

Theoretical Sciences

Information Geometry and Its Applications

Shun-ichi Amari

RIKEN Brain Science Institute

Information Geometry:

A Unifying Framework

**Statistical Inference,
Information Sciences**

Signal Processing,

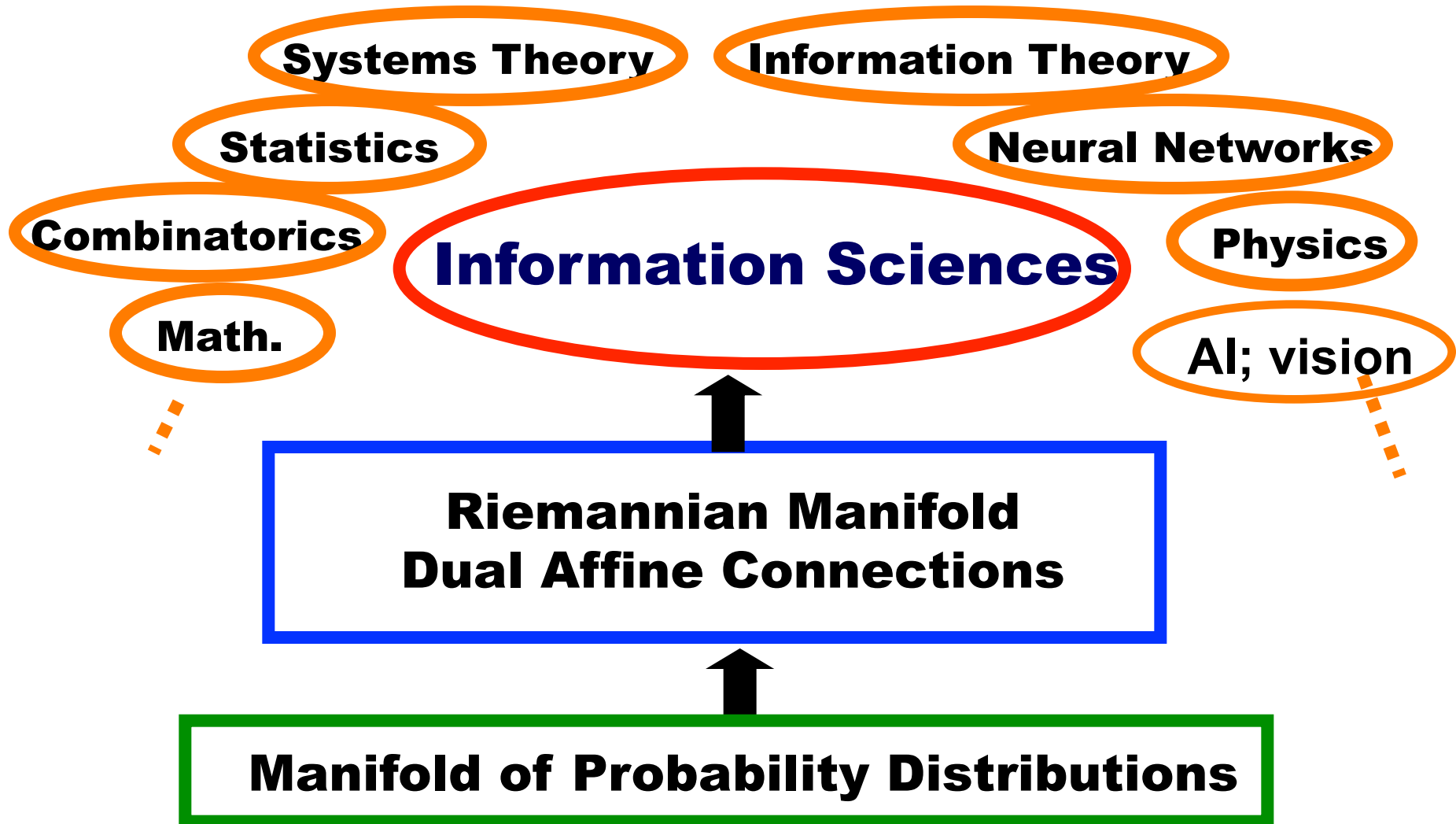
Machine Learning

Convex Analysis

Physics

Brain Science

Information Geometry



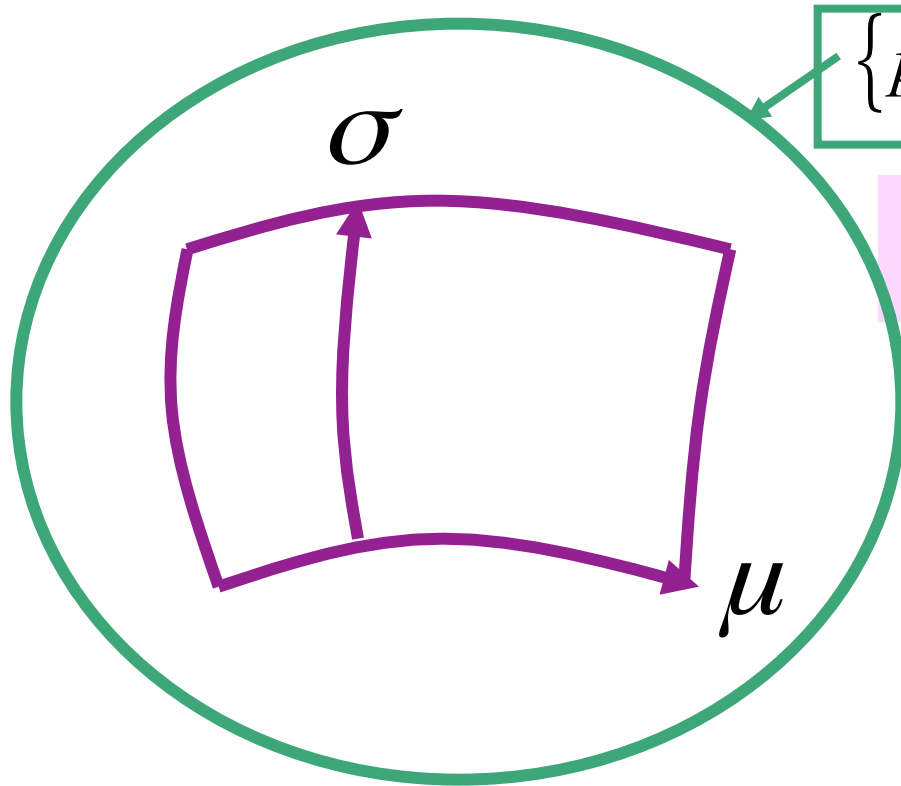
Information Geometry

**-- Manifolds of
Probability Distributions**

Information Geometry ?

Gaussian distributions

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



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

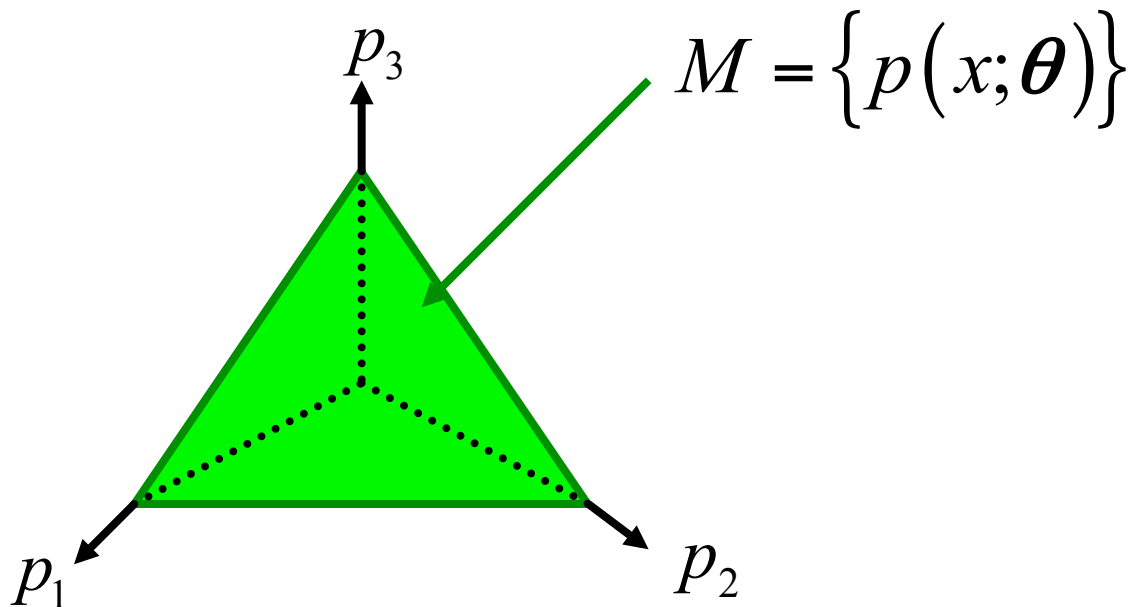
$$S = \{p(x; \theta)\}$$



Manifold of Probability Distributions

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

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



Invariance

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

1. *Invariant under reparameterization*

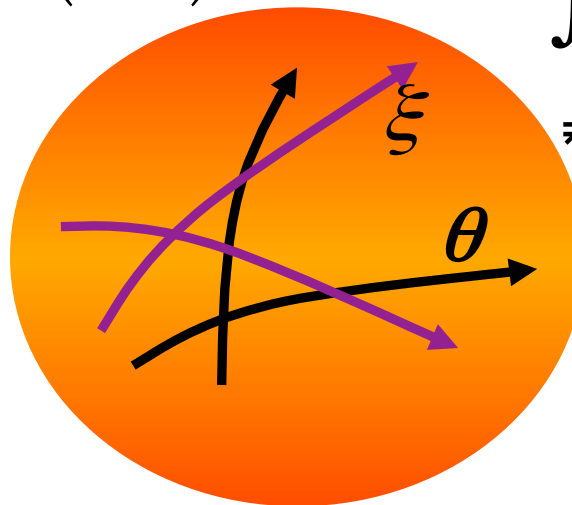
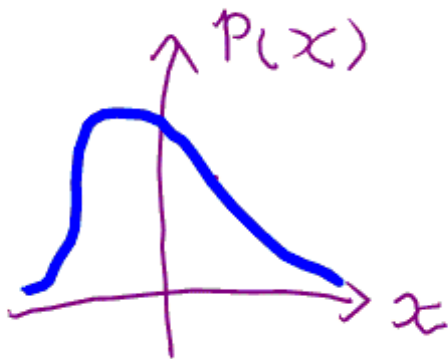
$$p(x, \theta) = \bar{p}(x, \xi) \quad D = \sum \theta_i^2 \neq \sum \xi_i^2$$

