

**IGAIA 4 Bohemia**

# **Information Geometry**

**— Historical Episodes and Future  
with Recent Developments**

**Shun-ichi Amari**

**RIKEN Brain Science Institute**

# Prehistory --- Riemannian Geometry

**H. Hotteling      1929**

Riemannian metric and Fisher information  
location-scale model : constant curvature

**P. Ch. Maharanobis    1936    Euclidean distance (multivariate-Gaussian)**

**C. R. Rao              1945    Cramer-Rao Theorem; Riemannian**

**H. Jeffreys            1946    Bayesian theory and Jeffreys invariant prior**

# Dual Geometry, Invariance

N. Chentsov 1972 invariance,  $\{g, T\}$ ,  $\alpha$ -connection

B. Efron 1975 (A. P. Dawid) statistical curvature; higher-order asymptotics

O. Barndorff-Nielsen 1976 exponential family; Legendre transform

S. Amari 1982 duality; curvature and statistics (M. Kumon)

H. Nagaoka and S. Amari 1982 duality, Pythagorean theorem

# Amari's personal history

**1958: statistics seminar (master course at U Tokyo)**

S. Kullback, "Information and Statistics"

Riemannian metric (suggested by S. Moriguti)

Gaussian  $N(\mu, \sigma^2)$  : geodesic, constant curvature (Poincare-half plane)

**beautiful structure → essential meaning?**

mathematical engineering

graph and topology of networks: homology

non-Riemannian geometry of materials manifold: dislocations

information systems, learning and neural networks

# Statistical curvature and higher-order inference

B. Efron, 1975

Fisher's idea; exponential connection and mixture connection

A. P. Dawid, 1975 e- and m-connections

S. Amari :  $\alpha$ -geometry 
$$Error^2 = \frac{1}{n}G^{-1} + \frac{1}{n^2}(H_e^2 + H_m^2 + \Gamma_m^2) + \frac{1}{n^3}K$$

(Rao, Kano K? Fisher's dream)

Amari and M.Kumon higher-order power of statistical test

# Amari paper : Ann. Statist. 1982:

Reviewers (S. Lauritzen and A.P. Dawid)

Chentsov work (handwritten manuscript)

H. Nagaoka and S. Amari 1982 (Technical Report)

Ann. Probability Theory 7 reviewers

Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete  
geometry has nothing to do with statistics

IEEE Trans. Inf. Theory Shannon Theory, now well-known

# **London Workshop: 1984 (D. Cox)**

**Cox visited Japan in 1983**

**patron of information geometry**

**Rao, Efron, Dawid, Barndorff-Nielsen, Lauritzen  
Kass, Eguchi many others**

**Dodson, Critchley, Marriot, Komaki, Zhang, Ay, Pistone, Giblisco, Nielsen, ...**

# Information Geometry --- lucky naming

## Applications area:

statistics, time-series and systems,

machine learning, signal processing, optimization

brain theory, consciousness

physics, economics, mathematics

(Banach manifold, affine differential geometry and beyond)

quantum information, Tsallis entropy



# International Conferences

**IGAIA series; GSIS series, ...**

**Many monographs**

**new journal (Jun Zhang); where to publish**

**mailing list and society**

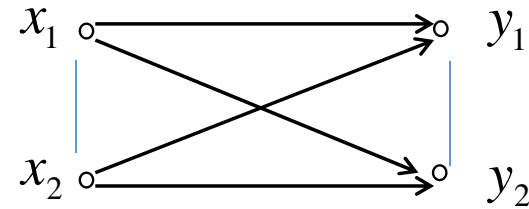
**still small community; united and cooperative, blessed by all**

# My recent works

1. **Systems complexity and consciousness (IIT)**
2. **Geometry of score matching (Hyvarinen score)**
3. **Natural gradient descent and topology of deep learning)**
4. **Canonical divergence**
5. **Multi-terminal statistical inference**
6. **Information geometry and Wasserstein distance**

# Information Integration and Complexity of Systems

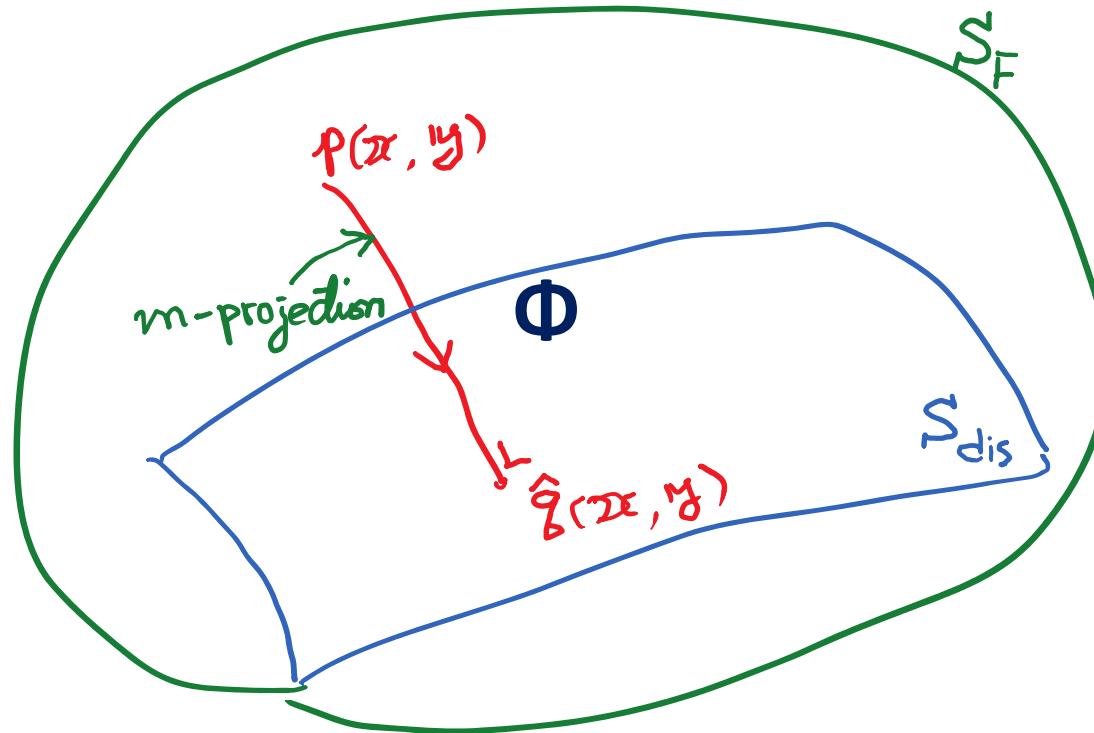
-- Stochastic approach



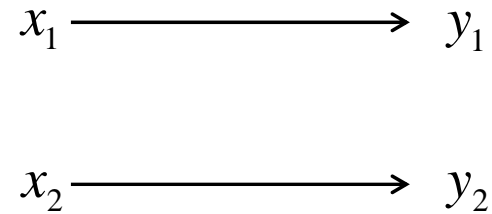
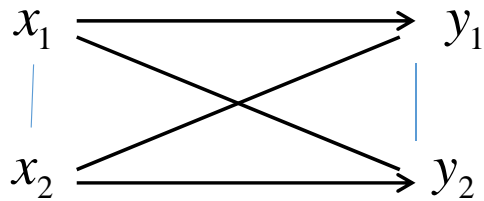
$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}) p(\mathbf{y} | \mathbf{x})$$

**x: state of the brain**      **y: next state of the brain**

# Integrated Information Theory G. Tononi



**Necessary condition; sufficient?**



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

Disconnected model:

