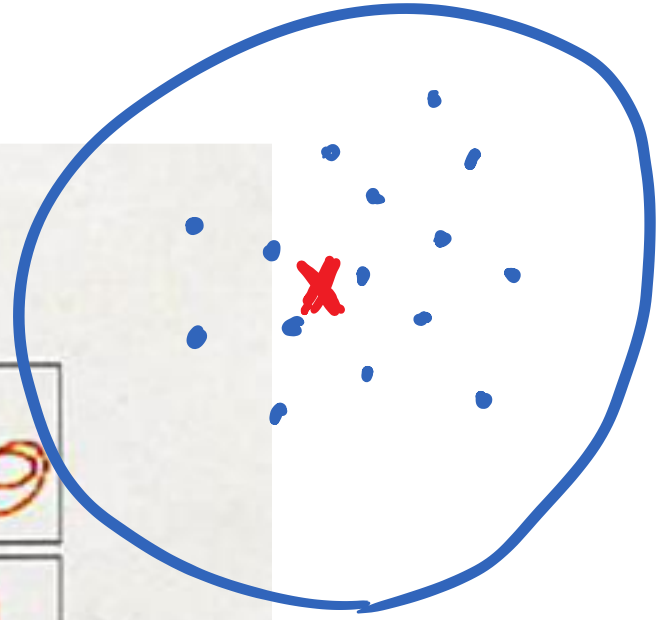
$$D_{KL} = \int p(x) \log \frac{p(x)}{q(x)} dx$$

$$D_{Hell} = \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx$$

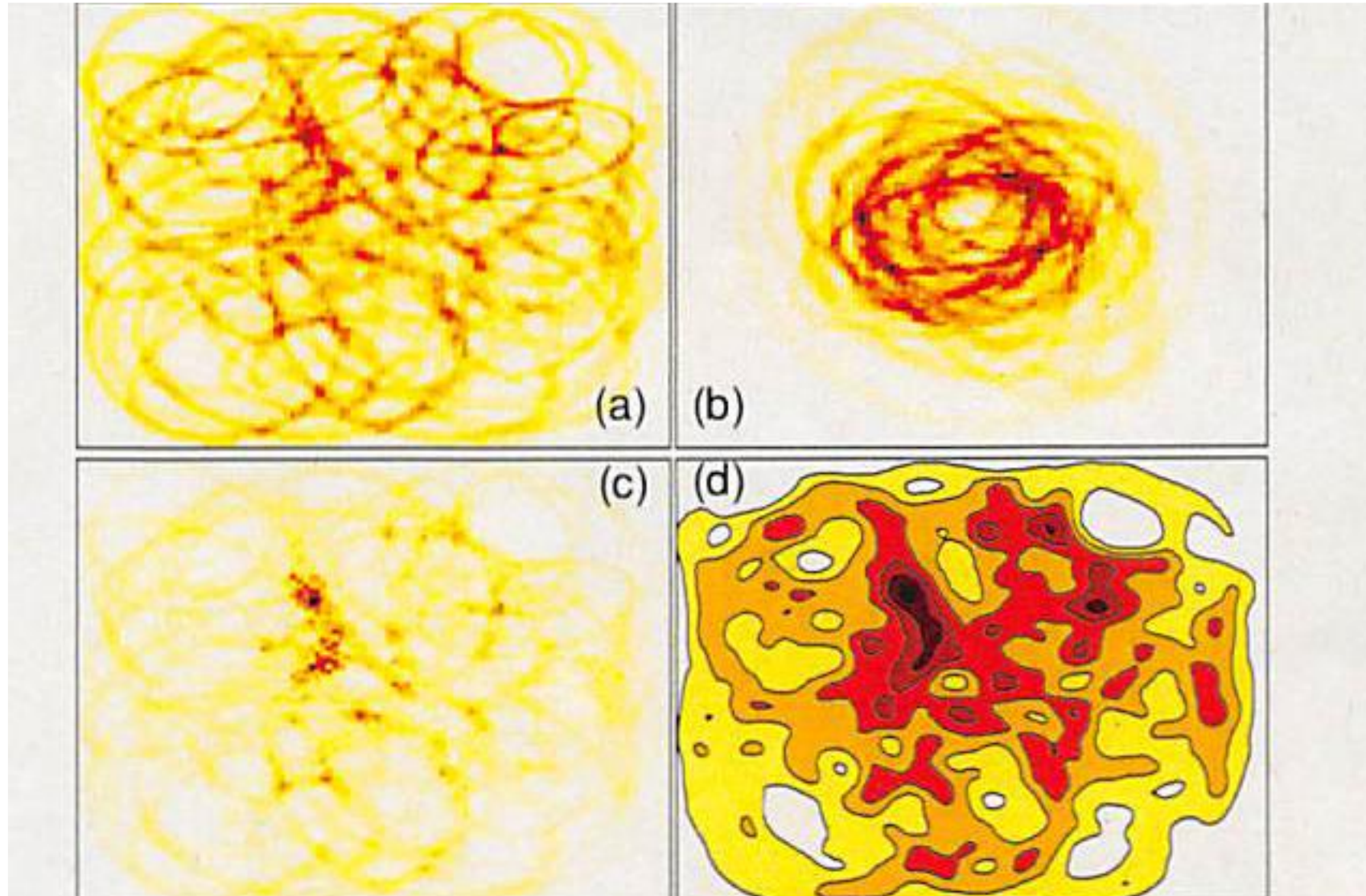
$$D_{Wass} = \int c(x, y) P(x, y) dx dy$$