2. *Invariant under different representation*

$$y = y(x), \quad \bar{p}(y, \theta)$$

$$\int |p(x, \theta_1) - p(x, \theta_2)|^2 dx$$

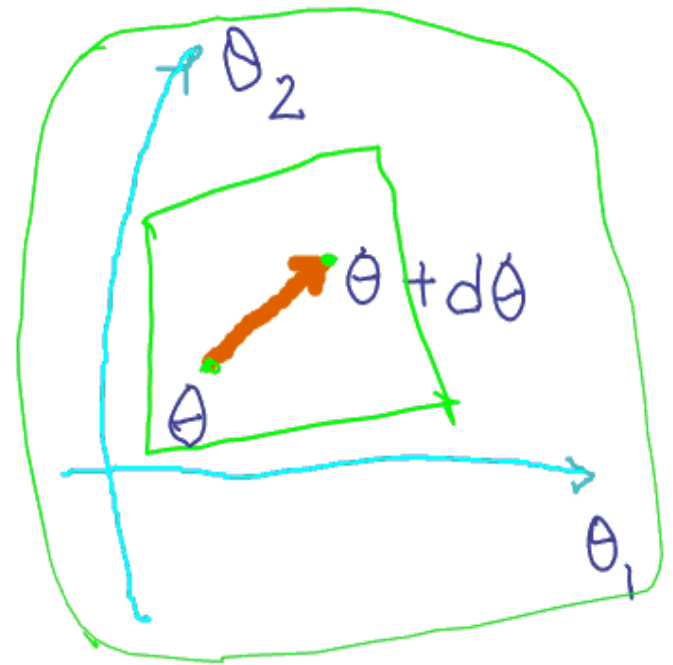
$$\neq \int |\bar{p}(y, \theta_1) - \bar{p}(y, \theta_2)|^2 dy$$



Two Structures

Riemannian metric

*affine connection ---
geodesic*



Riemannian Structure

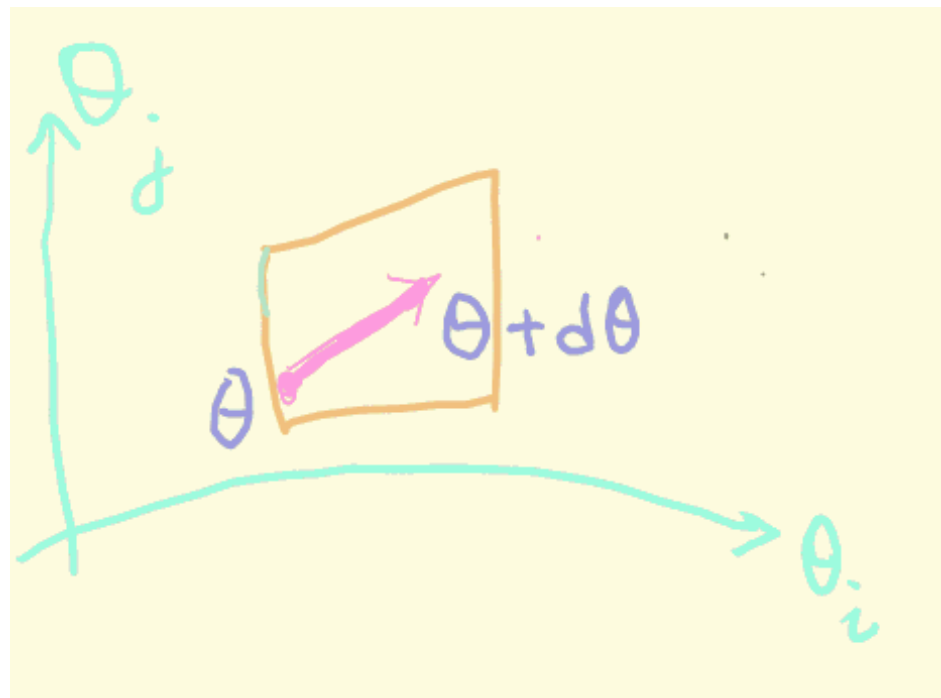
$$ds^2 = \sum g_{ij}(\theta) d\theta^i d\theta^j$$
$$= d\theta^T G(\theta) d\theta$$

$$G(\theta) = (g_{ij})$$

Euclidean $G = E$

Fisher information

$$g_{ij} = E \left[\frac{\partial}{\partial \theta_i} \log p \frac{\partial}{\partial \theta_j} \log p \right]$$