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

measure of interaction : N. Ay

information integration : Tononi

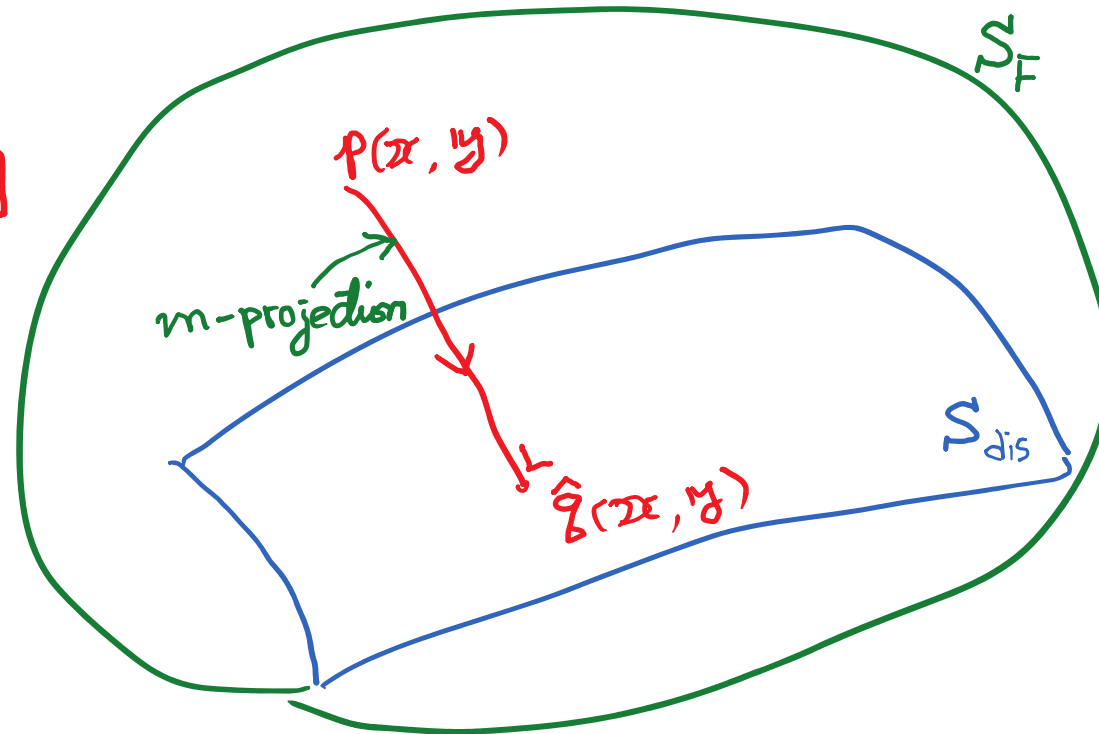
Barrett and Seth

Many other  $\Phi$

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

Information Geometry N. Ay

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



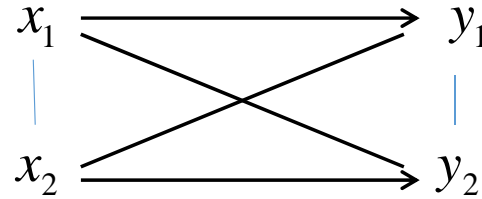
## Definition of $\Phi$ : Postulates

1) 
$$\Phi = \min_q D[p(x, y) : q(x, y)], \quad q \in \mathcal{S}_{dis}$$

2) 
$$D = D_{KL}[p : q] = \int p(x, y) \log \frac{p(x, y)}{q(x, y)}$$

3) **Disconnected model:  $\mathcal{S}_{dis}$  Markov conditions**

## Markov Condition



**(1→2) branch deleted: Markov condition:**  $x_1 \rightarrow x_2 \rightarrow y_2$

$$p(x_1, y_2 \mid x_2) = p(x_1 \mid x_2) p(y_2 \mid x_2)$$

$$X_1 - X_2 - Y_2$$

$$X_2 - X_1 - Y_1$$

$S_{dis}$  : all  $x_i \rightarrow y_j$  ( $i \neq j$ ) deleted



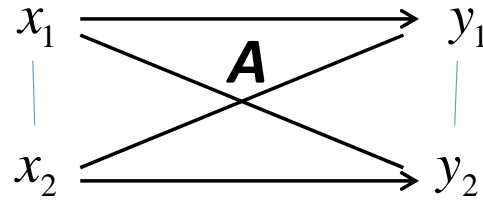
## Why KL-divergence?

- 1)  $D[p : q] \geq 0$ ,  $= 0$ , when and only when  $p = q$
- 2)  $D[p : q]$  invariant under transformations of  $\mathbf{X}$
- 3)  $D[p : q] = \sum_i d[p(x_i), q(x_i)]$
- 4)  $D[p : q]$  induces flat structure dually

## Geometric degree of information integration

$$\Phi_{geo} = \min_q D_{KL}[p(x, y) : q(x, y)], \quad q \in \mathcal{S}_{dis}$$

## Gaussian case

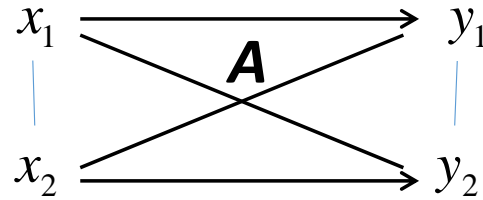


$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e} \quad \Sigma = \mathbf{E}[\mathbf{e}\mathbf{e}^T]$$

$$\mathbf{y} = \mathbf{A}'\mathbf{x} + \mathbf{e}' \quad \Sigma' = \mathbf{E}[\mathbf{e}'\mathbf{e}'^T] \quad \mathbf{A}': \text{diagonal}$$

$$\Phi_{geo} = \log \frac{|\Sigma'|}{|\Sigma|}$$

## Gaussian case



$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e} \quad \Sigma = \mathbf{E}[\mathbf{e}\mathbf{e}^T]$$

$$\mathbf{y} = \mathbf{A}'\mathbf{x} + \mathbf{e}' \quad \Sigma' = \mathbf{E}[\mathbf{e}'\mathbf{e}'^T] \quad \mathbf{A}': \text{diagonal}$$

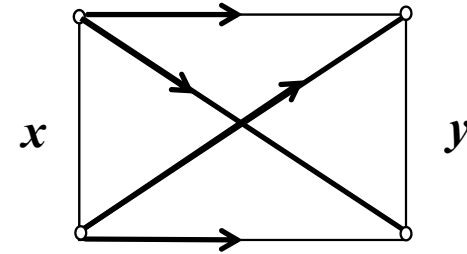
$$\Phi_{geo} = \log \frac{|\Sigma'|}{|\Sigma|}$$

## Many other definitions of $\Phi$

Full model

Disconnected models

**Full model  $S_F$  : graphical model**



$$p(\mathbf{x}, \mathbf{y}) = \exp\left\{\sum \theta_i^X x_i + \sum \theta_i^Y y_i + \theta_{12}^X x_1 x_2 + \theta_{12}^Y y_1 y_2 + \sum \theta_{ij}^{XY} x_i y_j + \text{higher-order terms} - \psi(\boldsymbol{\theta})\right\}$$

**exponential family**

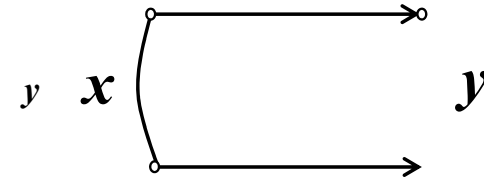
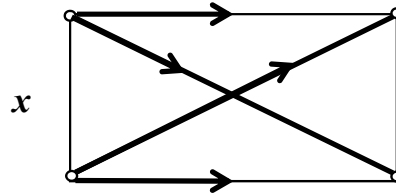
**$\theta$ -coordinates**  $\boldsymbol{\theta}$

**$\eta$ -coordinates**  $\eta_i^X = E[x_i], \dots, \eta_{ij}^{XY} = E[x_i y_j], \dots$

**There are many disconnected models!!**

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

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



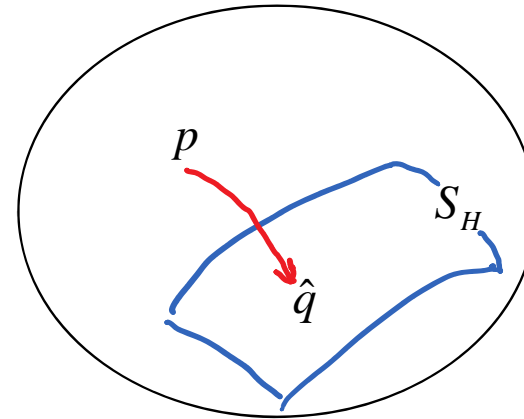
$$\theta_{12}^{XY} = \theta_{21}^{XY} = \theta_{12}^Y = 0$$

$$Y_1 - X_1 - X_2 - Y_2$$

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

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

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



$$S_H : \theta_{12}^{XY} = \theta_{21}^{XY} = \theta_{12}^Y = 0; \quad \hat{\eta}_q = \eta_q$$



## Mixed Coordinates :

$$\xi = (\eta_i^X, \eta_{12}^X, \eta_i^Y, \eta_{11}^{XY}, \eta_{22}^{XY}; \theta_{12}^{XY}, \theta_{21}^{XY}, \theta_{12}^Y)$$

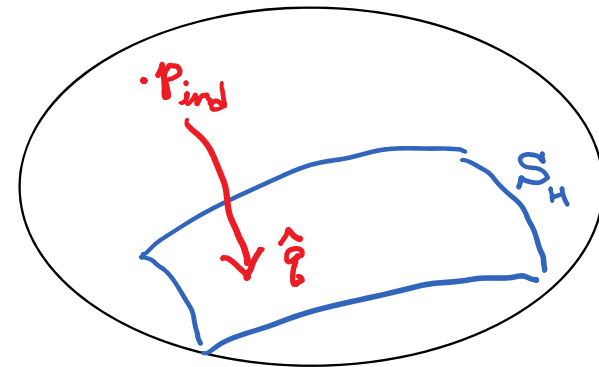
$$\hat{q} = \prod p = q(\mathbf{x}, \mathbf{y}; \hat{\xi})$$

## Markovian Condition

$$Y_1 - X_1 - X_2 - Y_2 \Rightarrow X_1 - X_2 - Y_2; Y_1 - X_1 - X_2$$

**problem**  $0 \leq \Phi[p] \leq I[X:Y]$

$$p_{ind} = p_X(\mathbf{x}) p_Y(\mathbf{y}) \Rightarrow I = 0, \quad \Phi > 0$$





# Split Model $S_{Gr}$

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

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

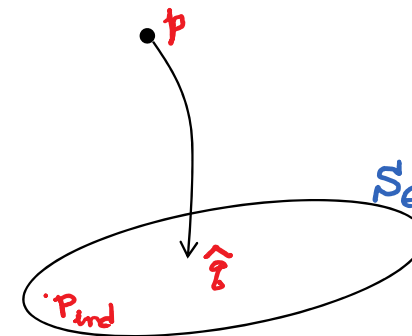
$$q(x_1, y_2 | x_2, y_1) = q(x_1 | x_2, y_1) q(y_2 | x_2, y_1)$$

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

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

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

graphical model

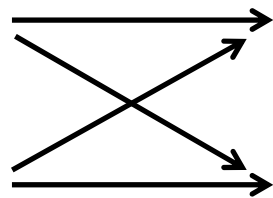


## Problem: Gaussian channel

$$p(\mathbf{x}, \mathbf{y}): \mathbf{y} = A\mathbf{x} + \boldsymbol{\varepsilon} \Rightarrow \hat{q}: \mathbf{y} = \hat{A}\mathbf{x} + \boldsymbol{\varepsilon}' \in S_G$$

$$p(\mathbf{x}, \mathbf{y}) = \exp\left[-\frac{1}{2}\left\{\mathbf{x}'\Sigma_x^{-1}\mathbf{x} + (\mathbf{y} - A\mathbf{x})'\Sigma_\varepsilon^{-1}(\mathbf{y} - A\mathbf{x}) - \psi\right\}\right]$$

$$\hat{A} \text{ is not diagonal} \iff \theta_{12}^{XY} = \theta_{21}^{XY} = 0$$



$A$

:



# Mismatched Decoding Model $S_M$

$$\mathbf{y} \rightarrow \hat{\mathbf{x}}$$

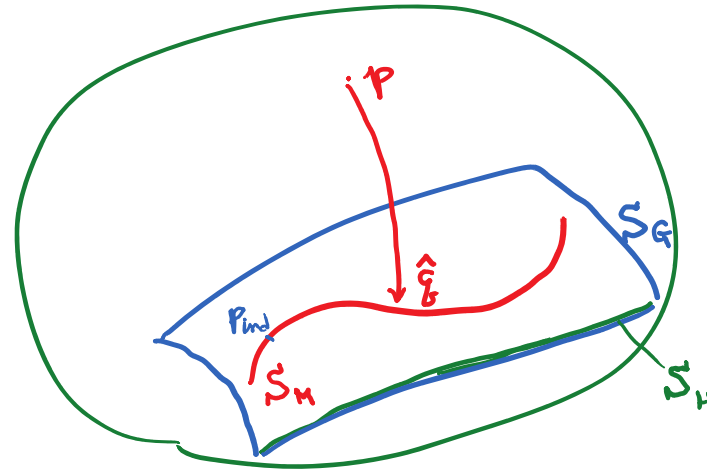
$$x_1 \longrightarrow y_1$$

$$x_2 \longrightarrow y_2$$

Best mismatched decoding from  $\mathbf{y}$  to  $\mathbf{x}$

$$S_M = \{q_\beta(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}) p_\beta(\mathbf{y}) \prod p(y_i | x_i)^\beta\}$$

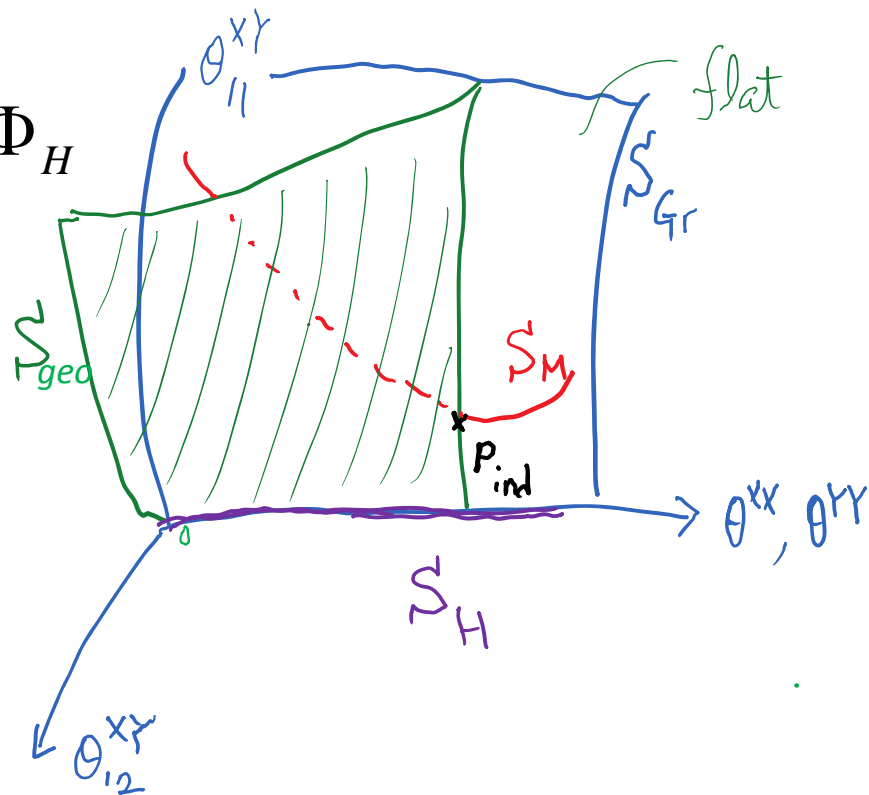
$$\Phi^* = D_{KL}[p : S_M]$$



$S_H \subset S_{Gr}, S_{Geo}$  :  $S_{Gr}, S_H$  dually flat

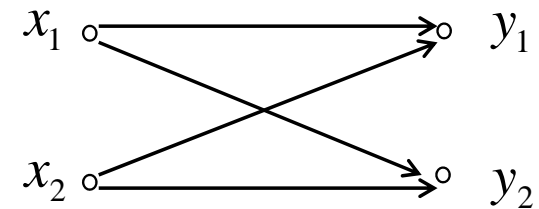
$S_M \subset S_{Gr}$  :  $S_M, S_{Geo}$  not flat

$\Phi_{Gr} \leq \Phi_H, \Phi_M$  ;  $\Phi_{Geo} \leq \Phi_H$



## Transfer Entropy; Granger causality

$$\begin{aligned} TE[x_i \rightarrow y_j] &= \min D_{KL}[p : q] \quad q \in (i \rightarrow j) \text{ disconnected} \\ &= H[Y_j | X - X_i] - H[Y_j | X] \end{aligned}$$

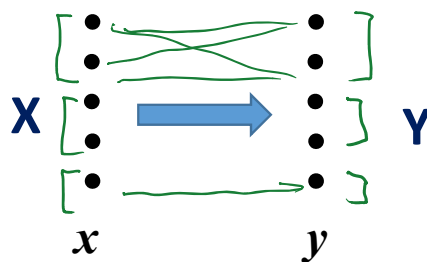


### Non-additive

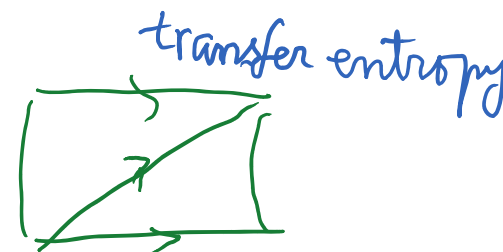
$$TE[x_i \rightarrow y_j; x_k \rightarrow y_m] \neq TE[x_i \rightarrow y_j] + TE[x_k \rightarrow y_m]$$

# Hierarchy: transfer entropy

Partition of X



cutting branches  
split models

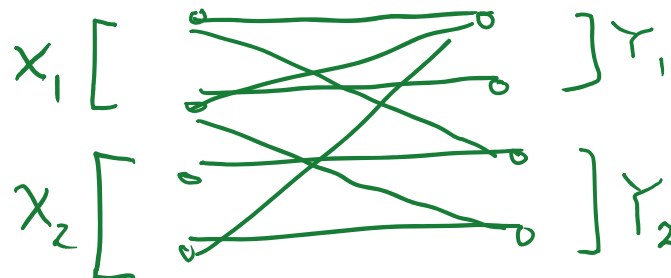


$$\cup X_i = X, \quad X_i \cap X_j = \phi$$

$$\cup Y_i = Y, \quad Y_i \cap Y_j = \phi$$

Partition

subadditivity



# Information Geometry of Hyvärinen Game Score

Following Grinwald, Dawid, Parry, Lauritzen, Hyvärinen

$$L(p, q) = E_{p(x)} [l(x, a)] : a = q(x)$$

$$S(x, q) = l(x, q) \quad l(x, q) = \log q(x)$$

$$S\text{-entropy} \quad H_S[p : q] = E_p [S(x, q)]$$

$$S\text{-divergence} \quad D_S[p : q] = H_S[p : q] - H_S[p : p]$$

# Hyvärinen score

$$S(x, q) = \ddot{l}(x, \xi) + \frac{1}{2} \{ \dot{l}(x, \xi) \}^2$$

$$D_S[p : q] = \frac{1}{2} E_p \left[ \frac{d}{dx} \{ \log p(x) - \log q(x, \theta) \}^2 \right]$$

$$D[p : cq] = D[p : q]$$

$$s(x, \xi) = \partial_\xi \ddot{l}(x, \xi) + \dot{l}(x, \xi) \partial_\xi \dot{l}(x, \xi)$$



**parametric case**  $M = \{q(\mathbf{x}, \boldsymbol{\xi})\} : A$

$$\min_{\boldsymbol{\xi}} D_S [p : q(\mathbf{x}, \boldsymbol{\xi})]$$

$$s(\mathbf{x}, \boldsymbol{\xi}) = \partial_{\boldsymbol{\xi}} S(\mathbf{x}, \boldsymbol{\xi}) \quad : \text{estimating function}$$

$$\sum_{i=1}^N s(\mathbf{x}_i, \boldsymbol{\xi}) = 0 \quad : \text{estimating equation}$$

**Information geometry of**  $D_S [p : q]$

## Asymptotic Analysis of estimator

$$E[\Delta\xi\Delta\xi^T] = \frac{1}{N} K^{-1}VK^{-T} \geq G^{-1}$$

$$K = E[\partial_\xi s(\mathbf{x}, \xi)]$$

$$V = E[s(\mathbf{x}, \xi)s(\mathbf{x}, \xi)^T] \rightarrow G$$

$$E[\Delta\xi\Delta\xi^T] = G^{-1}AG^{-1}$$

$$A = E[\mathbf{a}(\mathbf{x}, \xi)\mathbf{a}(\mathbf{x}, \xi)^T]$$

$$\mathbf{s}(\mathbf{x}, \xi) = c\{\partial_\xi \log(\mathbf{x}, \xi) + \mathbf{a}(\mathbf{x}, \xi)\}$$

$\mathbf{a} \perp \partial_\xi \log q$ : efficient

# Hyvärinen score

$$S(x, q) = \ddot{l}(x, \xi) + \frac{1}{2} \{ \dot{l}(x, \xi) \}^2$$

$$D_S[p : q] = \frac{1}{2} E_p \left[ \frac{d}{dx} \{ \log p(x) - \log q(x, \theta) \}^2 \right]$$

$$D[p : cq] = D[p : q]$$

$$s(x, \xi) = \partial_\xi \ddot{l}(x, \xi) + \dot{l}(x, \xi) \partial_\xi \dot{l}(x, \xi)$$