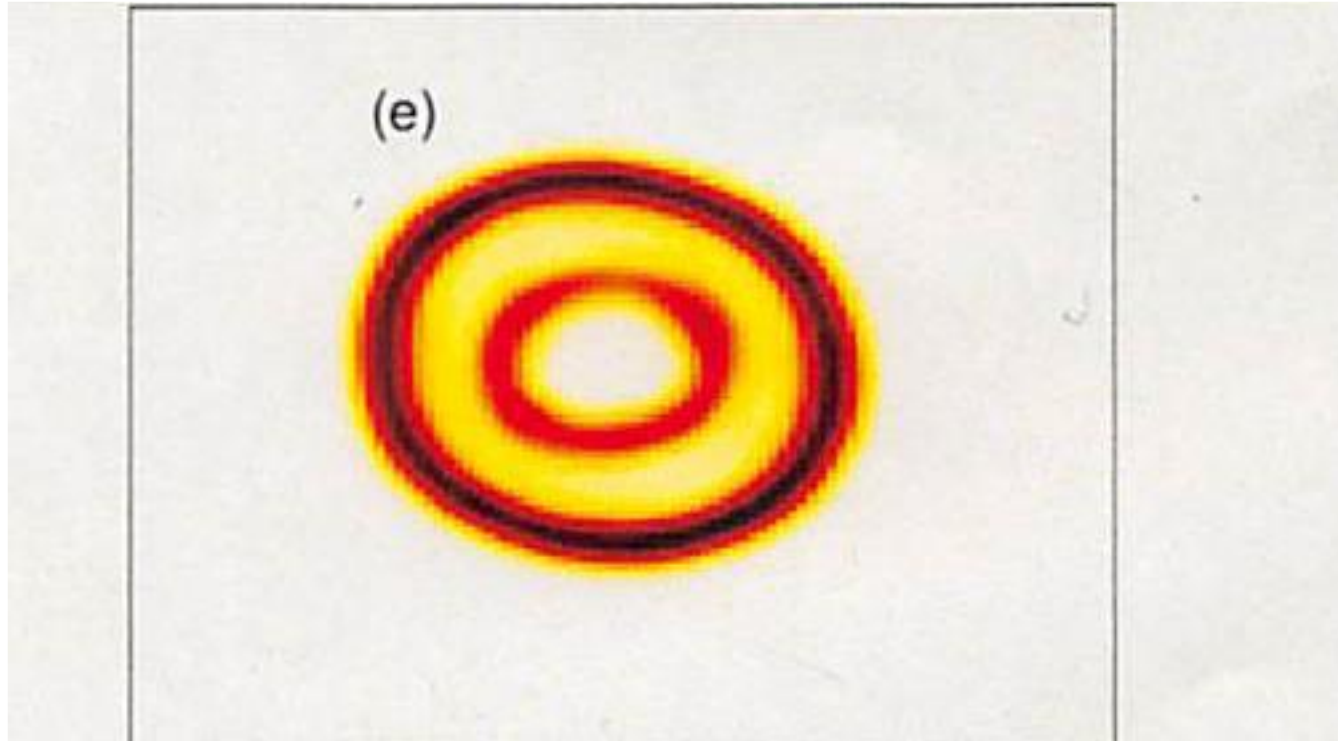


# Cuturi: cluster center of double circles



# Surprising Results!! Cuturi





$D_{\text{Wass}}$

# Discrete Case

Minimize F  
constraints

$$F_\lambda(\mathbf{P}) = \frac{1}{1+\lambda} \langle \mathbf{c}, \mathbf{P} \rangle - \frac{\lambda}{1+\lambda} H(\mathbf{P}).$$

$$c(\mathbf{P}) = \langle \mathbf{c}, \mathbf{P} \rangle = \sum c_{ij} P_{ij}.$$

$$\sum_j P_{ij} = p_i, \quad \sum_i P_{ij} = q_j, \quad \sum_{ij} P_{ij} = 1.$$

# Optimal Transportation Plan

$$L_\lambda(\mathbf{P}) = \frac{1}{1+\lambda} \langle \mathbf{c}, \mathbf{P} \rangle - \frac{\lambda}{1+\lambda} H(\mathbf{P}) - \sum_{i,j} (\alpha_i + \beta_j) P_{ij}.$$

$$\frac{\partial}{\partial P_{ij}} L_\lambda(\mathbf{P}) = \frac{1}{1+\lambda} c_{ij} + \frac{\lambda}{1+\lambda} \log P_{ij} - \alpha_i - \beta_j + \frac{\lambda}{1+\lambda}.$$

$$P_{ij} = \exp \left\{ -\frac{c_{ij}}{\lambda} + \frac{1+\lambda}{\lambda} (\alpha_i + \beta_j + 1) \right\}.$$

$$K_{ij} = \exp \left\{ -\frac{c_{ij}}{\lambda} \right\},$$

$$a_i = \exp \left( \frac{1 + \lambda}{\lambda} \alpha_i \right) \quad b_j = \exp \left( \frac{1 + \lambda}{\lambda} \beta_j \right),$$

the optimal solution is written as

$$P_{ij}^* = a_i b_j K_{ij},$$

# Exponential Family of Optimal Transportation Plans

$$P(x) = \sum_{i,j=1}^n P_{ij} \delta_{ij}(x). \quad \theta^{ij} = \log \frac{P_{ij}}{P_{nn}}, \quad \theta = (\theta^{ij}),$$

$$P(x, \theta) = \exp \left\{ \sum_{i,j} \theta^{ij} \delta_{ij}(x) + \log P_{nn} \right\}$$

$$P(x, \alpha, \beta) = \exp \left[ \sum_{i,j} \left\{ \frac{\lambda+1}{\lambda} (\alpha_i + \beta_j) - \frac{c_{ij}}{\lambda} \right\} \delta_{ij}(x) - \frac{(\lambda+1)}{\lambda} \psi \right]$$

$$\psi(\alpha, \beta) = \frac{\lambda}{1+\lambda} \log \sum_{i,j} \exp \left\{ \frac{\lambda+1}{\lambda} (\alpha_i + \beta_j) - \frac{1}{\lambda} (c_{ij}) \right\}$$



$$\theta^{ij} = \frac{1 + \lambda}{\lambda} (\alpha_i + \beta_j) - \frac{c_{ij}}{\lambda}$$