Affine Connection

covariant derivative; parallel transport

$$\nabla_X Y, \quad \Pi_c X = Y$$

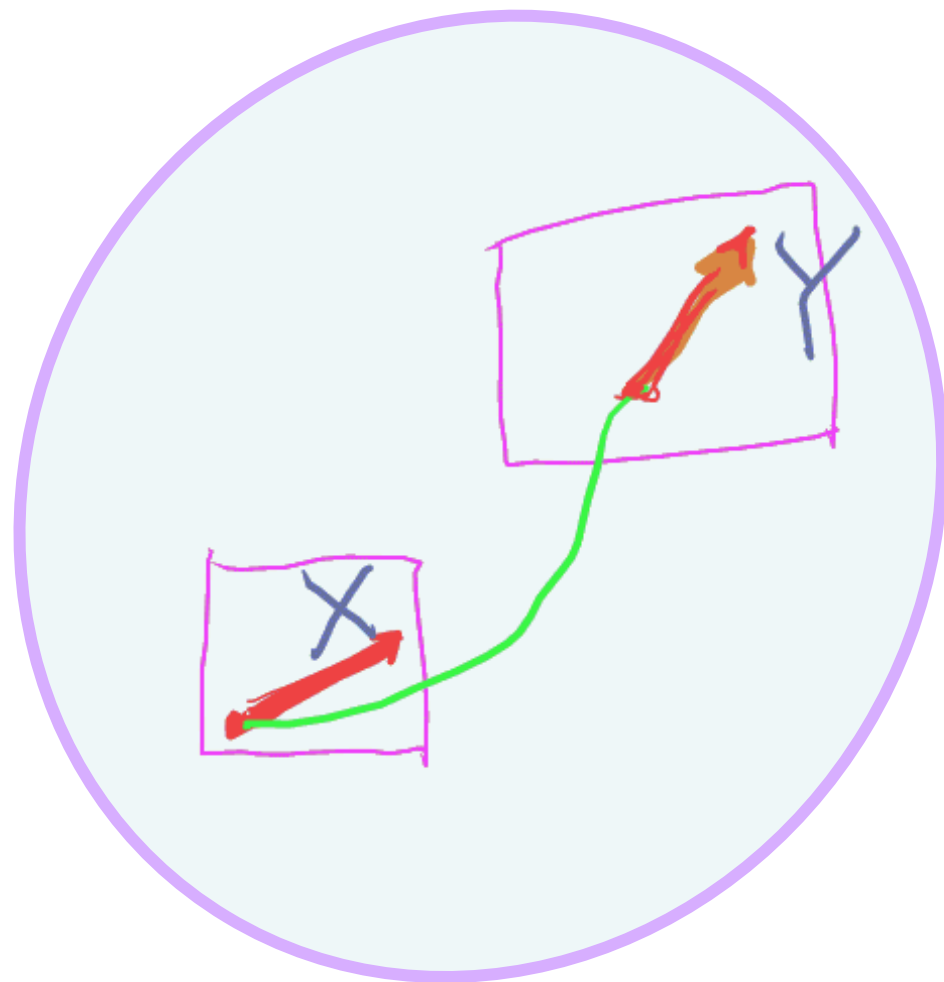
geodesic ~~$\Pi X = X$~~ $X = X(t)$

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

minimal distance

straight line

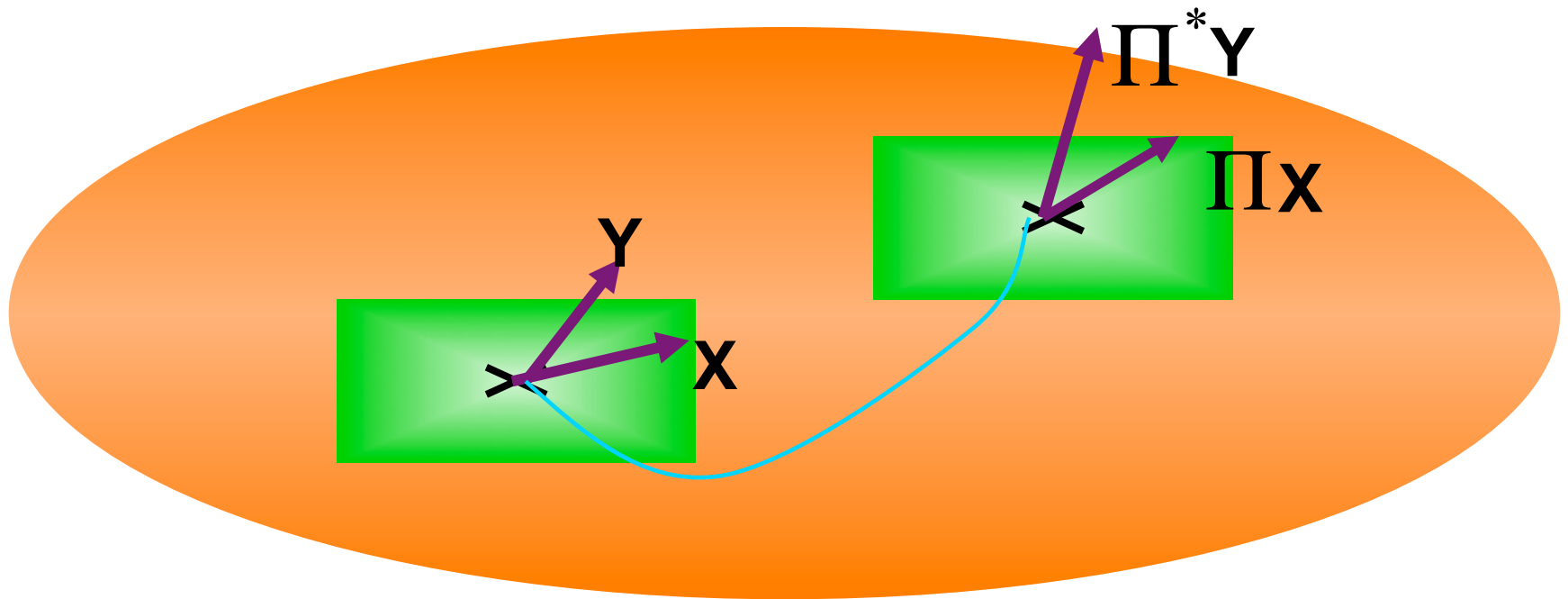
different concepts



Duality: two affine connections

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

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



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

Dual Affine Connections (∇, ∇^*)

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

e-geodesic

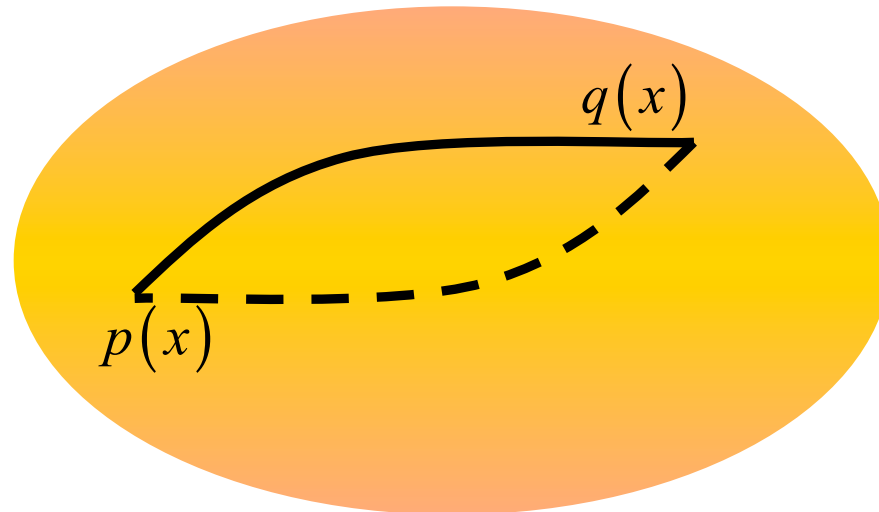
$$\log r(x, t) = t \log p(x) + (1-t) \log q(x) + c(t)$$

m-geodesic

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

$$\nabla_{\dot{x}} \dot{x}(t) = 0$$

$$\nabla^*_{\dot{x}} \dot{x}(t) = 0$$



Mathematical structure of $S = \{p(x, \xi)\}$

$$(S, g, T) \quad g_{ij}(\xi) = E[\partial_i l \partial_j l]$$
$$T_{ijk}(\xi) = E[\partial_i l \partial_j l \partial_k l]$$

$$l = \log p(x, \xi); \quad \partial_i = \frac{\partial}{\partial \xi^i}$$

α -connection

$$\Gamma_{ijk}^\alpha = \{i, j; k\} - \alpha T_{ijk}$$

$\nabla^\alpha \leftrightarrow \nabla^{-\alpha}$: dually coupled

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

Dually flat

Dually flat manifold

θ -coordinates \leftrightarrow η -coordinates

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

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

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

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

ψ : cumulant generating function

φ : negative entropy

$$\text{canonical divergence } D(P: P') = \psi(\theta) + \varphi(\eta') - \sum \theta_i \eta_i'$$

Information Geometry

-- Dually Flat Manifold

Convex Analysis

Legendre transformation

Divergence

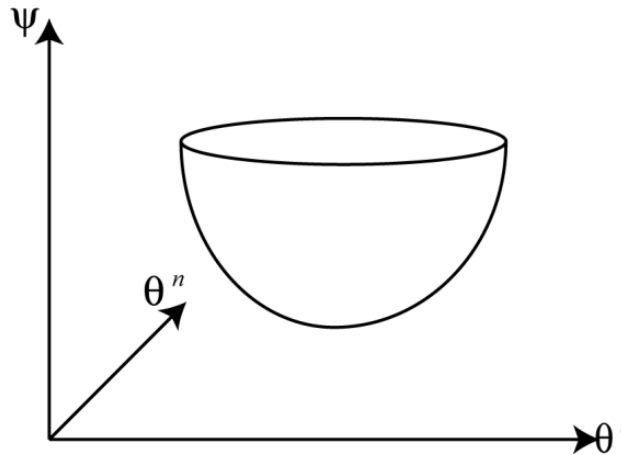
Pythagorean theorem

I-projection

Manifold with Convex Function

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

$\psi(\theta)$: convex function



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

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

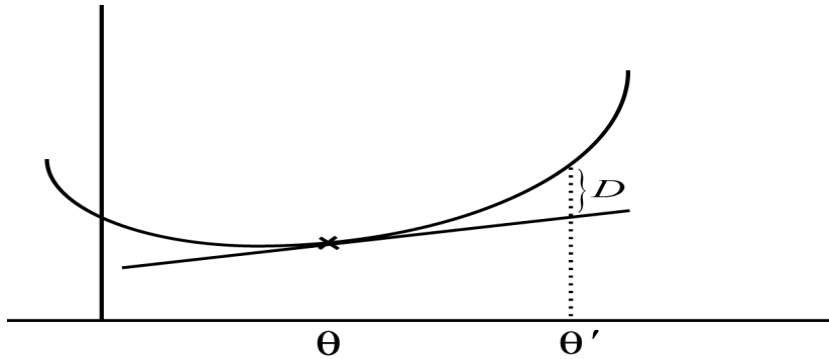
**negative entropy
energy**

**mathematical programming, control systems
physics, engineering**

Riemannian metric and flatness

Bregman divergence

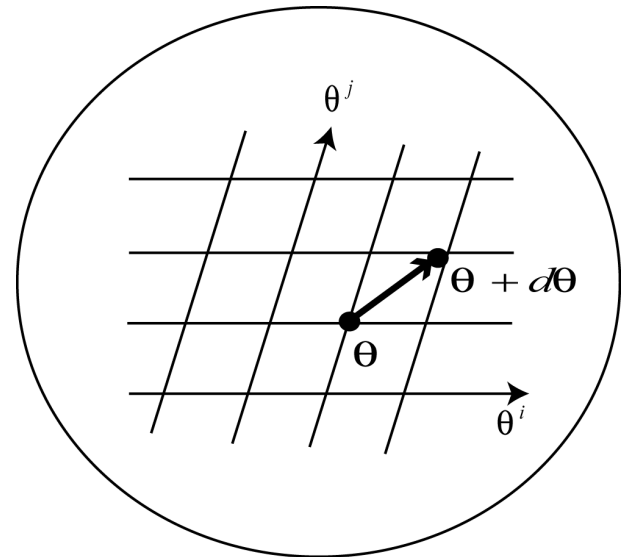
$$D(\boldsymbol{\theta}', \boldsymbol{\theta}) = \psi(\boldsymbol{\theta}') - (\boldsymbol{\theta}' - \boldsymbol{\theta}) \cdot \text{grad } \psi(\boldsymbol{\theta})$$



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

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

$\boldsymbol{\theta}$: geodesic



Legendre Transformation

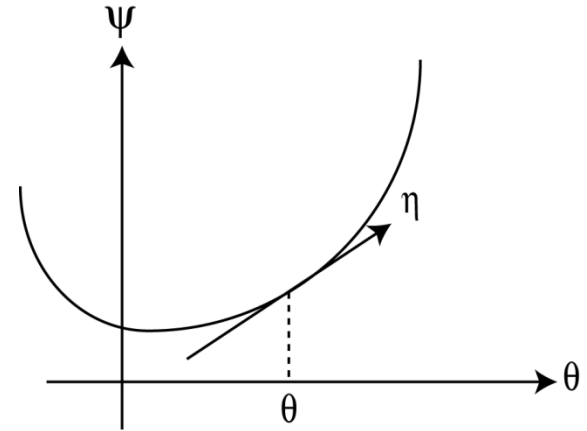
$$\eta_i = \partial_i \psi(\boldsymbol{\theta})$$

$\boldsymbol{\theta} \leftrightarrow \boldsymbol{\eta}$: one-to-one

$$\varphi(\boldsymbol{\eta}) + \psi(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \boldsymbol{\eta} = 0$$

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

Divergence $D(\boldsymbol{\theta}, \boldsymbol{\theta}') = \psi(\boldsymbol{\theta}) + \varphi(\boldsymbol{\eta}') - \boldsymbol{\theta} \cdot \boldsymbol{\eta}'$