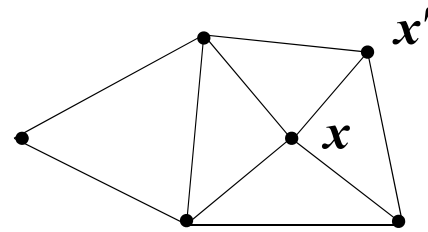
# Hyvärinen estimator

Fisher efficient  $\longleftrightarrow$   $q$  multivariate Gaussian

**Discrete case :  $\mathbf{x} \in$  graph nodes**

$$\Delta f(\mathbf{x}) = \frac{1}{N_x} \sum_{\mathbf{x}' \in N} \{f(\mathbf{x}) - f(\mathbf{x}')\}$$

$$N_x = \text{const}$$



$$S(\mathbf{x}, q) = \left\{ \frac{\Delta q(\mathbf{x})}{q(\mathbf{x})} \right\}^2 - 2\Delta \left\{ \frac{\Delta q(\mathbf{x})}{q(\mathbf{x})} \right\}$$

$$D_S[\xi : \xi'] = E_{q(\mathbf{x}, \xi)} \left\{ \frac{\Delta q(\mathbf{x}, \xi)}{q(\mathbf{x}, \xi)} - \frac{\Delta q(\mathbf{x}, \xi')}{q(\mathbf{x}, \xi')} \right\}^2$$

$$\sum_x \Delta f(\mathbf{x}) h(\mathbf{x}) = \sum_x f(\mathbf{x}) \Delta' h(\mathbf{x})$$

$$f, h \rightarrow 0$$

$$\int f'(x) h(x) dx = - \int f(x) h'(x) dx$$

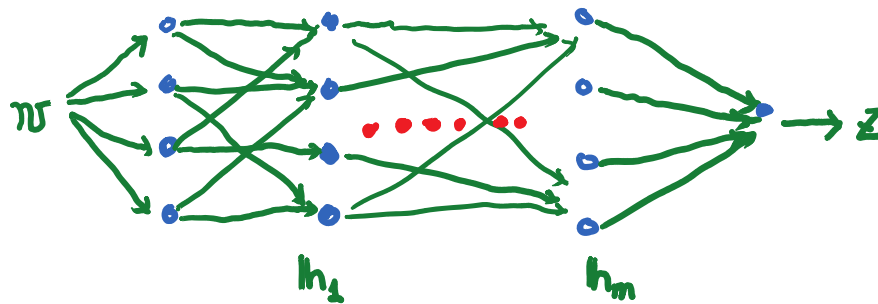
$$x \rightarrow \pm\infty$$

# Deep Learning

## Self-Organization + Supervised Learning

RBM: Restricted Boltzmann Machine

Auto-Encoder, Recurrent Net



tricks!!  
ideas!

Dropout

Contrastive divergence

bi-directional

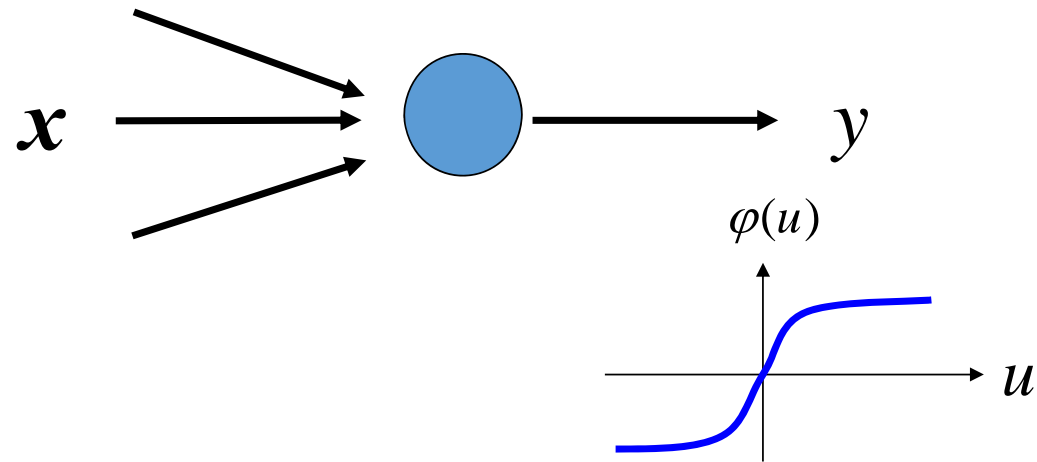
convolution

.



# Mathematical Neurons

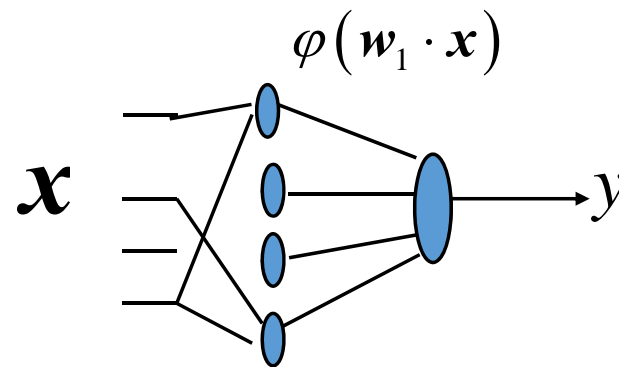
$$y = \varphi\left(\sum w_i x_i - h\right) = \varphi(\mathbf{w} \cdot \mathbf{x})$$



# Multilayer Perceptrons

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

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



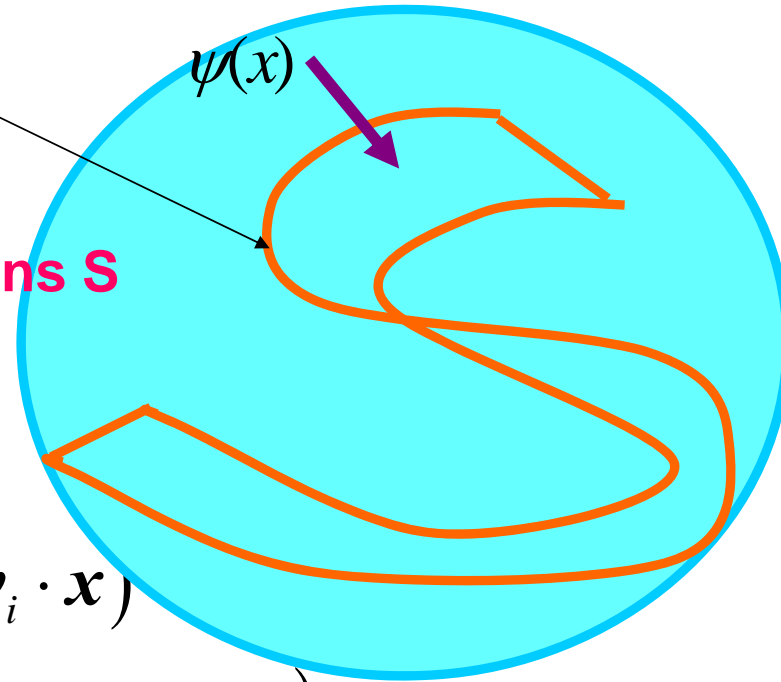
$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum v_i \varphi(\mathbf{w}_i \cdot \mathbf{x})$$

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

# Multilayer Perceptron

neuromanifold

space of functions  $\mathcal{S}$



$$y = f(\mathbf{x}, \boldsymbol{\theta})$$

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

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

## Backpropagation --- stochastic gradient learning

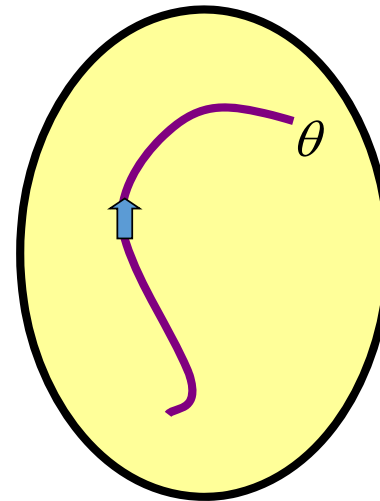
examples :  $(y_1, \mathbf{x}_1), \dots, (y_t, \mathbf{x}_t)$  --- training set

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

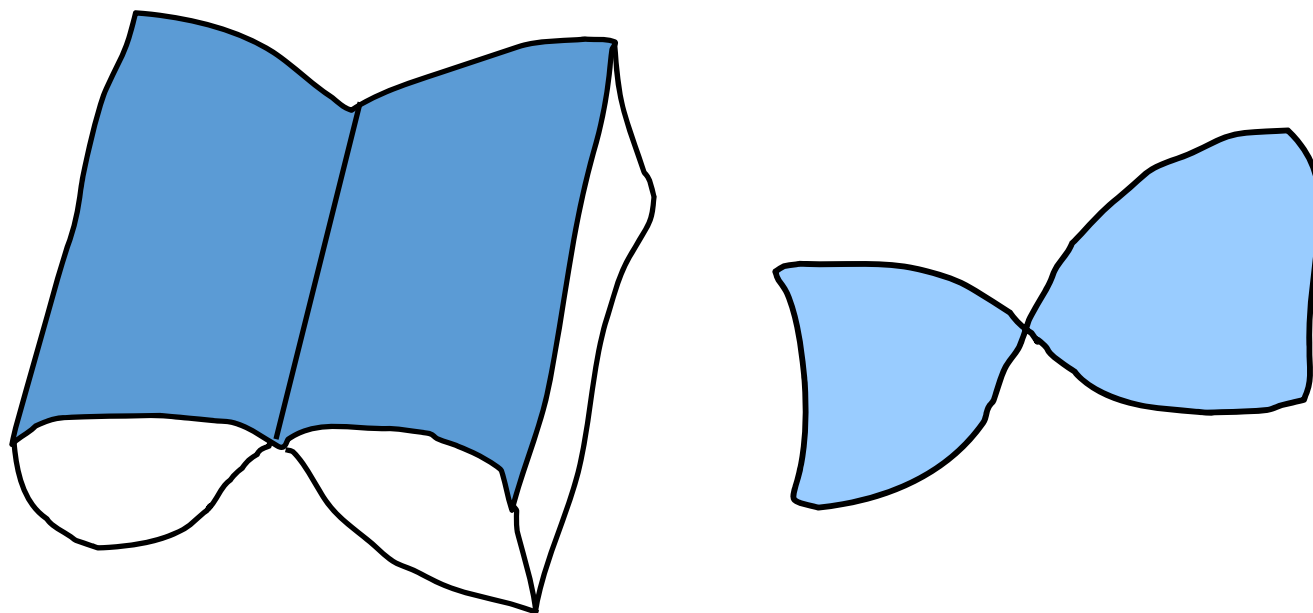
$$= -\log p(y, \mathbf{x}; \theta)$$

$$\Delta \theta_t = -\eta_t \frac{\partial l(y_t, x_t; \theta_t)}{\partial \theta}$$