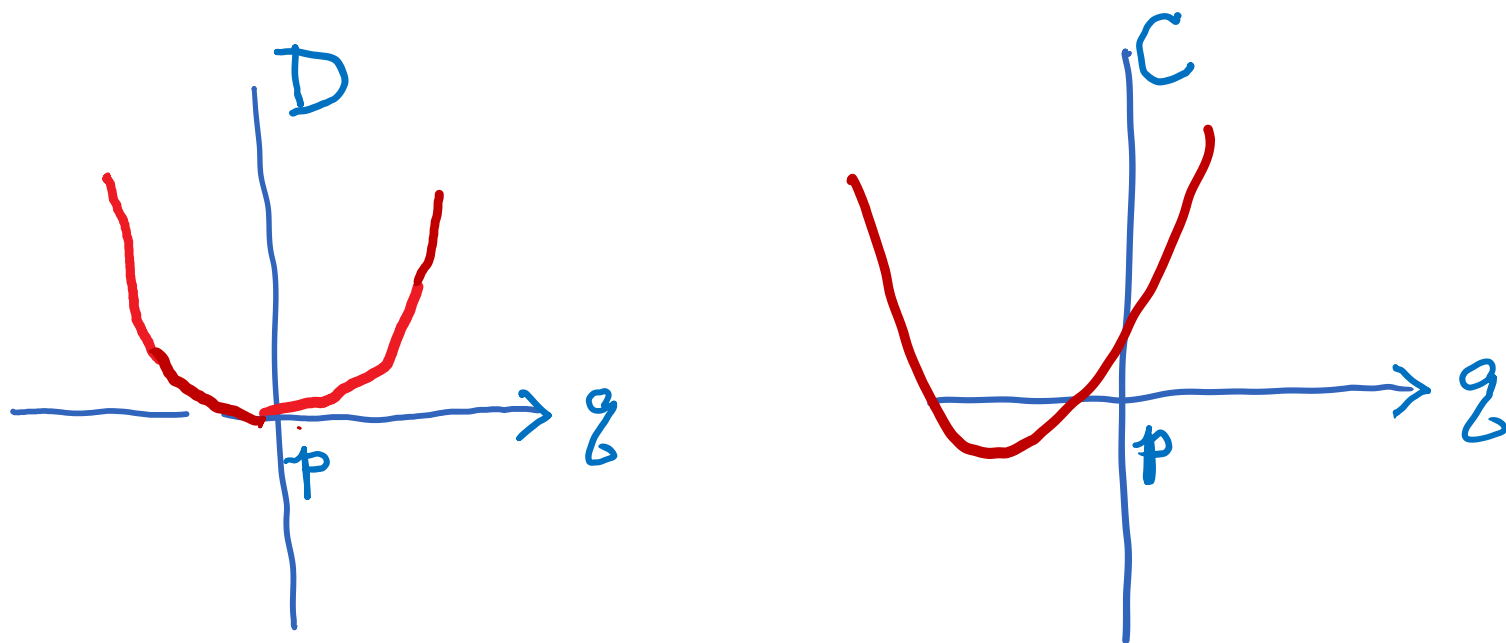
$$\begin{aligned} \varphi_\lambda(\mathbf{p}, \mathbf{q}) &= \frac{1}{1 + \lambda} \langle \mathbf{c}, \mathbf{P} \rangle + \frac{\lambda}{1 + \lambda} \sum_{i,j} P_{ij} \left\{ \frac{1 + \lambda}{\lambda} (\alpha_i + \beta_j) - \frac{c_{ij}}{\lambda} - \frac{(1 + \lambda)}{\lambda} \psi_\lambda \right\} \\ &= \mathbf{p} \cdot \boldsymbol{\alpha} + \mathbf{q} \cdot \boldsymbol{\beta} - \psi_\lambda(\boldsymbol{\alpha}, \boldsymbol{\beta}). \end{aligned} \quad (33)$$

$$\psi_\lambda(\boldsymbol{\theta}) + \varphi_\lambda(\boldsymbol{\eta}) = \boldsymbol{\theta} \cdot \boldsymbol{\eta}, \quad \boldsymbol{\eta} = (\mathbf{p}, \mathbf{q})^T, \boldsymbol{\theta} = (\boldsymbol{\alpha}, \boldsymbol{\beta})^T$$

# C関数には問題あり

$$C_\lambda(p, q) \Rightarrow D_\lambda(p, q)$$

$q = p$  is not the minimizer of  $C_\lambda(p, q)$



# 新しいダイバージェンス：その幾何学

$$D_\lambda(p : q) = C_\lambda(p : K_\lambda q) - C_\lambda(p : K_\lambda p)$$

$K_\lambda$ : diffusion operator

$$\tilde{D}_\lambda(p : q) = C_\lambda(p : q) - \frac{1}{2} \{C_\lambda(p : p) - C_\lambda(q : q)\}$$