Two coordinate systems of S

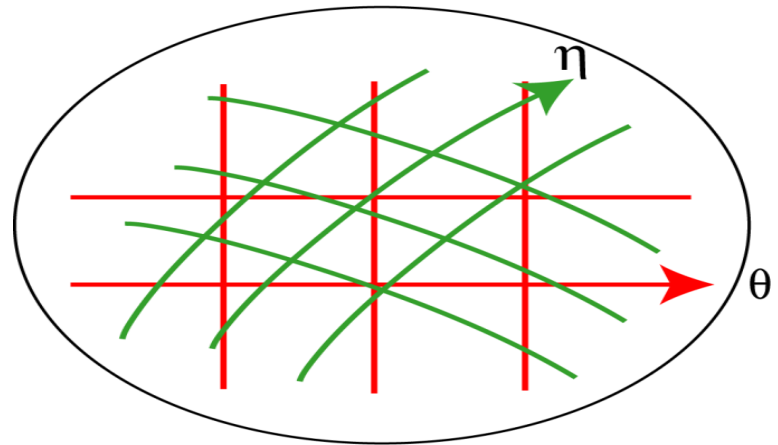
(θ, η)

θ : geodesic (e-geodesic)

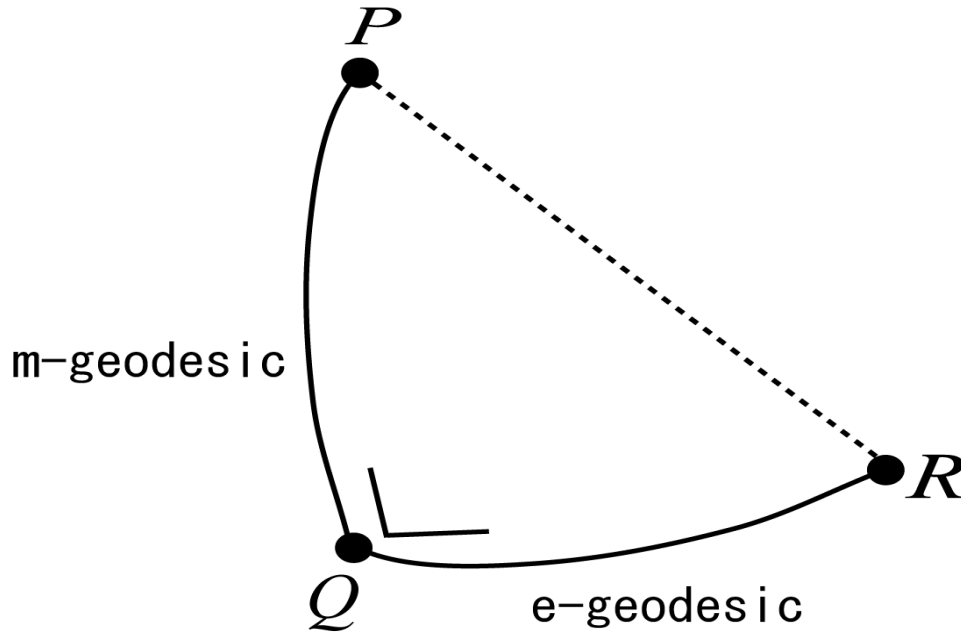
η : dual geodesic (m-geodesic)

“orthogonal”

$$\langle \partial_i, \partial^j \rangle = \delta_i^j$$



Pythagorean Theorem



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

Euclidean space: self-dual

$$\theta = \eta$$

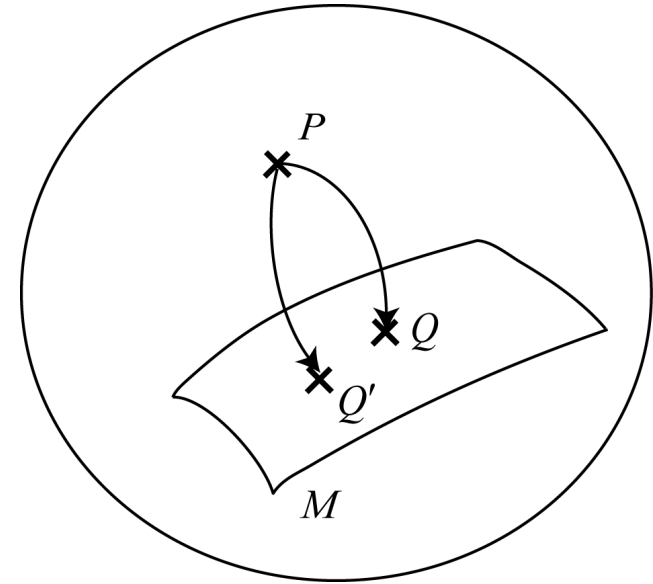
$$\psi(\theta) = \frac{1}{2} \sum (\theta_i)^2$$

Projection Theorem

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

$Q =$ **m-geodesic projection** of

P to M



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

$Q =$ **e-geodesic projection** of P to M

Information Geometry

Dually flat manifold; curved submanifold

convex potential functions $\psi(\boldsymbol{\theta}), \varphi(\boldsymbol{\eta})$

Euclidean space : self-dual

$$\psi(\boldsymbol{\theta}) = \frac{1}{2} \sum (\theta^i)^2, \quad \theta_i = \eta^i$$

Probability distributions $S = \{p(x)\}$

Exponential family : $p(x, \boldsymbol{\theta}) = \exp\{\boldsymbol{\theta} \cdot \boldsymbol{x} - \psi(\boldsymbol{\theta})\}$

$$\boldsymbol{\eta} = E[\boldsymbol{x}]$$

$\varphi(\boldsymbol{\eta})$: **negentropy**

Dually flat manifold

θ -coordinates \leftrightarrow η -coordinates

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

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

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

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

ψ : cumulant generating function

φ : negative entropy

Divergence

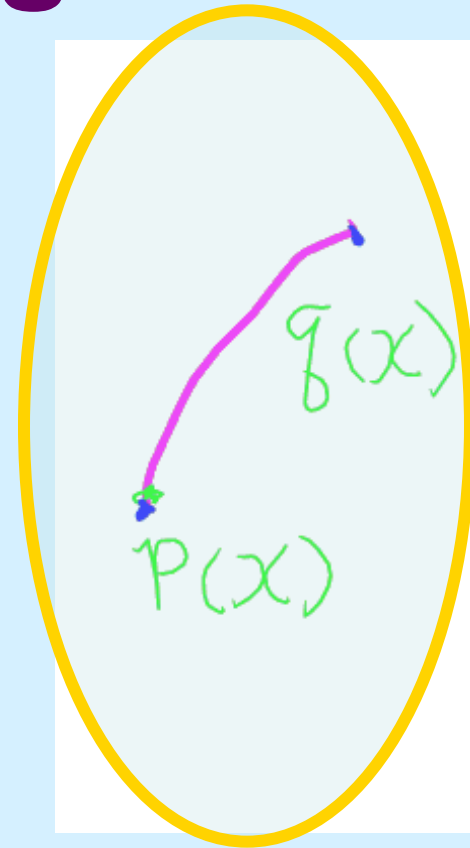
Kullback-Leibler Divergence

quasi-distance

$$D[p(x) : q(x)] = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

$$D[p(x) : q(x)] \geq 0 \quad =0 \text{ iff } p(x) = q(x)$$

$$D[p : q] \neq D[q : p]$$



Dually Flat Manifold

1. *Potential Functions*

---convex (Legendre transformation)

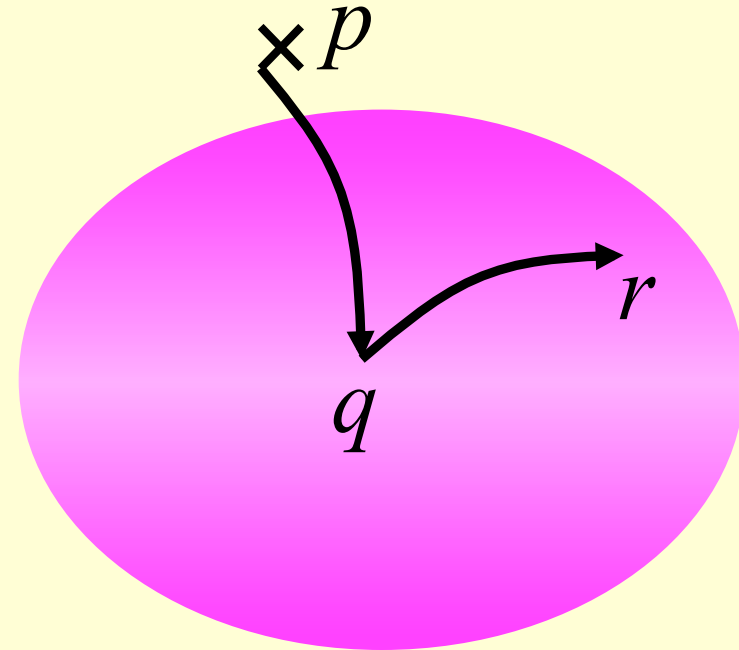
2. *Divergence* $D[p:q]$ Bregman divergence

3. *Pythagoras Theorem*

$$D[p:q] + D[q:r] = D[p:r]$$

4. *Projection Theorem*

5. *Dual foliation*



Applications to Statistics

curved exponential family:

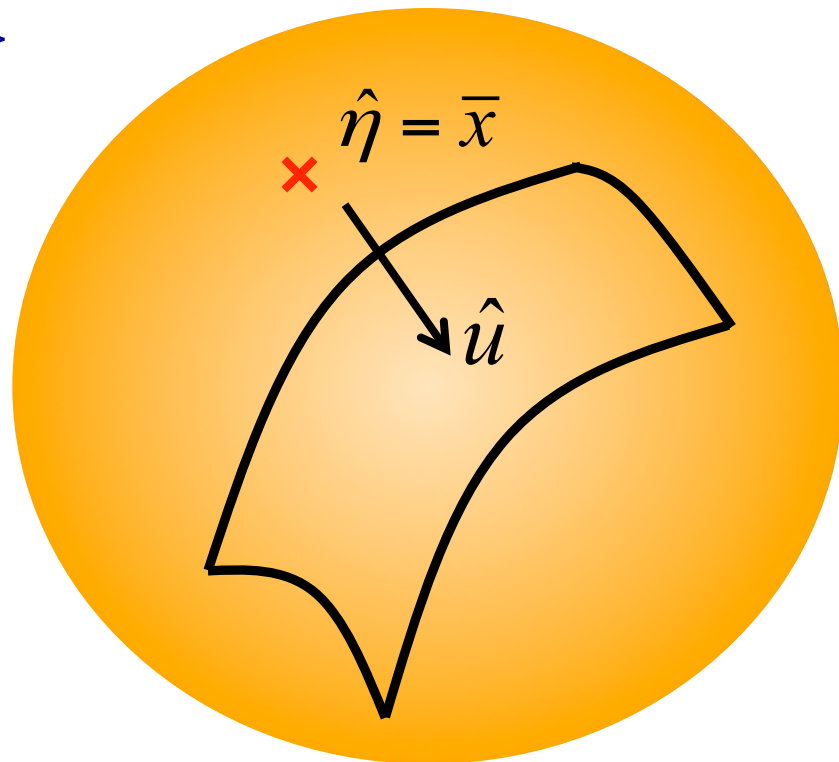
$$p(x, u) : x_1, x_2, \dots, x_n \quad p(x, \theta) = \exp \{ \theta \cdot x - \psi(\theta) \}$$

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

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

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

$H_0 : u = u_0$: testing



High-Order Asymptotics

$$p(x, \theta(u)) \quad : x_1, L, x_n$$

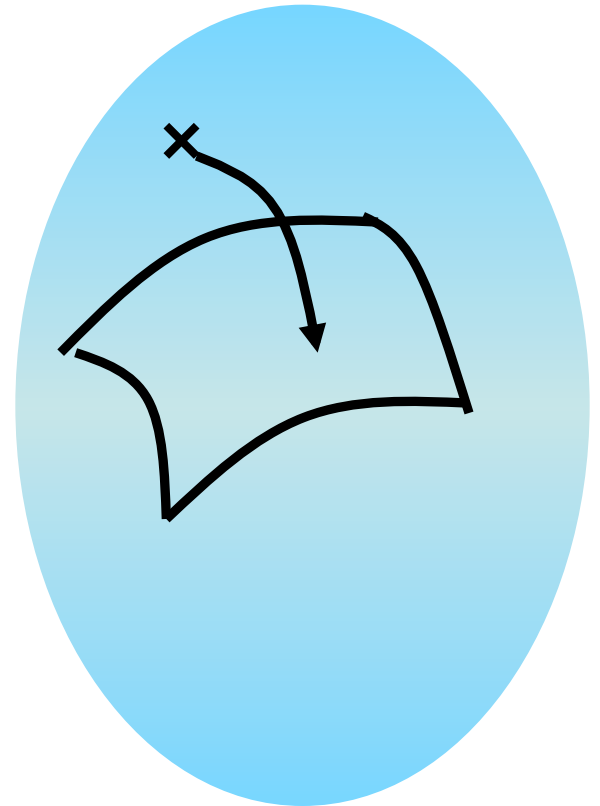
$$\hat{u} = u(x_1, L, x_n)$$

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

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

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

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

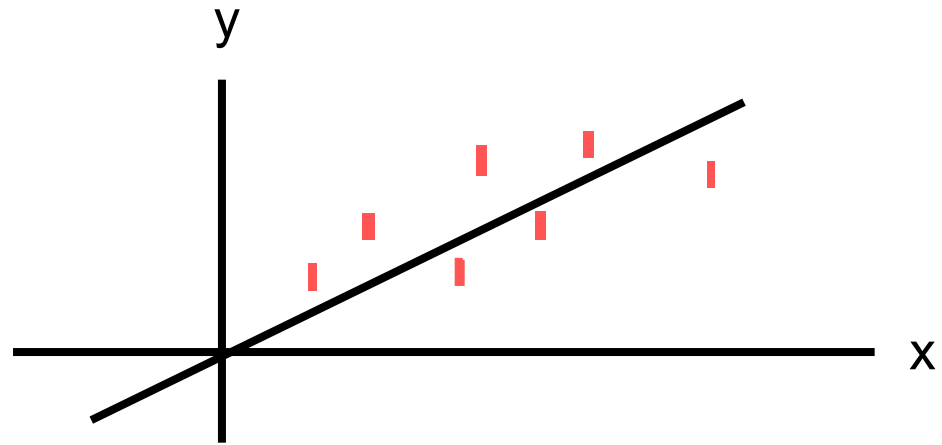


Semiparametric Statistical Model

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

linear relation

$$y = \theta x$$



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

mle, least square, total least square

Linear Regression: Semiparametrics

$$(x_1, y_1)$$

$$x_i = \xi_i + \varepsilon_i$$

$$(x_2, y_2)$$

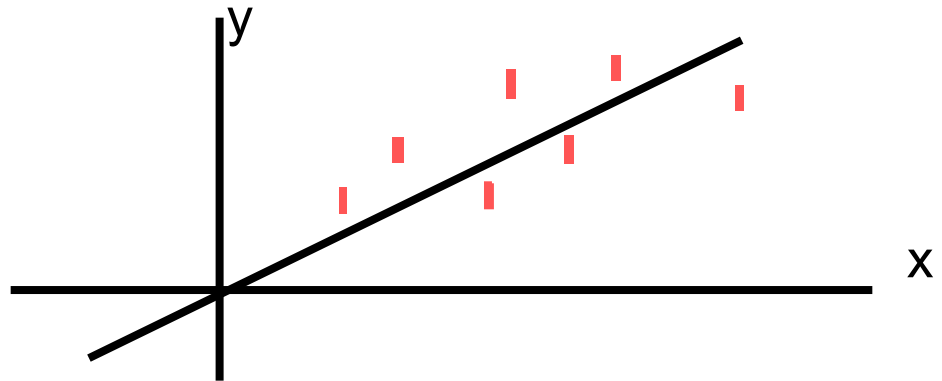
$$y_i = \theta \xi_i + \varepsilon_i'$$

M

$$\varepsilon_i, \varepsilon_i' : N(0, \sigma^2)$$

$$(x_n, y_n)$$

$$y = \theta x$$



Statistical Model

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

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

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

———— **semiparametric**

Least squares?

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

$$\frac{1}{n} \sum \frac{y_i}{x_i}, \quad \frac{\sum y_i}{\sum x_i}$$

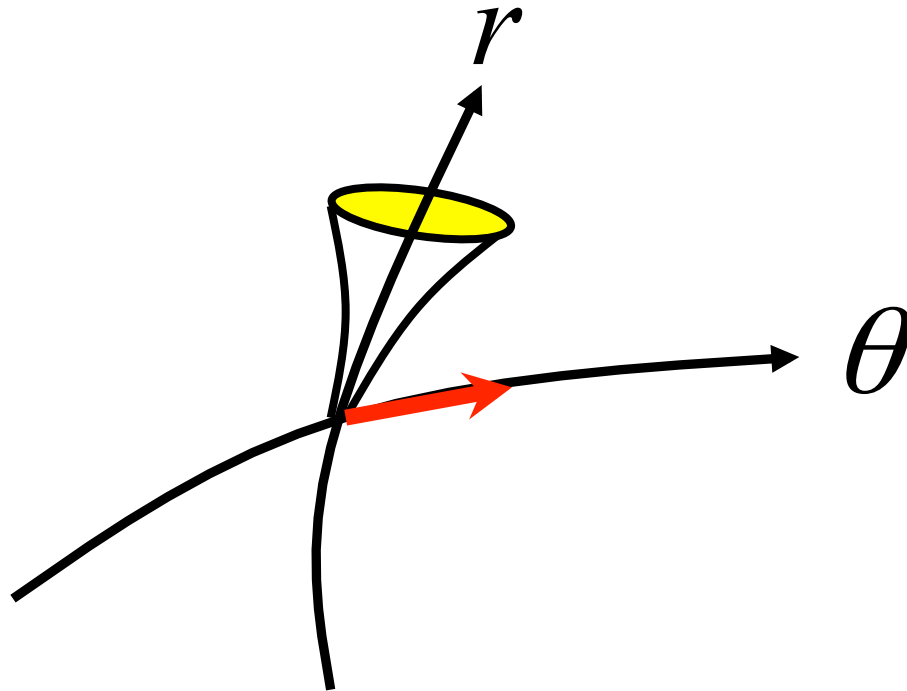
mle, TLS

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

Neyman-Scott

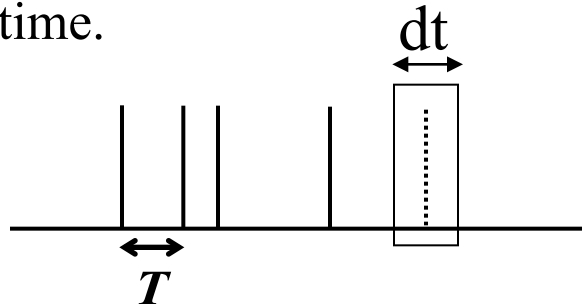
Fibre bundle

function space



Poisson process

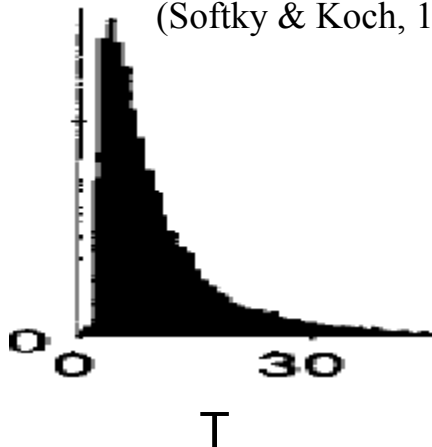
Poisson Process: Instantaneous firing rate is constant over time.