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



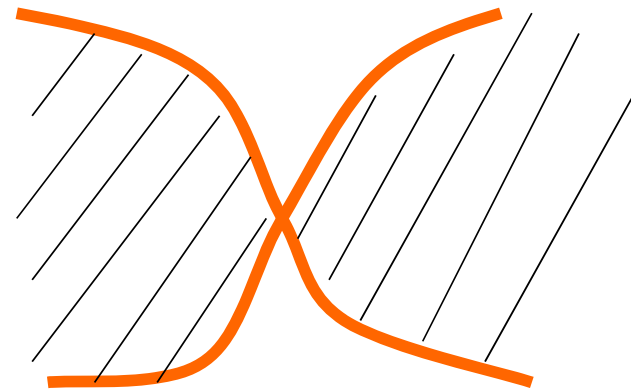
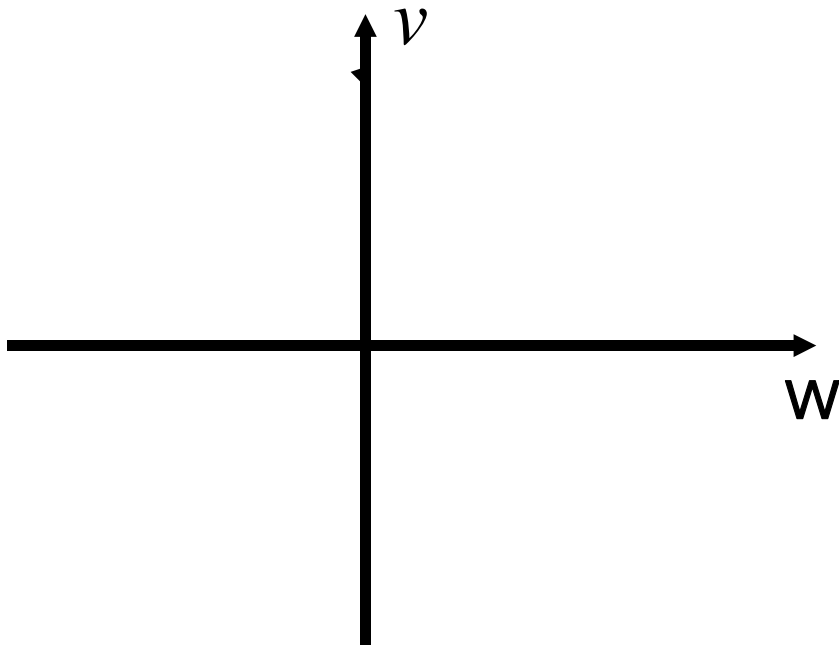
# singularities



## Geometry of singular model

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

$$\mathbf{v} \perp \mathbf{w} \Rightarrow 0$$



# model: 2 hidden neurons

$$f(\mathbf{x}, \boldsymbol{\theta}) = w_1 \varphi(\mathbf{J}_1 \cdot \mathbf{x}) + w_2 \varphi(\mathbf{J}_2 \cdot \mathbf{x})$$

$$y = f(\mathbf{x}, \boldsymbol{\theta}) + \varepsilon$$

$$\varphi(u) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^u e^{-\frac{t^2}{2}} dt$$

**loss function:**  $l(\mathbf{x}, y; \boldsymbol{\theta}) = \frac{1}{2} \{y - f(\mathbf{x}, \boldsymbol{\theta})\}^2$

$y$ : teacher signal :  $\boldsymbol{\theta}_0$  **stochastic descent learning**

$\dot{\boldsymbol{\theta}} = -\eta \left\langle \frac{\partial l(\mathbf{x}_t, y_t, \boldsymbol{\theta}_t)}{\partial \boldsymbol{\theta}} \right\rangle$  **backprop : vanilla gradient**



## Natural Gradient Stochastic Descent

$$\dot{\boldsymbol{\theta}} = -\eta G^{-1}(\boldsymbol{\theta}_t) \langle \nabla_{\boldsymbol{\theta}}(\mathbf{x}_t, y_t, \boldsymbol{\theta}_t) \rangle$$

$$\nabla_{\boldsymbol{\theta}} = \frac{\partial}{\partial \boldsymbol{\theta}}$$

$$G(\boldsymbol{\theta}) = \langle \nabla_{\boldsymbol{\theta}} l \nabla_{\boldsymbol{\theta}} l \rangle \quad : \quad \text{Fisher Information Matrix}$$

**invariant; steepest descent**

# Natural gradient is superior

Steepest descent; invariant    Yan Ollivier

Fisher-efficient

Natural gradient is non-vanishing even in multiple layers

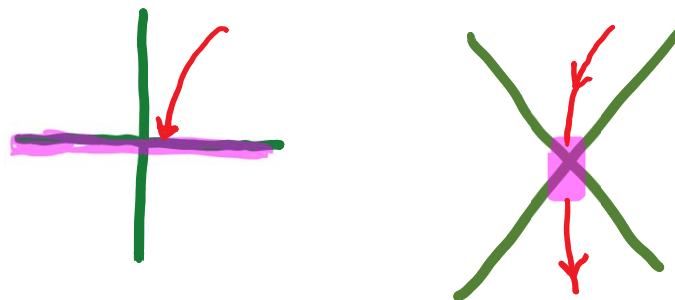
Good at singular regions (avoid plateaus: Milnor attractor)

## Adaptive Natural Gradient

$$G_{t+1}^{-1} = (1 + \varepsilon) G_t^{-1} - \varepsilon G_t^{-1} \nabla l(x_t) \nabla l(x_t)^\top G_t^{-1}$$

$G^{-1} \rightarrow \infty$ ,  $\nabla l \rightarrow 0$  at singularities

$G^{-1} \nabla l$



## Singular Region in Parameter Space

$$R(w, \mathbf{J}) = \{\boldsymbol{\theta} \mid \mathbf{J}_1 = \mathbf{J}_2 = \mathbf{J}, w_1 + w_2 = w\}$$

$$\cup \{\boldsymbol{\theta} \mid w_1 = 0, w_2 = w, \mathbf{J}_2 = \mathbf{J}\}$$

$$\cup \{\boldsymbol{\theta} \mid w_1 = w, w_2 = 0, \mathbf{J}_1 = \mathbf{J}\}$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = w_1 \varphi(\mathbf{J}_1 \cdot \mathbf{x}) + w_2 \varphi(\mathbf{J}_2 \cdot \mathbf{x})$$

# Coordinate transformation

$$\mathbf{v} = \frac{w_1 \mathbf{J}_1 + w_2 \mathbf{J}_2}{w_1 + w_2},$$

$$w = w_1 + w_2,$$

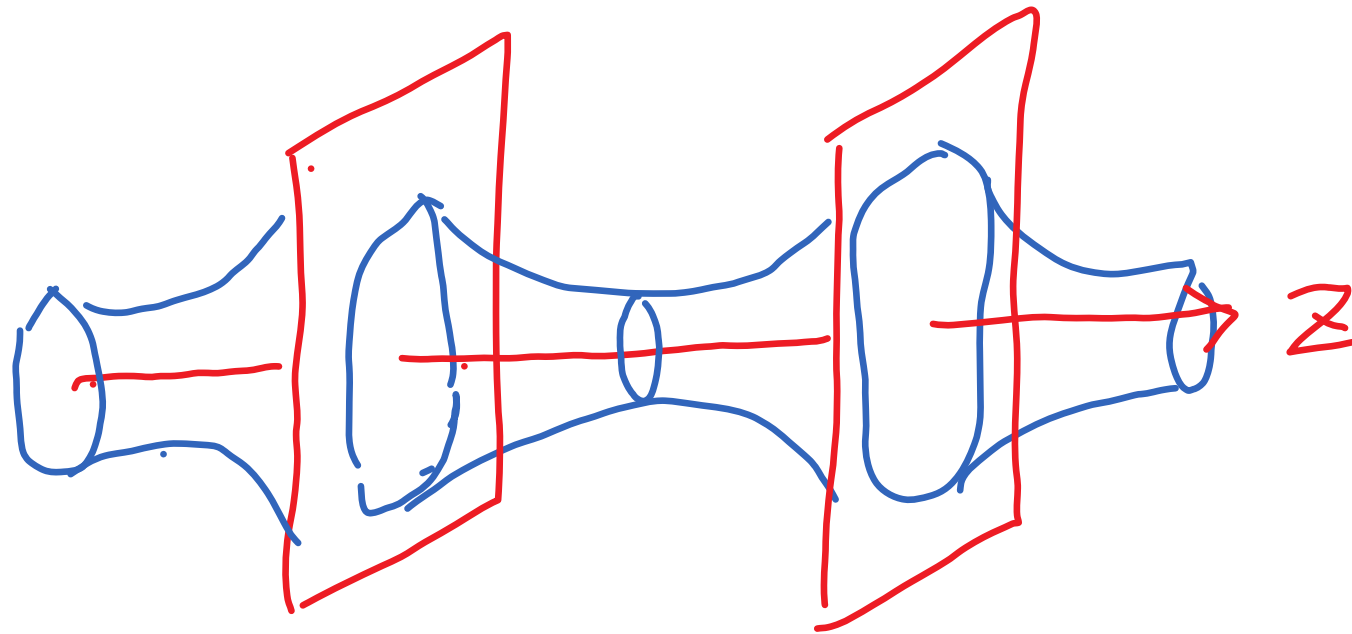
$$\mathbf{u} = \mathbf{J}_2 - \mathbf{J}_1,$$