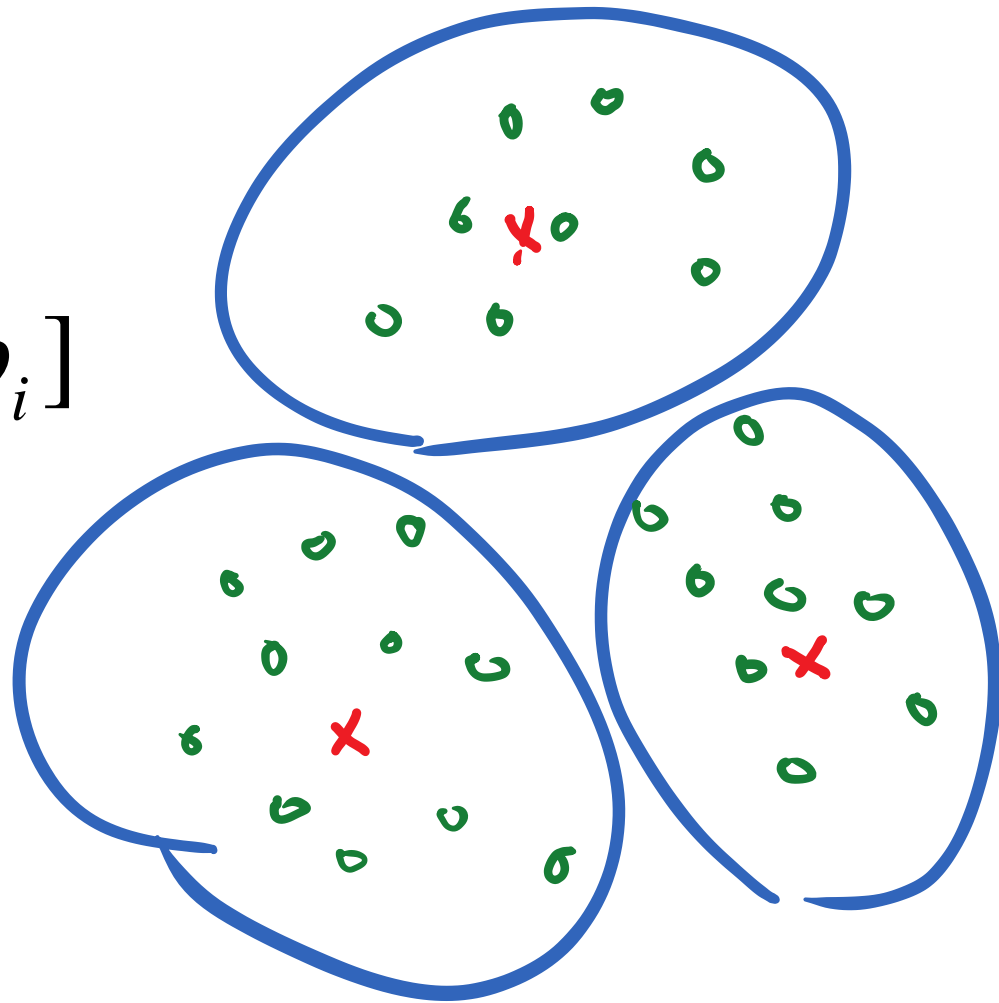
# エントロピー制約の一般化

q-エントロピー

$$\min F = \langle c, P \rangle - \lambda H[P(x, y)]$$

# clustering

$$p_{center} = \arg \min_p \sum D[p : p_i]$$



# W-GAN

$$D[p_r : p_g] = E_{p_r} [\log D(x)] + E_{p_g} [1 - \log D(G(z))]$$

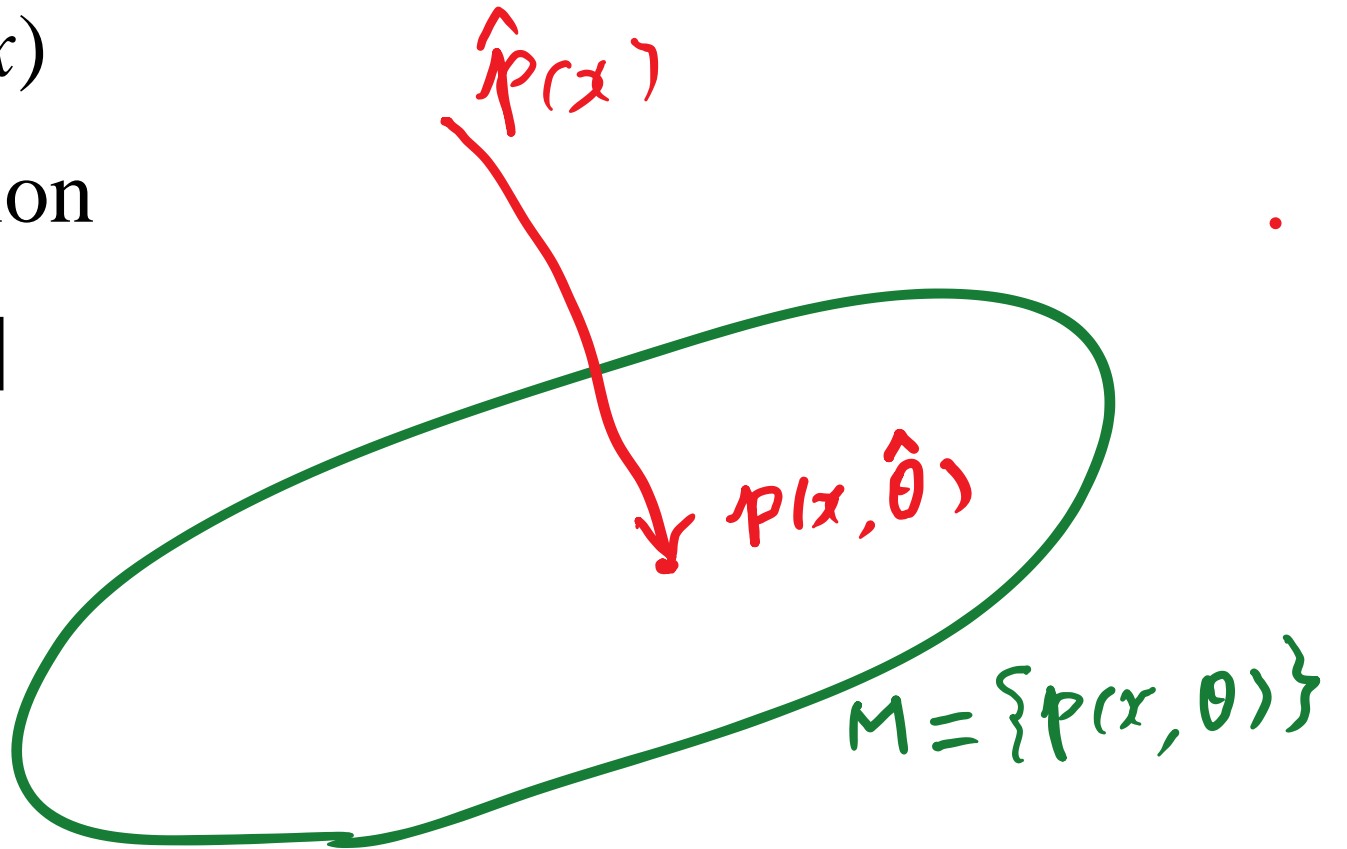
$$\text{KL}[p_r : p_m] + \text{KL}[p_g : p_m]$$

$$D_{\text{Wass}}$$

# W-statistics vs likelihood statistics

$p(x, \theta) : x_1, x_2, \dots, x_N \rightarrow \hat{p}(x)$   
empirical distribution

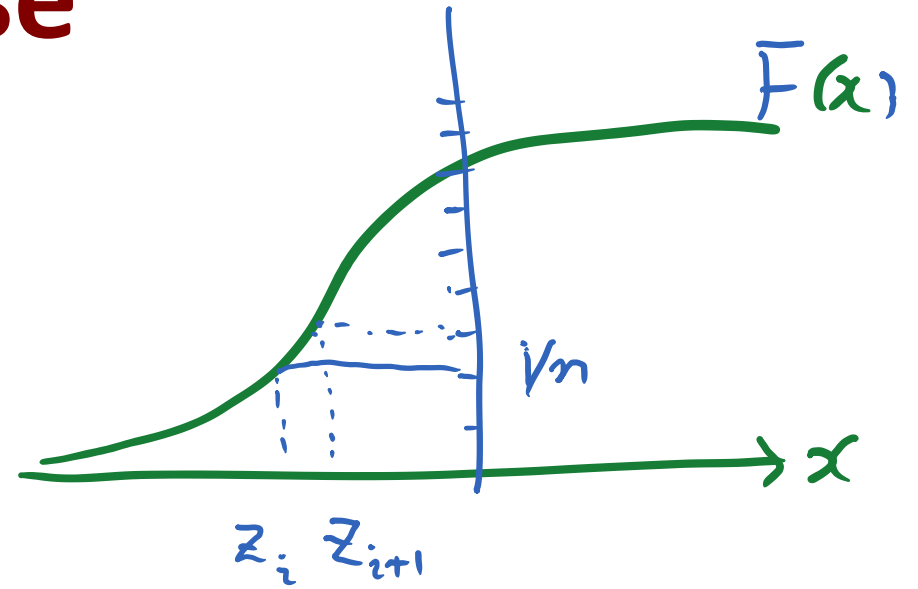
$$\hat{\theta} = \operatorname{argmin} D[\hat{p}(x) : p(x, \theta)]$$