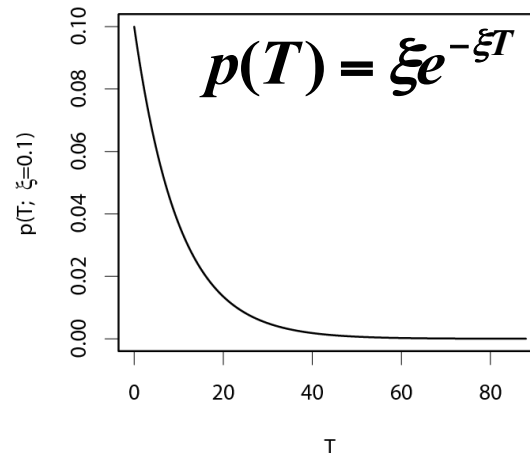
For every small time window dt , generate a spike with probability ξdt .

Cortical Neuron

(Softky & Koch, 1993)



Poisson Process



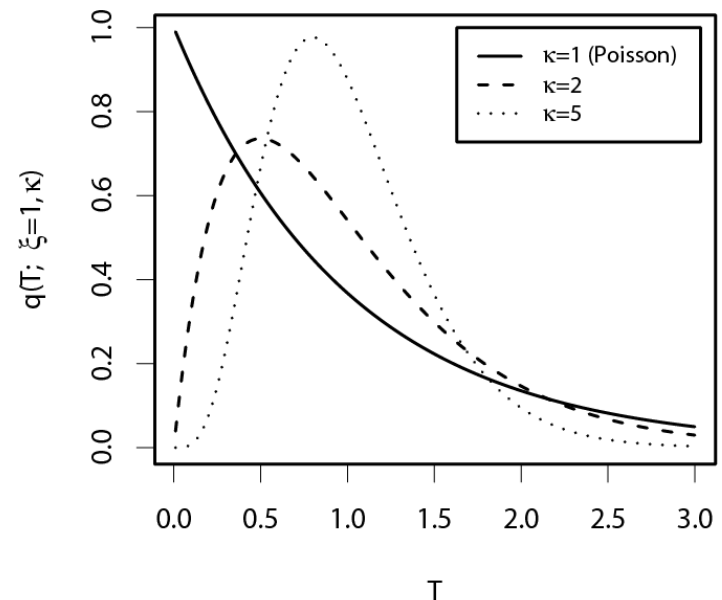
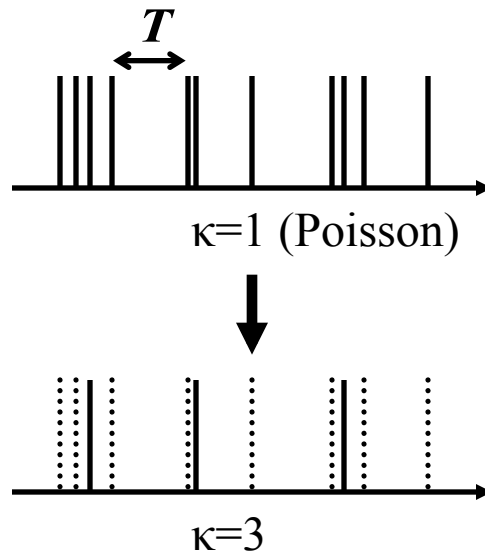
Poisson process cannot explain inter-spike interval distributions.

Gamma distribution

Gamma Distribution: Every κ -th spike of the Poisson process is left.

$$q(T; \xi, \kappa) = \frac{(\xi \kappa)^\kappa}{\Gamma(\kappa)} T^{\kappa-1} e^{-\xi \kappa T}.$$

Two parameters $\left\{ \begin{array}{l} \xi: \text{Firing rate} \\ \kappa: \text{Irregularity} \end{array} \right.$



Gamma distribution

$$f(T) = \frac{(r\kappa)^\kappa}{\Gamma(\kappa)} T^{\kappa-1} \exp\{-r\kappa T\}$$

$\kappa = 1$: Poisson

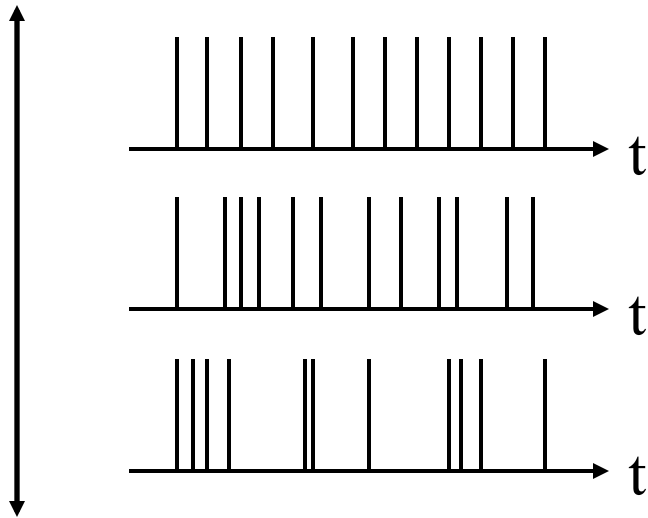
$\kappa \rightarrow \infty$: regular

Integrate-and fire

Markov model

Irregularity κ is unique to individual neurons.

Regular (large κ)



Irregular (small κ)



Irregularity varies among neurons.
(Baker & Lemon 2000; Shinomoto et.al., 2003)

➔ We assume that κ is independent of time.

Bayesian Information Geometry

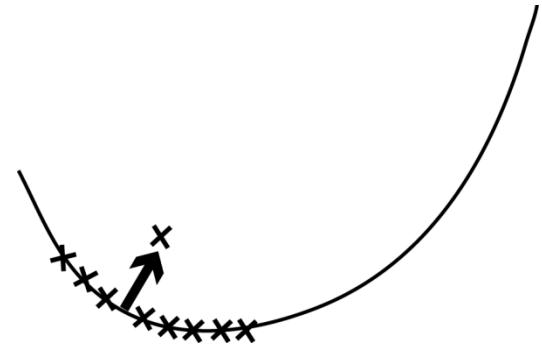
$$p(\mathbf{x}, \boldsymbol{\theta}) = \exp\{\boldsymbol{\theta} \cdot \mathbf{x} + k(\mathbf{x}) - \psi(\boldsymbol{\theta}) + \log \pi(\boldsymbol{\theta})\}$$

$$p(\mathbf{x}|\boldsymbol{\theta}), \quad p(\boldsymbol{\theta}|\mathbf{x})$$

Predictive Distribution

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

$$\begin{aligned} p(\mathbf{x}|D) &= \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta}|D)d\boldsymbol{\theta} \\ &= p(\mathbf{x}, \hat{\boldsymbol{\theta}}) + \frac{1}{N} H_{ijk} \hat{v}^k \end{aligned}$$



$$\text{Min}_p E \left[KL \left[p(\mathbf{x}, \boldsymbol{\theta}_0) : p(\mathbf{x}|D) \right] \right]$$

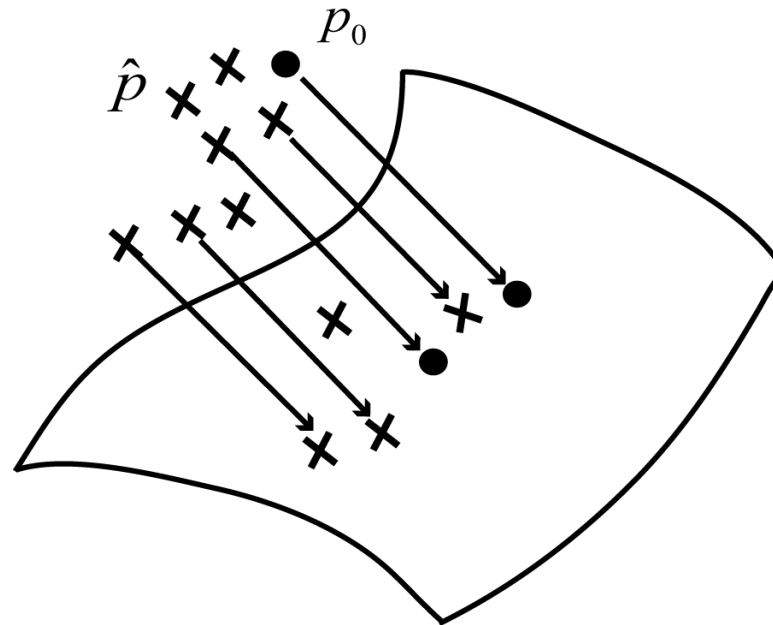
$$\text{Min} ED_\alpha \left[p(\mathbf{x}, \boldsymbol{\theta}_0) : p(\mathbf{x}|D) \right]$$

α -predictive distribution

Bootstrap Method

$$D = \{x_1, \dots, x_N\} \Rightarrow \hat{p}(x)$$

resampling :



EM algorithm

hidden variables

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

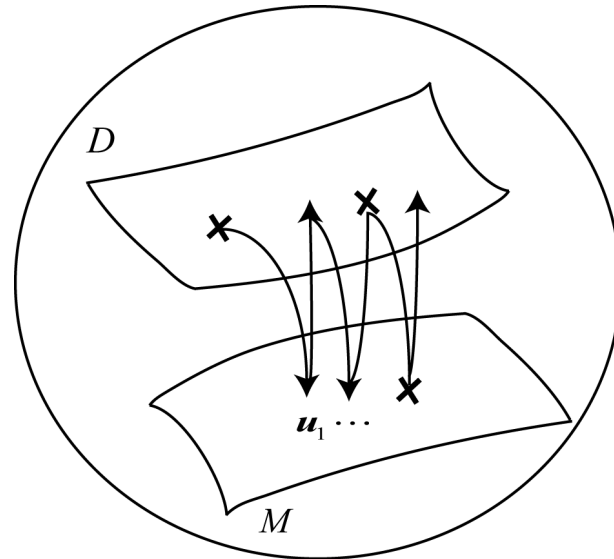
$$D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$