$$z = \frac{w_2 - w_1}{w_1 + w_2}$$

$$\boldsymbol{\xi} = (\mathbf{v}, w, \mathbf{u}, z)$$

## Singular Region

$$R(w, \mathbf{J}) = \{\mathbf{u} = 0\} \cup \{z = \pm 1\}$$



## *Milnor attractor*

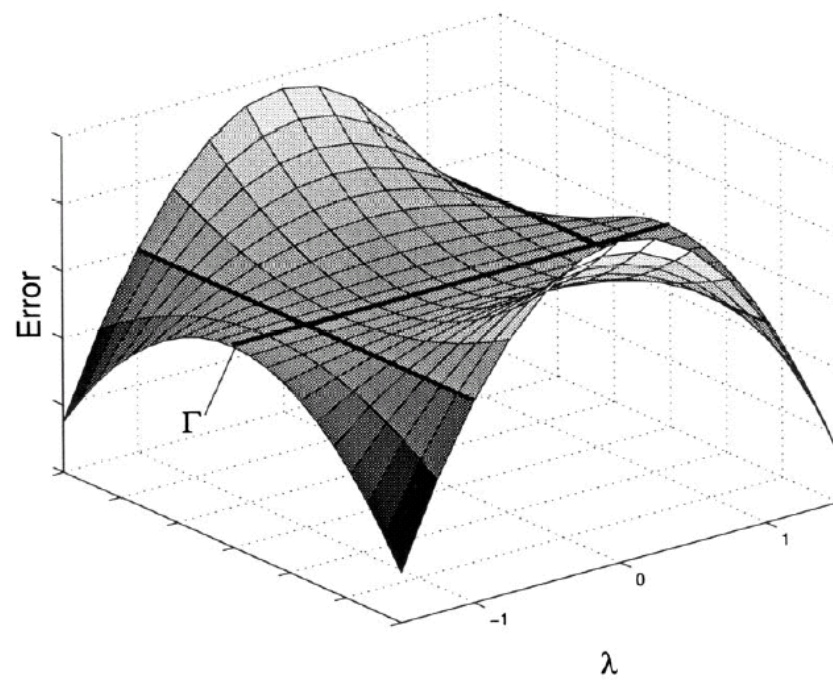


Fig. 5. Critical set with local minima and plateaus.

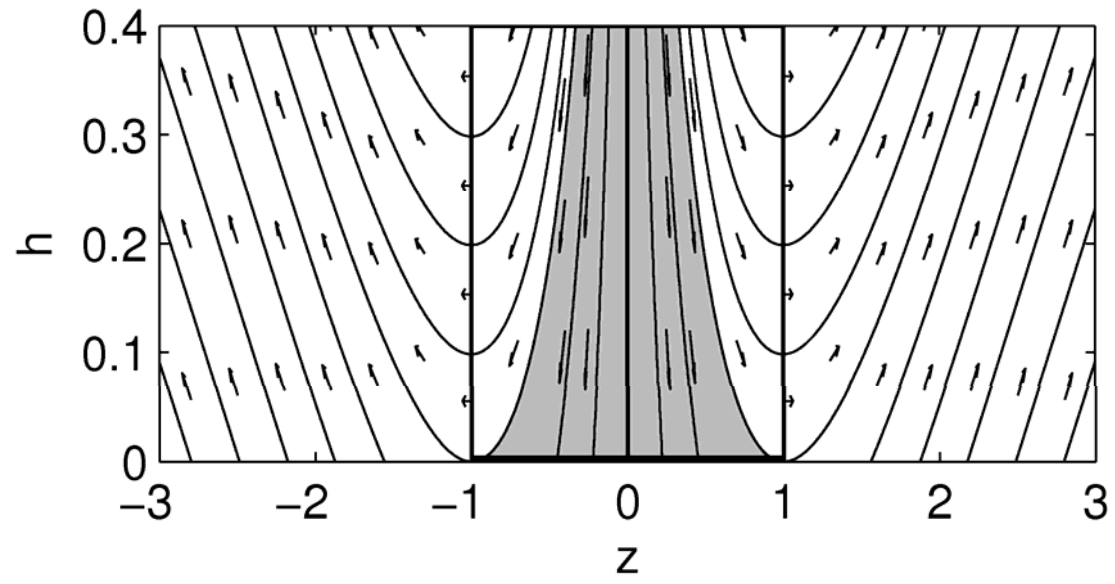
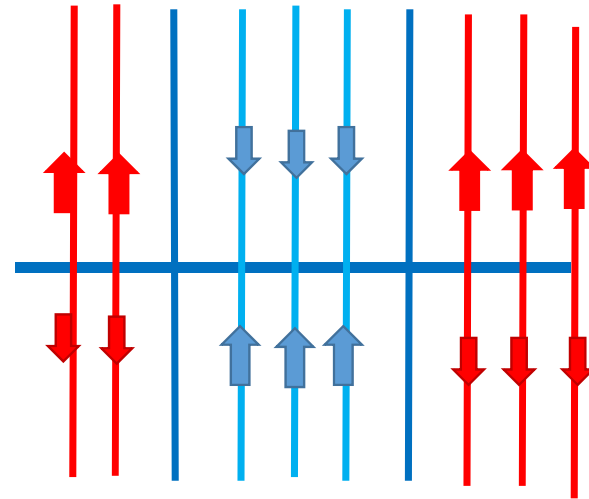
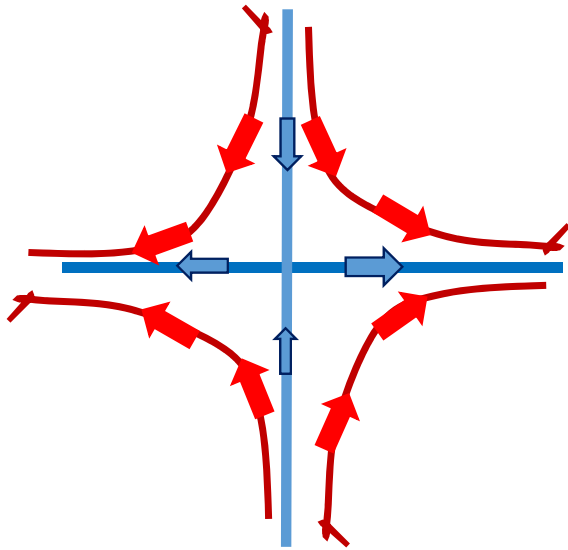


Fig. 2: trajectories



# Saddle and plateau



# Topology of singular R

blow-down coordinates :  $\alpha = (\tau, \sigma, \mathbf{e})$

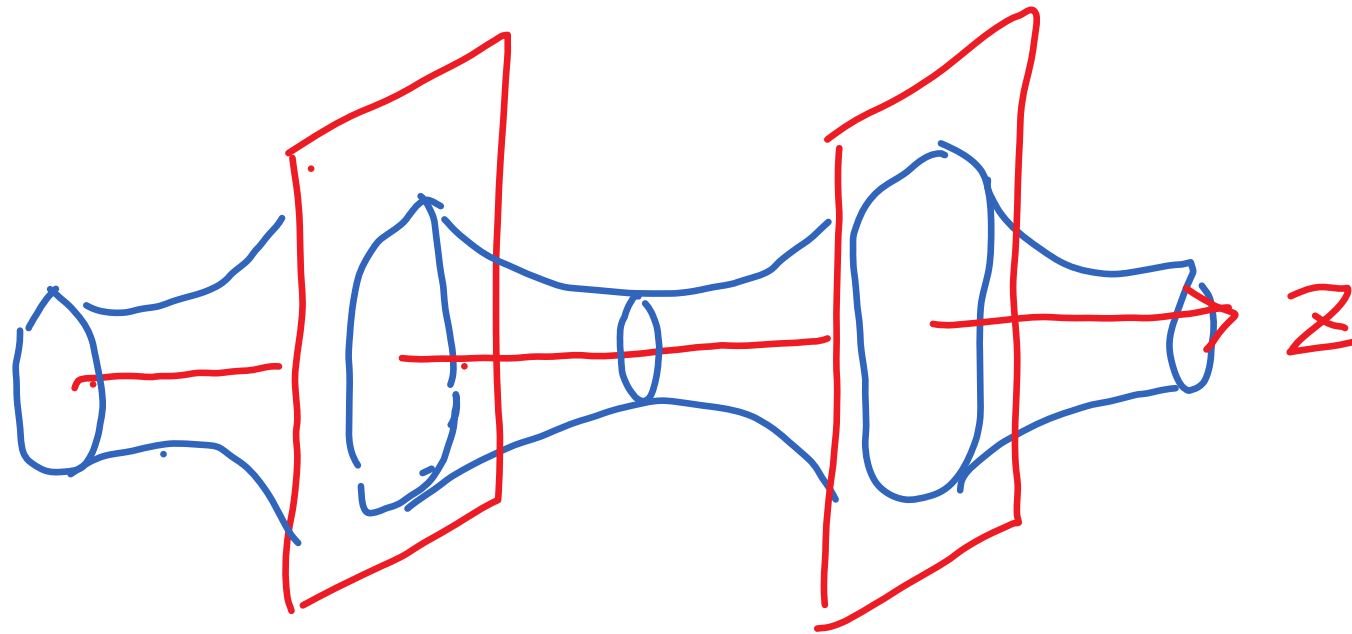
$$\tau = c_1 (1 - z^2) u^2, \quad u = |\mathbf{u}|$$