# Estimation : Gaussian case

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

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



$$KL: \bar{x} = \frac{1}{N} \sum x_i; \quad \hat{\sigma}^2 = \frac{1}{N} \sum (x_i - \mu)^2$$

$$D_W: \bar{x} = \frac{1}{N} \sum x_i; \quad \hat{\sigma} = \frac{1}{N} \sum z_i x_i$$

$$z_i = F^{-1}\left(\frac{i}{N}\right)$$



# W-statistics: $\lambda = 0, \quad X = \mathbb{R}$

Model:  $p(x, \xi) \quad x_1 \leq x_2 \leq \dots \leq x_n$

Observed data

Partition points  $z_i(\xi) = P\left(\frac{i}{n}\right) = \int_{-\infty}^{i/n} p(x, \xi) dx$

# Optimal transport plan

$$x_i \rightarrow z_i(\xi)$$

**Cost**

$$C(\xi) = \frac{1}{n} \sum_i |x_i - z_i(\xi)|^2$$

**Estimating equation**

$$\sum_i \{x_i - z_i(\xi)\} \partial_\xi z_i(\xi) = 0$$

# Consistency and efficiency

$$\lim_{n \rightarrow \infty} E[\hat{\xi}] = \xi$$

$$V[\hat{\xi}] = \frac{1}{n} G^{-1} H G^{-1}$$

$$\begin{aligned} G(\xi) &= \int \partial_{\xi} P(x, \xi) \{ \partial_{\xi} P(x, \xi) \}^T dx \\ &= \partial_{\xi} \partial_{\xi} C(\xi', \xi) |_{\xi' = \xi} \end{aligned}$$

$$H(\xi) = \int P^2(x, \xi) \partial_{\xi} P(x, \xi) \{ \partial_{\xi} P(x, \xi) \}^T dx$$

# Information Geometry of Sinkhorn Algorithm

Obtaining  $a$  and  $b$  in  $P^* = ca b K$

$$M_p = \{P_{ij} \mid \sum_j P_{ij} = p_i\}$$

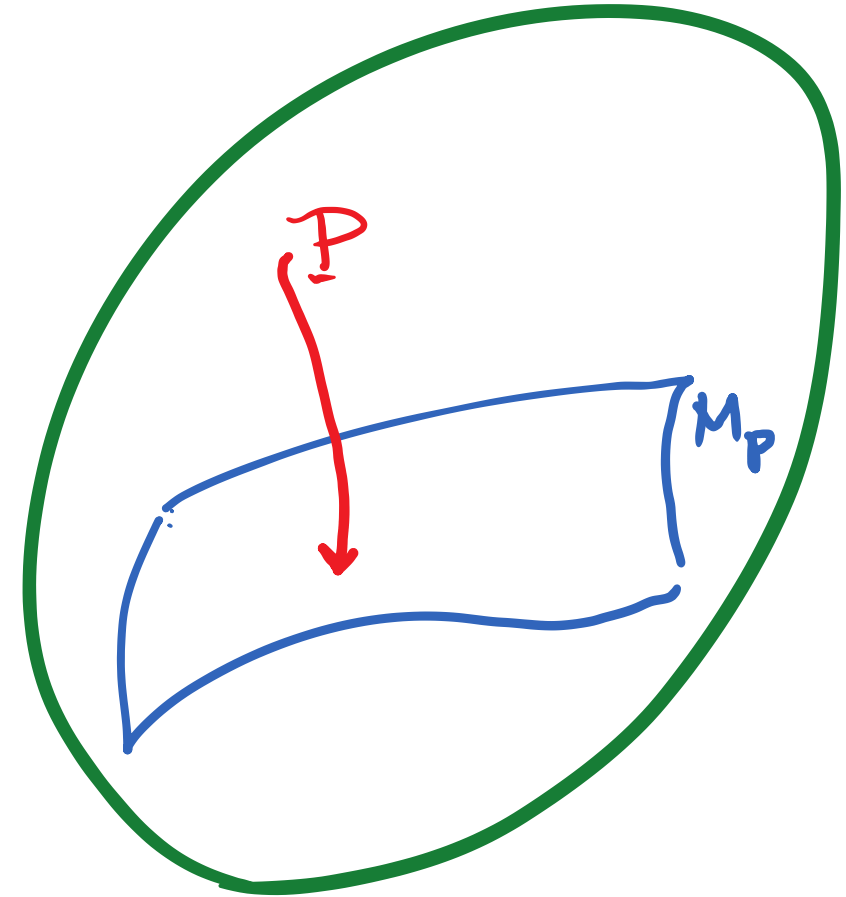
e-projection of  $P$  to  $M$  :