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

$$D_M = \{p(\mathbf{x}, \mathbf{y}) \mid p(\mathbf{x}) = p_D(\mathbf{x})\}$$

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

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



Tsallis q-entropy

conformal information geometry

$$H_T = E \left[\ln_q \frac{1}{p(x)} \right] = \frac{1}{1-q} \left\{ \int p(x)^q dx - 1 \right\}$$

Computer vision:

$$s(x, y) \geq 0$$

Divergence

Clustering

Center

Retrieval

Total Bregman Divergence and its Applications to Shape Retrieval

• **Baba C. Vemuri, Meizhu Liu, Shun-ichi Amari,
Frank Nielsen**

IEEE Conference on Computer Vision and Pattern Recognition (CVPR),
2010

3. Conformal change of divergence

$$D^{\sigma}(p:q) = \sigma(p)\sigma(q)D[p:q]$$

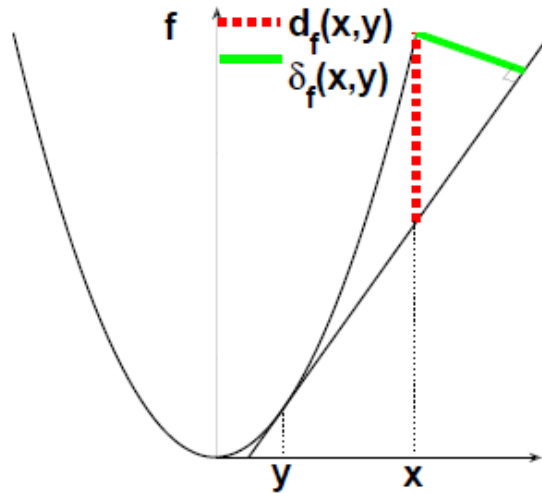
$$g_{ij}^{\sigma} = \sigma(p)\sigma(p)g_{ij}$$

$$\Gamma_{ijk}^{\sigma} = \sigma_1\sigma_2\Gamma_{ijk} + \sigma_2\partial_i\sigma_1g_{ijk} + \sigma_2\partial_j\sigma_1g_{ik} \\ - \sigma_1\partial_k\sigma_2g_{ij}$$

Conformally flat \rightarrow canonical divergence

Total Bregman Divergence

$$TD[x : y] = \frac{D[x : y]}{\sqrt{1 + \|\nabla f\|^2}}$$



- rotational invariance
- conformal geometry

Figure: $d_f(x, y)$ (dotted red line) is BD, $\delta_f(x, y)$ (bold green line) is TBD, and the two arrows indicate the coordinate system. Note that $d_f(x, y)$ changes with rotation unlike $\delta_f(x, y)$ which is invariant to rotation.

t-center \mathbf{x}^*

$$\nabla f(\mathbf{x}^*) = \frac{\sum w_i \nabla f(\mathbf{x}_i)}{\sum w_i}$$

$$w_i = \frac{1}{\sqrt{1 + \|\nabla f(\mathbf{x}_i)\|^2}}$$

t-center is robust

$$E^* = \{\mathbf{x}_1, \dots, \mathbf{x}_n; \mathbf{y}\}$$

$$\hat{\theta}_0^* = \mathbf{x}^* + \varepsilon \mathbf{z}(\mathbf{x}^*; \mathbf{y}), \quad \varepsilon = \frac{1}{n}$$

influence function $\mathbf{z}(\mathbf{x}^*; \mathbf{y})$

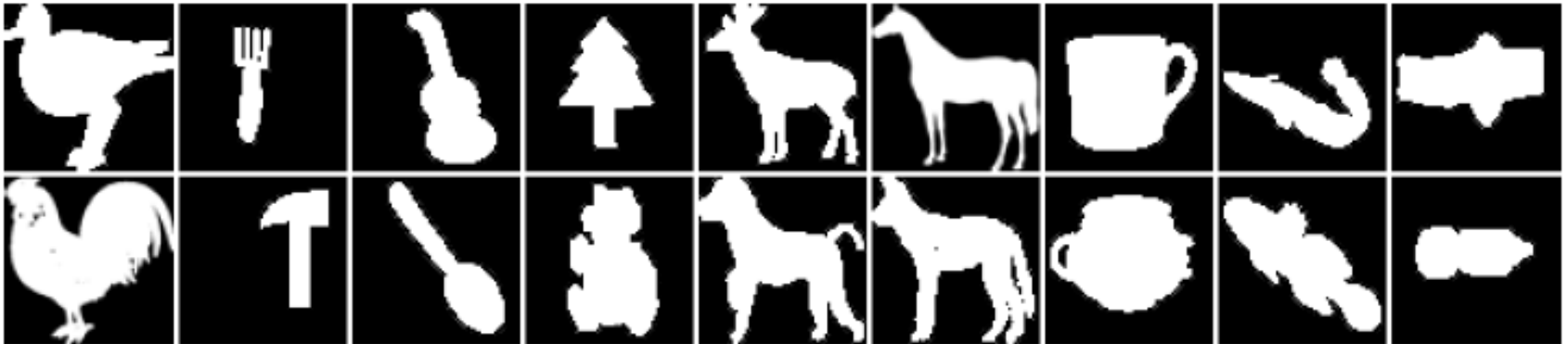
$|\mathbf{z}| < c$ as $|\mathbf{y}| \rightarrow \infty$: robust

TBD application-shape retrieval

- **Using MPEG7 database;**
- **70 classes, with 20 shapes each class**
(Meizhu Liu)

MPEG7 database

- **Great intraclass variability, and small interclass dissimilarity.**



Shape representation

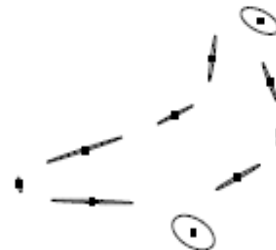
A shape is represented using a mixture of Gaussians from the aligned boundary points.



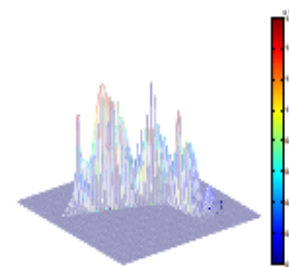
(a)



(b)



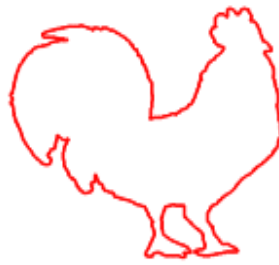
(c)



(d)



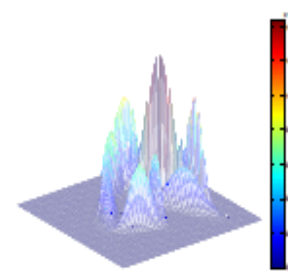
(e)



(f)



(g)



(h)

Experimental results

Technique	Recognition rate (%)
Mixture of Gaussians + tSL	89.1
Mixture of Gaussians + χ^2	63.3
Mixture of Gaussians + SL	56.7
Shape-tree[6]	87.7
IDSC + DP + EMD[14]	86.56
Hierarchical Procrustes [15]	86.35
IDSC + DP [13]	85.4
Shape L'Âne Rouge[18]	85.25
Generative Models [21]	80.03
Curve Edit [19]	78.14
SC + TPS [3] [3]	76.51
Visual Parts [10]	76.45
CSS [16]	75.44

Table 3. Recognition rates comparison.