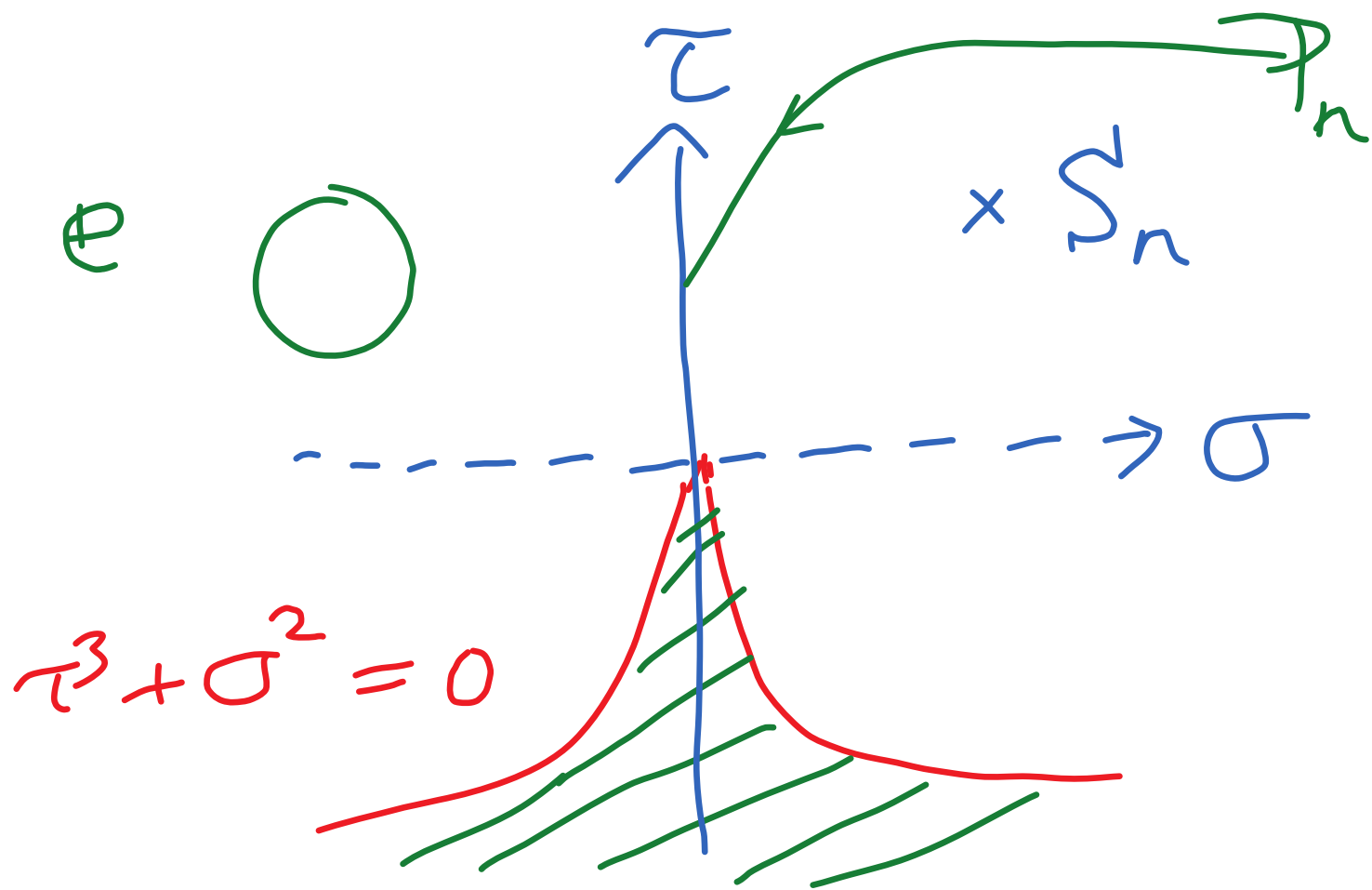
$$\sigma = c_2 z (1 - z^2) u^3,$$

$$\mathbf{e} = \frac{\mathbf{u}}{|\mathbf{u}|} \in S_n, \quad |\mathbf{e}| = 1$$

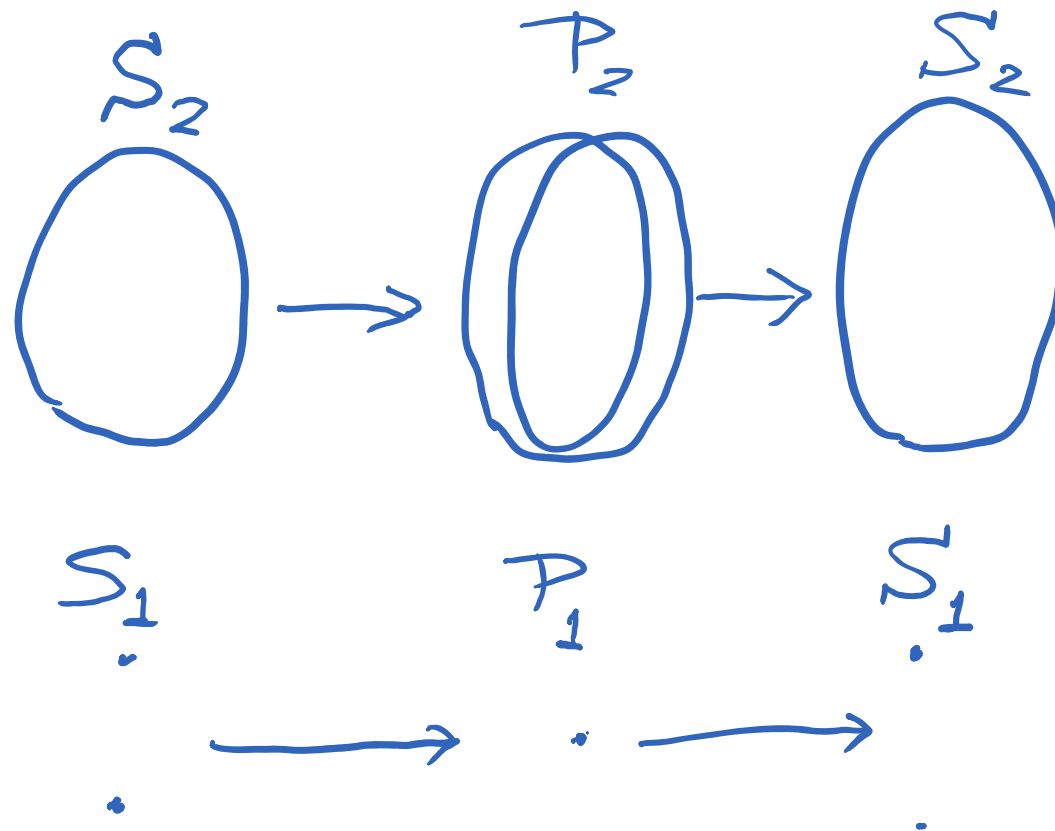
## Singular Region

$$R(w, \mathbf{J}) = \{\mathbf{u} = 0\} \cup \{z = \pm 1\}$$





# Sphere $S_n$ and Projective space $P_n$



# natural gradient learning near singularity

$$\frac{d}{dt} \begin{pmatrix} \tau \\ \sigma \end{pmatrix} = -\eta \begin{pmatrix} \tau \\ \sigma \end{pmatrix} \quad : \quad \text{true model} \in R$$

$$\frac{d}{dt} \begin{pmatrix} \tau \\ \sigma \end{pmatrix} = O(1) \quad : \quad \text{true model} \notin R$$

Milnor attractor

# Canonical Divergence in Manifold of Dual Affine Connections

*Nihat Ay and S. Amari*

# Divergence and metric

$$D[p : q] \geq 0$$

$$D[\xi : \xi + d\xi] = \frac{1}{2} g_{ij}(\xi) d\xi^i d\xi^j + O(|d\xi|^3)$$

**$G$  : Riemannian metric, positive-definite**



# Divergence and dual affine connections

$$\Gamma_{ijk} \sim \nabla \quad \Gamma_{ijk}^* \sim \nabla^*$$

$$\Gamma_{ijk} = -\partial_i \partial_j \partial'_k D[\xi : \xi']_{\xi'=\xi}$$

$$\Gamma_{ijk}^* = -\partial'_i \partial'_j \partial_k D[\xi : \xi']_{\xi'=\xi}$$

$$\partial_i = \frac{\partial}{\partial \xi^i}; \quad \partial'_j = \frac{\partial}{\partial \xi'^j}$$

# Dual geometry

$$\{M, g, \nabla, \nabla^*\}$$

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

$$\{M, g, T\}, \quad T_{ijk} = \Gamma_{ijk}^* - \Gamma_{ijk}$$

$$\Gamma_{ijk}^{\pm\alpha} = \Gamma_{ijk}^o \mp \frac{\alpha}{2} T_{ijk} \quad \Gamma^o: \text{Levi-Civita connection}$$

## Dual geometry $\rightarrow$ canonical divergence

$M$  : dually flat :  $\exists \psi(\boldsymbol{\theta}), \varphi(\boldsymbol{\eta})$

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

**Bregman divergence**

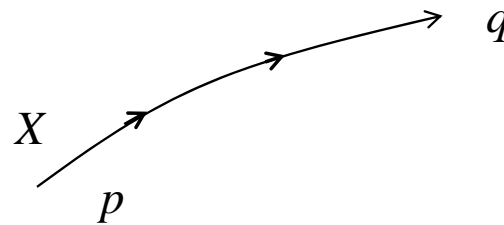
# Exponential map : $\xi(t)$ geodesic

$$\nabla_{\dot{\xi}} \dot{\xi} = 0$$

$$\xi(0) = p$$

$$\xi(1) = q$$

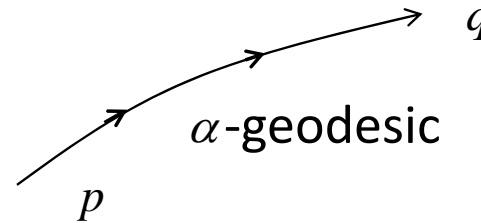
$$\dot{\xi}(0) = X = \log_p q$$



# Exponential map divergence

$$D[p : q] = \|X(p : q)\|^2$$

$\alpha$ -divergence



$$D_\alpha[p : q] = \|X_\alpha(p : q)\|^2$$

**Theorem 1.** Exponential map divergence induces  $\alpha = -3$  geometry

**Theorem 2.**  $\alpha = -\frac{1}{3}$  exponential map divergence recovers the original geometry

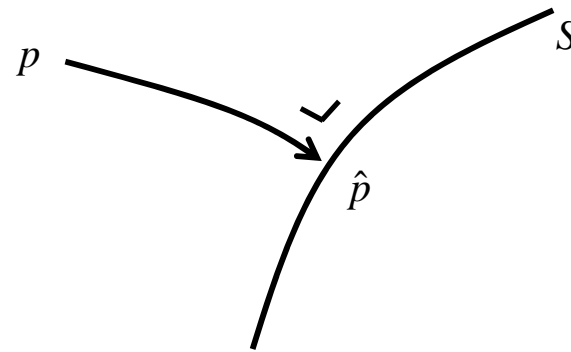
**Standard divergence:**  $D_{\text{stan}} [p : q] = \|X_{-1/3}(p, q)\|^2$

$$\begin{aligned} D[p : q] &= \int_0^1 \langle X_t(q, p), \dot{\xi}_{q,p}(t) \rangle dt \\ &= \int_0^1 t \|\dot{\xi}_{q,p}(t)\|^2 dt \end{aligned}$$