$$M_p = \{P_{ij} \mid \sum_j P_{ij} = p_i\}$$

$$M_q = \{P_{ij} \mid \sum_i P_{ij} = q_j\}$$

$$P_{ij} \rightarrow a_i P_{ij} \quad a_i = \frac{p_i}{\sum_j P_{ij}}$$

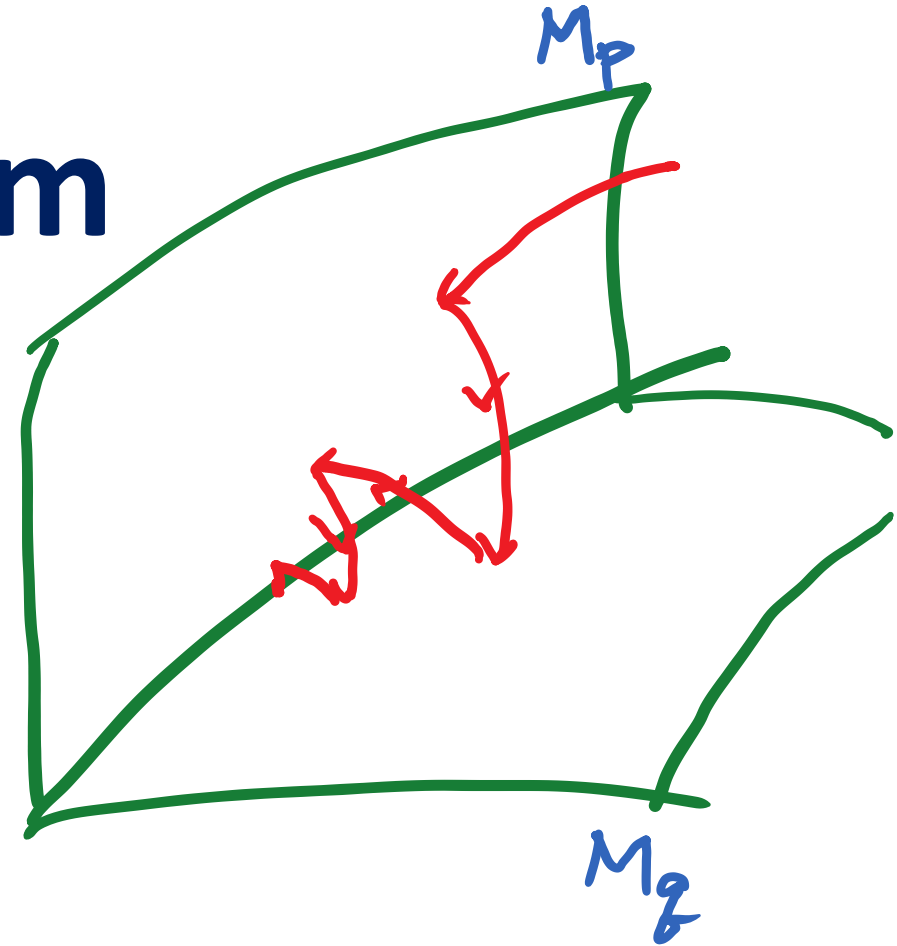
$$P_{ij} \rightarrow b_j P_{ij} \quad b_j = \frac{q_j}{\sum_i P_{ij}}$$



# Iterative Algorithm

e-projection to  $M$   
e-projection to  $M$

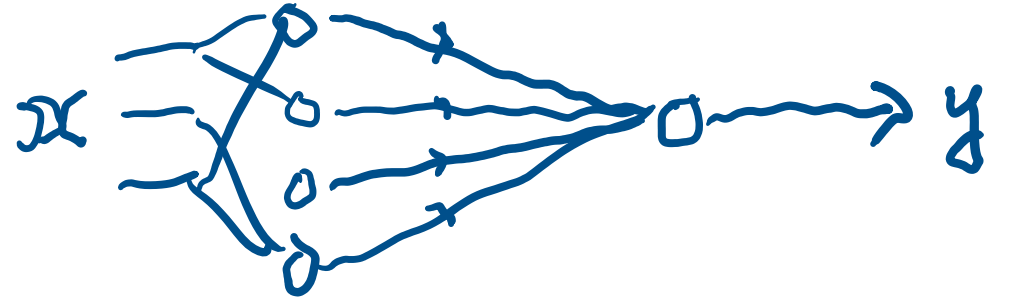
easy to solve (not LP)



# 3層パーセプトロン学習のW幾何

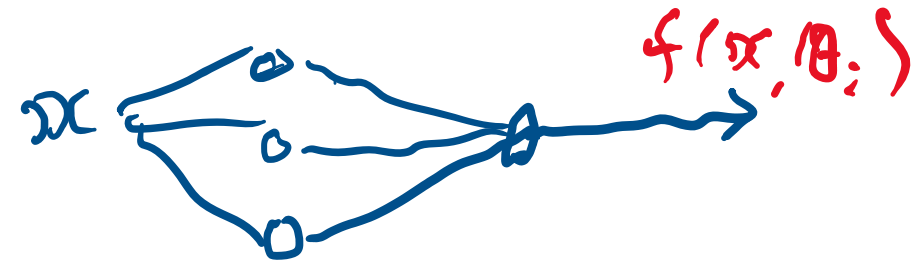
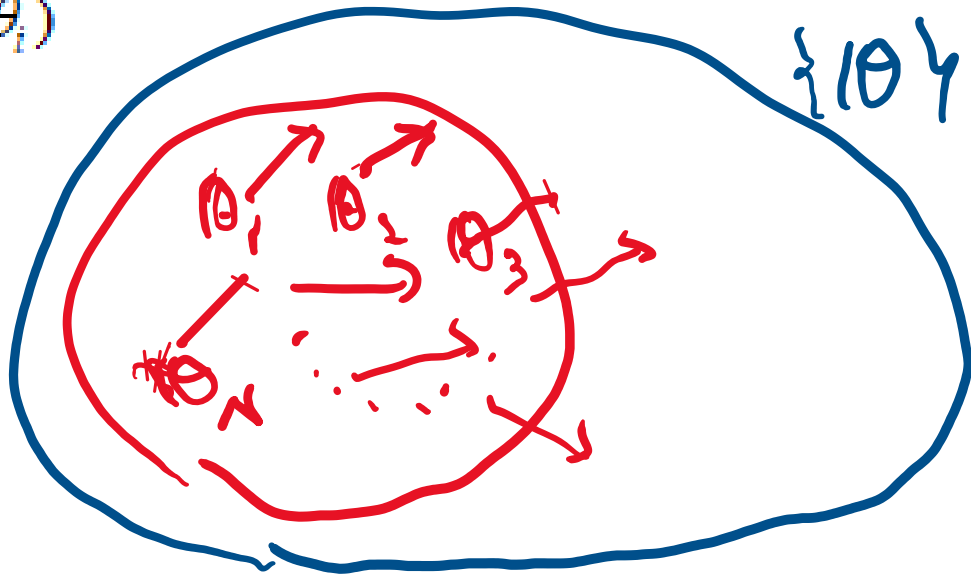
$$y = f(\mathbf{x}, \Theta) = \frac{1}{N} \sum f(\mathbf{x}, \theta_i);$$

$$f(\mathbf{x}, \theta_i) = v_i \phi(\mathbf{w}_i \cdot \mathbf{x} + b_i)$$



$$l = \frac{1}{2} \{y - f(\mathbf{x}, \Theta)\}^2 = \frac{1}{2} y^2 - \frac{1}{N} \sum y f(\mathbf{x}, \theta_i) + \frac{1}{2N^2} \sum f(\mathbf{x}, \theta_i) f(\mathbf{x}, \theta_j)$$

$$\dot{\theta}_i = \frac{\partial l}{\partial \theta_i} = v_i(\mathbf{x}, y, \theta_i)$$



$$l = \frac{1}{2} \{y - f(\mathbf{x}, \Theta)\}^2 = \frac{1}{2} y^2 - \frac{1}{N} \sum y f(\mathbf{x}, \theta_i) + \frac{1}{2N^2} \sum f(\mathbf{x}, \theta_i) f(\mathbf{x}, \theta_j)$$

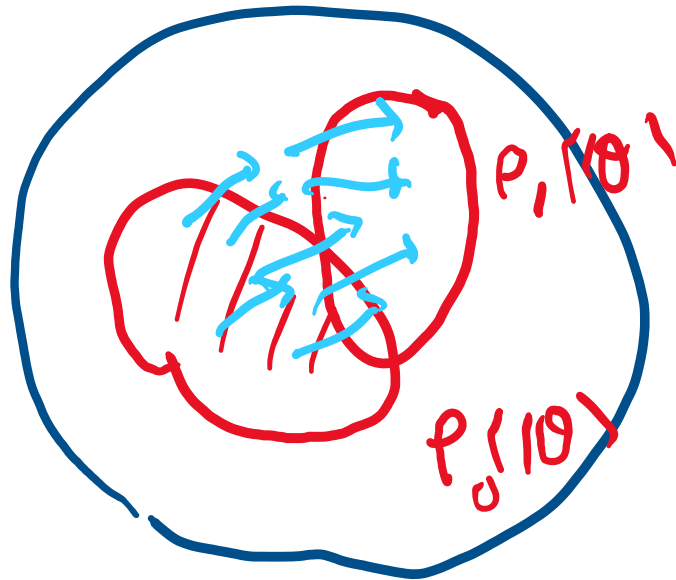
$$\rho(\theta) = \frac{1}{N} \sum \delta(\theta - \theta_i)$$

$$V(\theta) = \langle y f(\mathbf{x}, \theta_i) \rangle_{\rho}; \quad U(\theta) = \frac{1}{2} \langle f(\mathbf{x}, \theta) f(\mathbf{x}, \theta') \rangle_{\rho, \rho'}$$

$$\Psi(\theta) = V(\theta) + \langle U(\theta, \theta') \rangle_{\rho, \rho'}$$

$$v_i(\mathbf{x}, y, \Theta) = \nabla \Psi(\Theta)$$

$$\dot{\rho}_t(\theta) = \eta \nabla_{\theta} \cdot \{ \rho_t(\theta) \nabla \Psi(\Theta) \}$$



**深層学習とW幾何 拡散モデル**

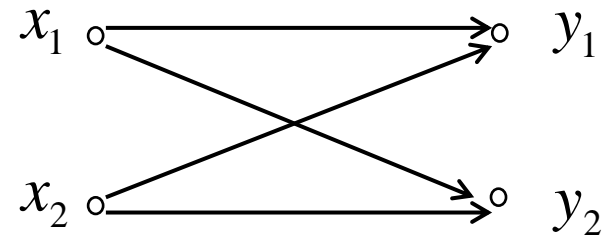
**3層パーセプトロン学習のW幾何**

**アファイン変換モデルのW幾何**



# Information Integration and Complexity of Systems

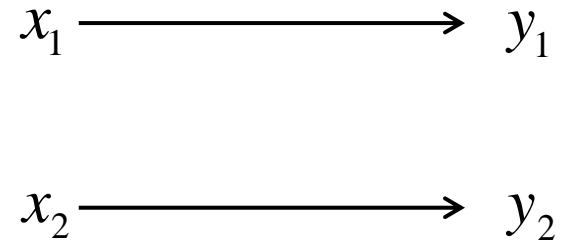
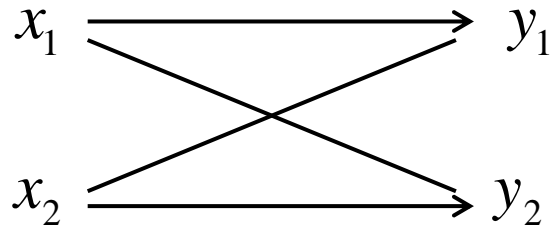
$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}) p(\mathbf{y} | \mathbf{x})$$



Shun-ichi Amari (*RIKEN Brain Science Institute*)

Masafumi Oizumi (RIKEN BSI, Monash U.)

Naotsugu Tsuchiya (Monash U.)