$$\int_0^1 w(t) \|\dot{\xi}_{q,p}(t)\|^2 dt$$

# Divergence and projection

$$\hat{p} = \arg \min_{q \in \mathcal{S}} D[p : q]$$



**projection theorem:**

$$X = c \operatorname{grad}_q D[p : q]$$



**IEEE ISIT-2011 Sankt Petersburg**

**Data Compression  
in Multiterminal Statistical Inference**

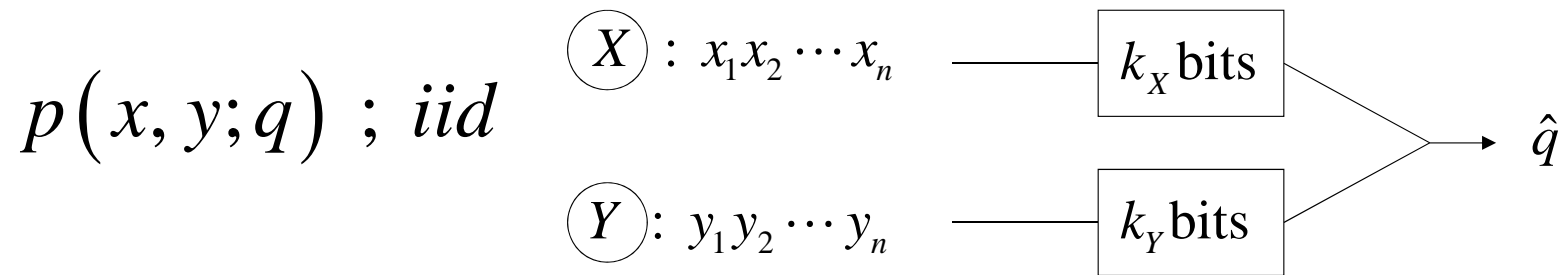
**Shun-ichi Amari  
RIKEN Brain Science Institute**

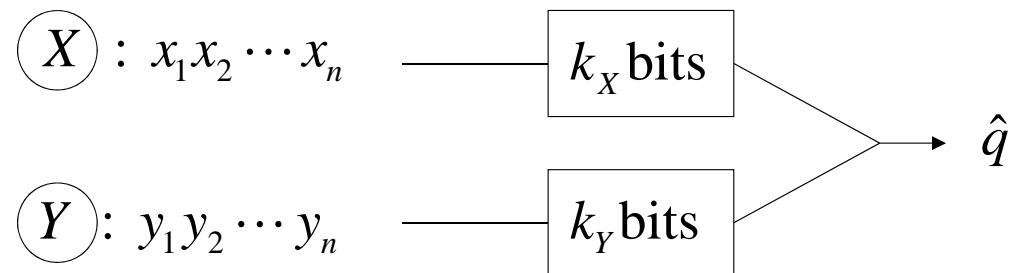
# A long standing problem

T.Berger; Csiszar, Ahlswede, Burnashev, Han, Amari

correlated sources  $X, Y$

data compression and statistical inference

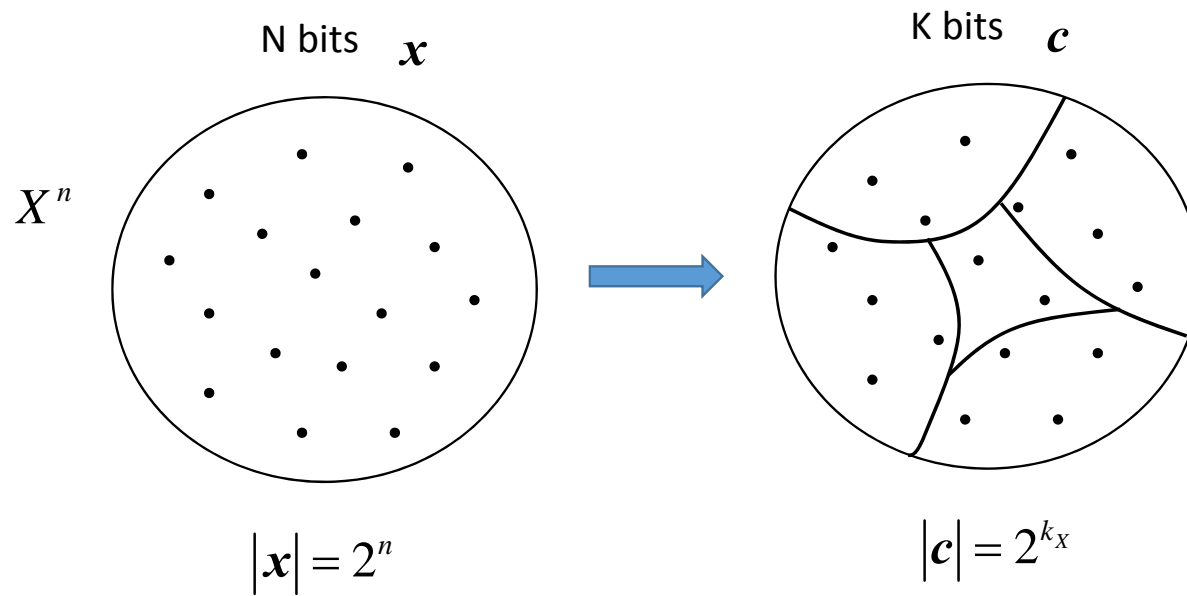




$$p(x, y; q) ; iid \qquad \text{Prob} \{x = y\} = q$$

$$\text{binary} : x, y = 0, 1 ; \text{Prob} \{x = 1\} = \text{Prob} \{y = 1\} = \frac{1}{2}$$

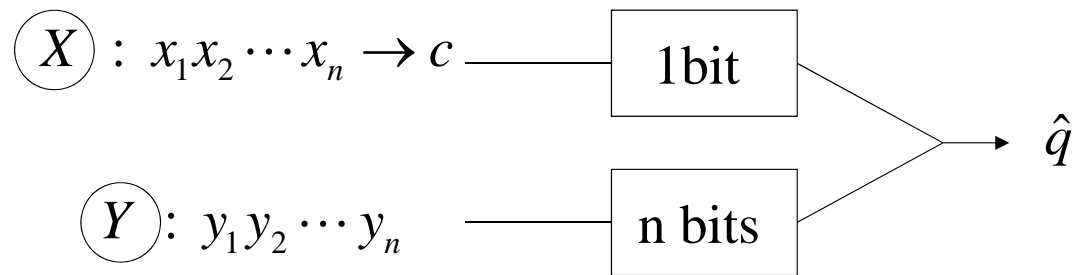
# Encoding : data compression



## One-bit helper case

$$k_X = 1, \quad k_Y = n \quad c = \text{sgn}(\mathbf{a} \cdot \mathbf{x})$$

$c$        $\mathbf{y}$



# Is single-bit encoding optimal?

It is optimal

when  $q = 1/2$  ( $x, y$  independent),

but not for general  $q$ .

# Fisher information: $k_X = 1, k_Y = n$

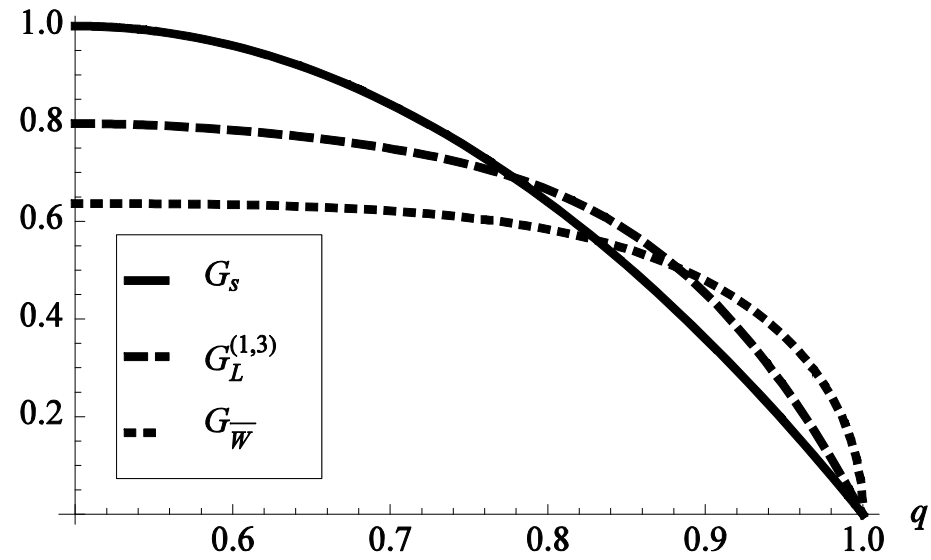
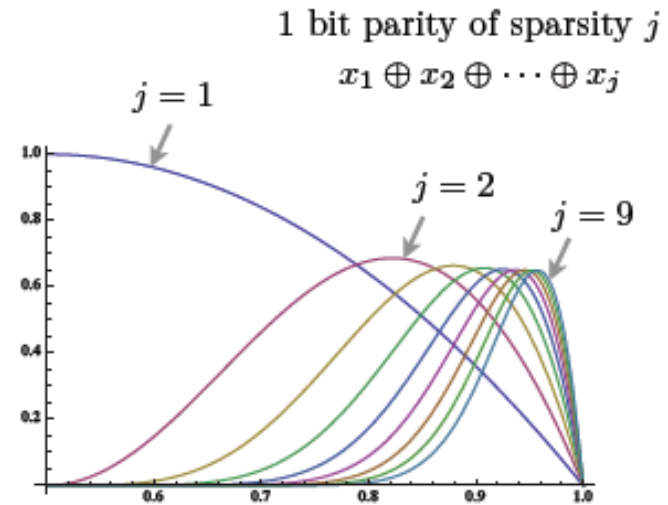


Fig. 2

# Kingo Kobayashi: parity encoding

$$x_1 \oplus x_2 \oplus \dots \oplus x_s$$

in progress !





# Information Geometry and Transportation Problem (Wasserstein distance)

entropic relaxation :  $\min \langle c, p \rangle - \alpha \{-H(p)\}$  dual

New Paper

S. Pal and T-K L. Wong,

Exponentially concave function and a new information geometry

Portfolio theory, transportation problem and information geometry  
(dually projectively flat)