full model:  $S_F = \{p(\mathbf{x}, \mathbf{y})\}$

split model:  $S_S = \{q(\mathbf{x}, \mathbf{y})\}$

$$q(\mathbf{y} | \mathbf{x}) = \prod q(y_i | x_i)$$

**measure of interaction : N. Ay**

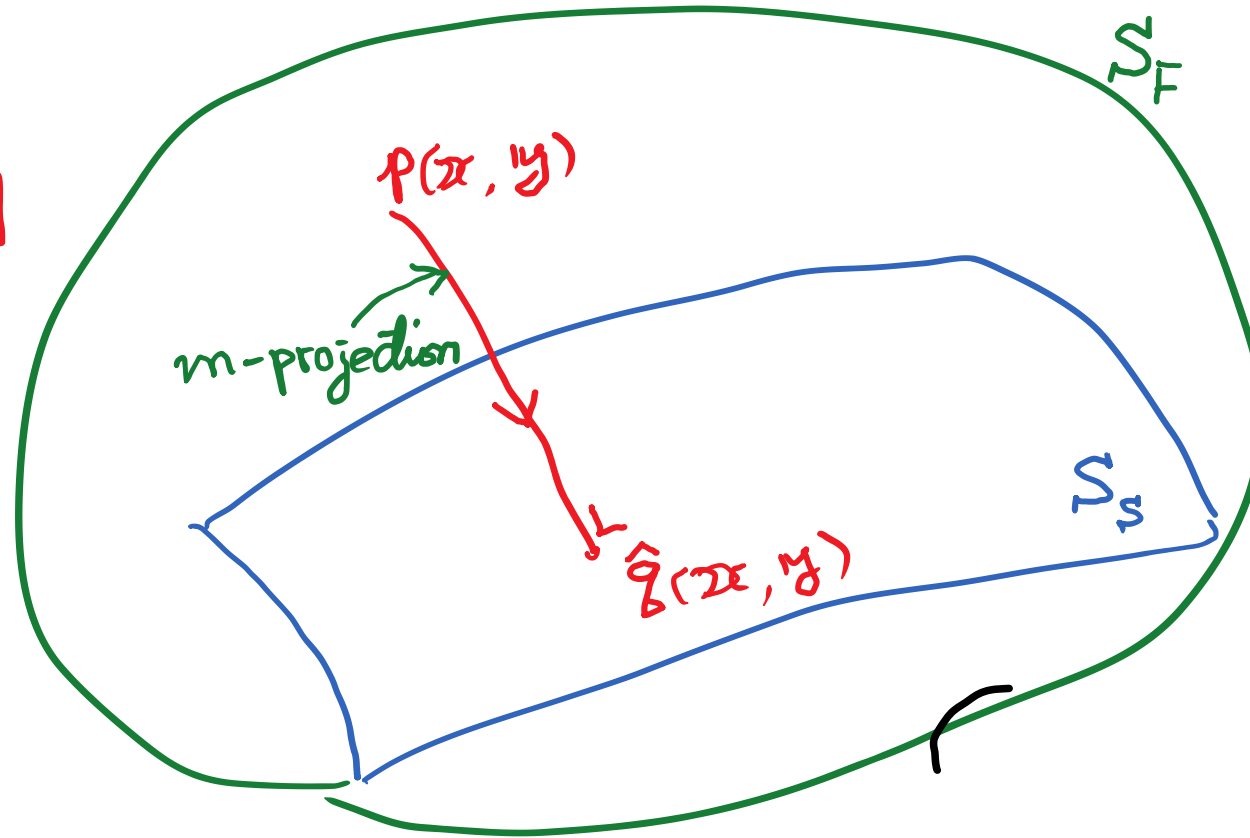
**information integration : Tononi**

**Barrett and Seth**

# Measure of information integration, or system complexity $\Phi$

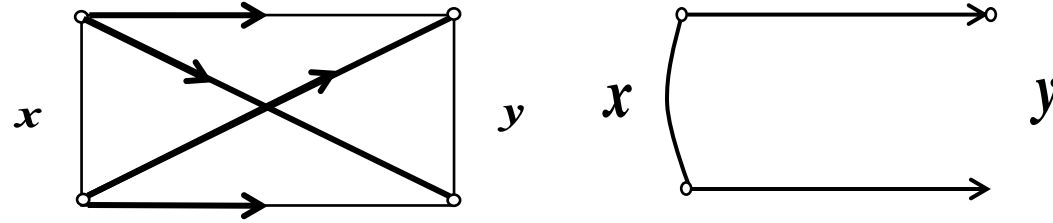
Information Geometry N. Ay

$$\Phi = \mathcal{D}_{KL} [P : \hat{P}]$$



# Split Model $S_H$ : Ay, Barrett & Seth

$$q(\mathbf{x}, \mathbf{y}) = q(\mathbf{x}) \prod q(y_i | x_i)$$

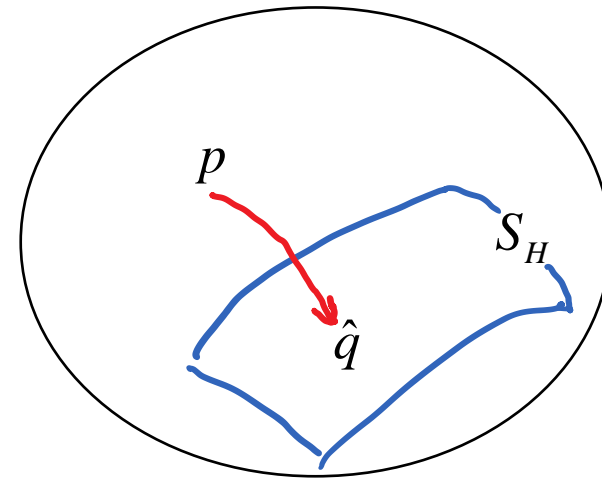


$$\theta_{12}^{XY} = \theta_{21}^{XY} = \theta_{12}^Y = 0$$

$$\Phi_H = D_{KL}[p : S_H] = \min_{q \in S_H} D_{KL}[p : q]$$

$$\hat{q} = \prod_{M_S} p : \hat{q}(y|x) = \prod p(y_i | x_i)$$

$$\Phi_H = \sum H[Y_i | X_i] - H[\mathbf{Y} | \mathbf{X}]$$



# Split Model $S_G$

$$q(\mathbf{x}, \mathbf{y}) = q_X(\mathbf{x}) \tilde{q}_Y(\mathbf{y}) \prod q(y_i | x_i)$$

$$\theta_{12}^{XY} = \theta_{21}^{XY} = 0$$

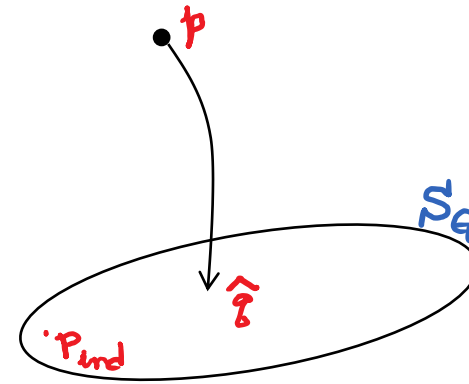
$$q(x_1, y_2 | x_2, y_1) = q(x_1 | x_2, y_1) q(y_2 | x_2, y_1)$$

$$0 \leq \Phi \leq I(X : Y)$$

$$\hat{q}_X(\mathbf{x}) = p_Y(\mathbf{x}), \quad \hat{q}_Y(\mathbf{y}) = p_Y(\mathbf{y})$$

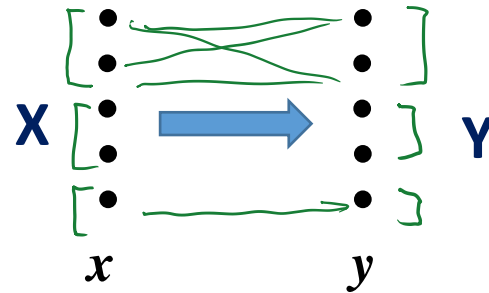
$$\hat{q}(y_i | x_i) = p(y_i | x_i)$$

graphical model

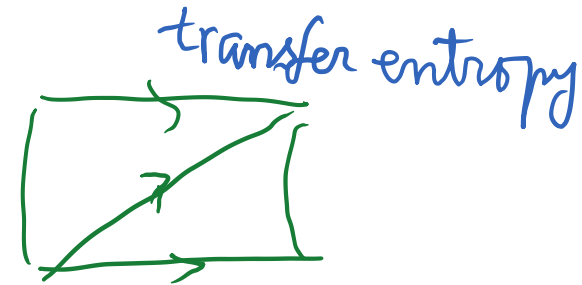


# Hierarchy: transfer entropy

Partition of X



cutting branches  
split models



$$\cup X_i = X,$$

$$X_i \cap X_j = \phi$$

$$\cup Y_i = Y,$$

$$Y_i \cap Y_j = \phi$$

Partition