Neural Networks

Multilayer Perceptron

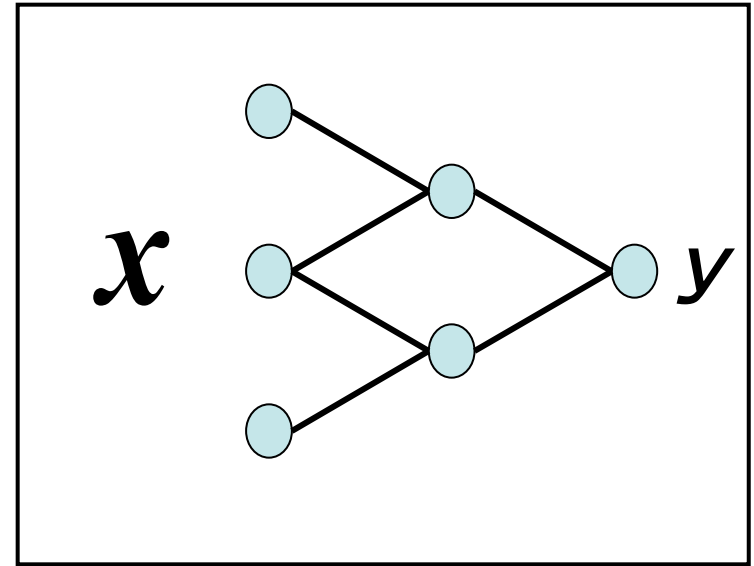
Higher-order correlations

Synchronous firing

Multilayer Perceptrons

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

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



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

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

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

Multilayer Perceptron

neuromanifold

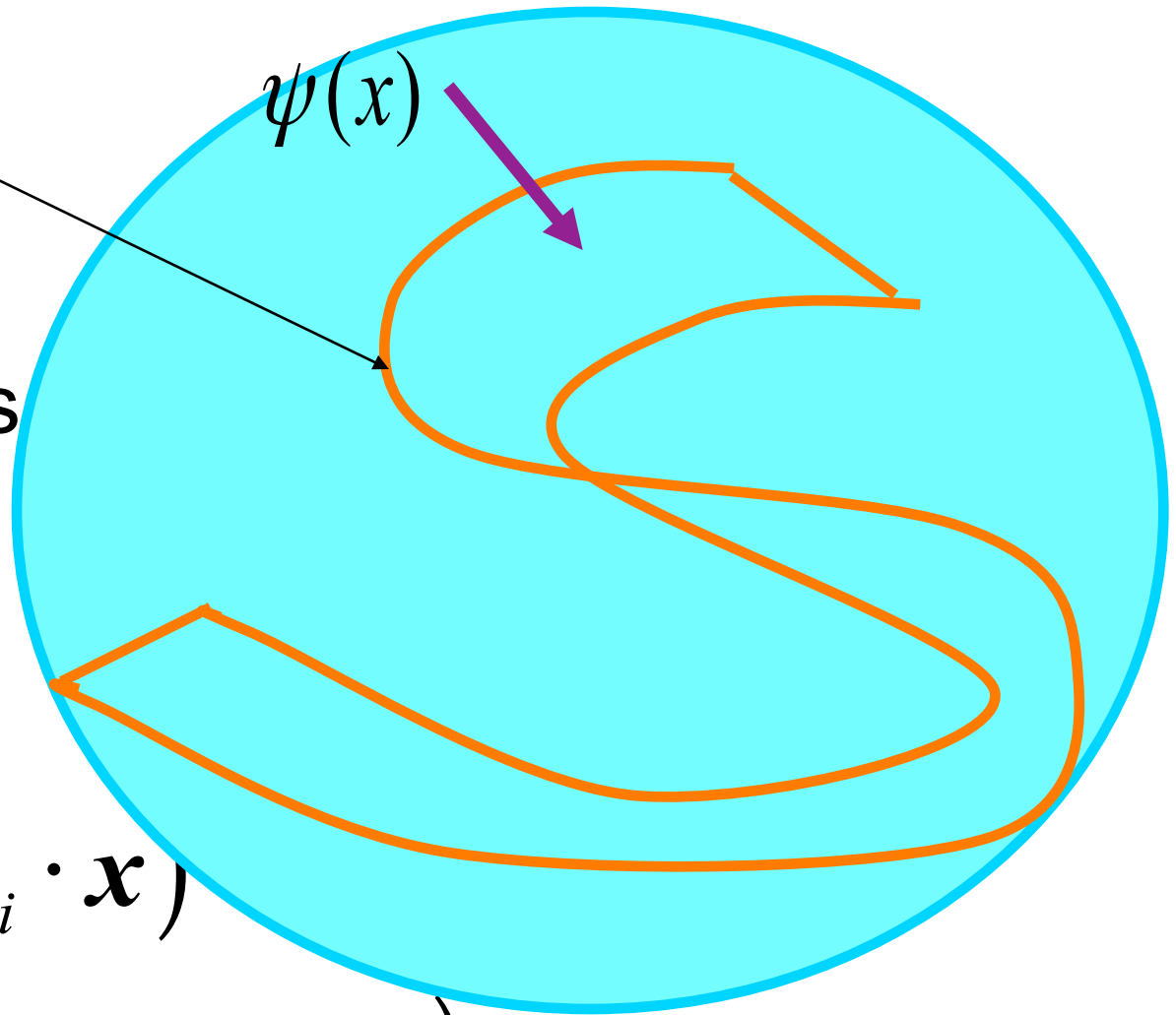
$\psi(x)$

space of functions

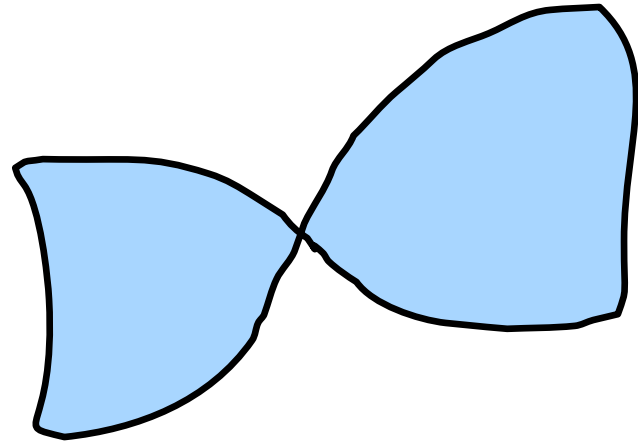
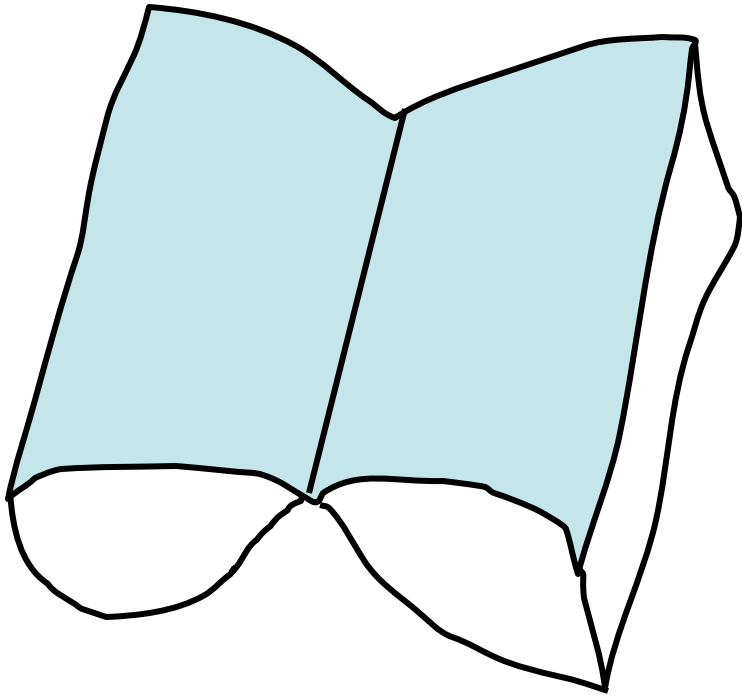
$$y = f(\mathbf{x}, \boldsymbol{\theta})$$

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

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



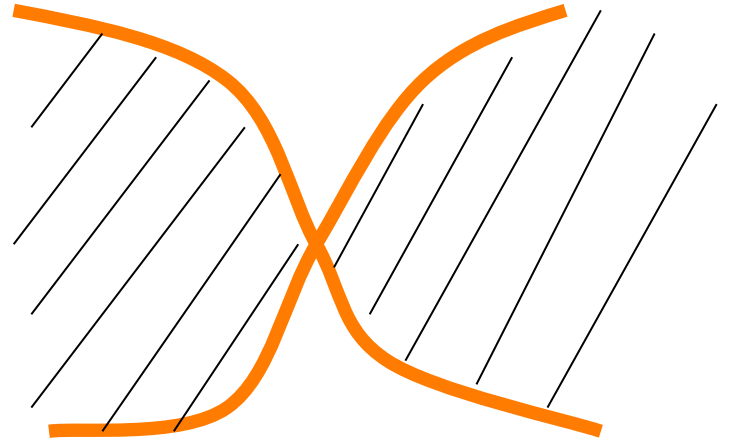
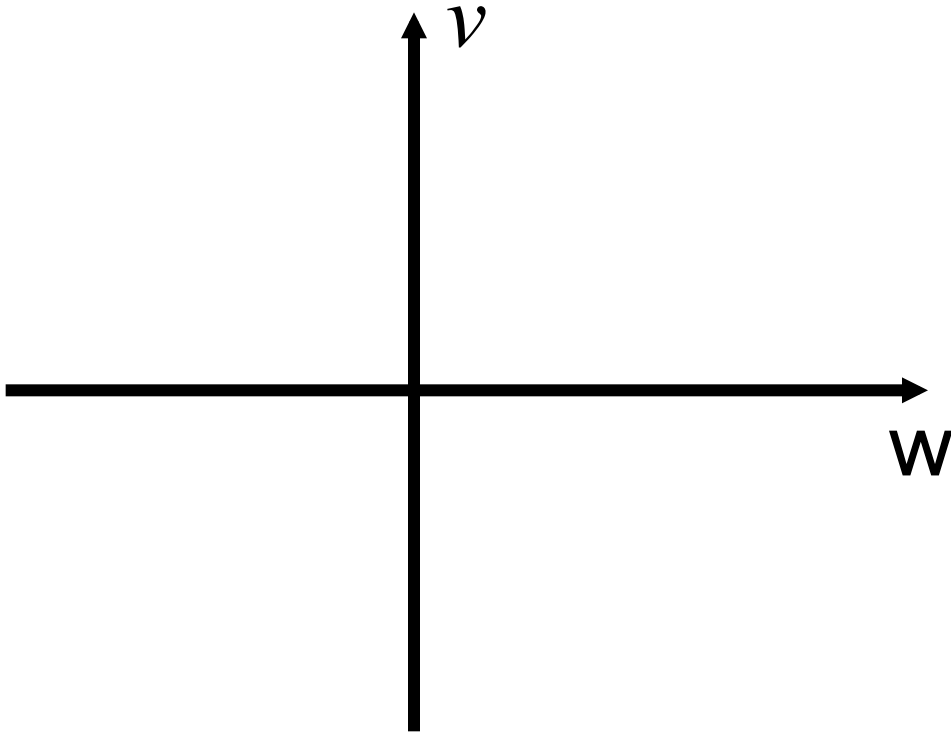
singularities



Geometry of singular model

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

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



Backpropagation ---gradient learning

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

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

natural gradient (Riemannian)

$$\nabla_{\boldsymbol{\theta}} E = G^{-1} \nabla E \text{ ---steepest descent}$$

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

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

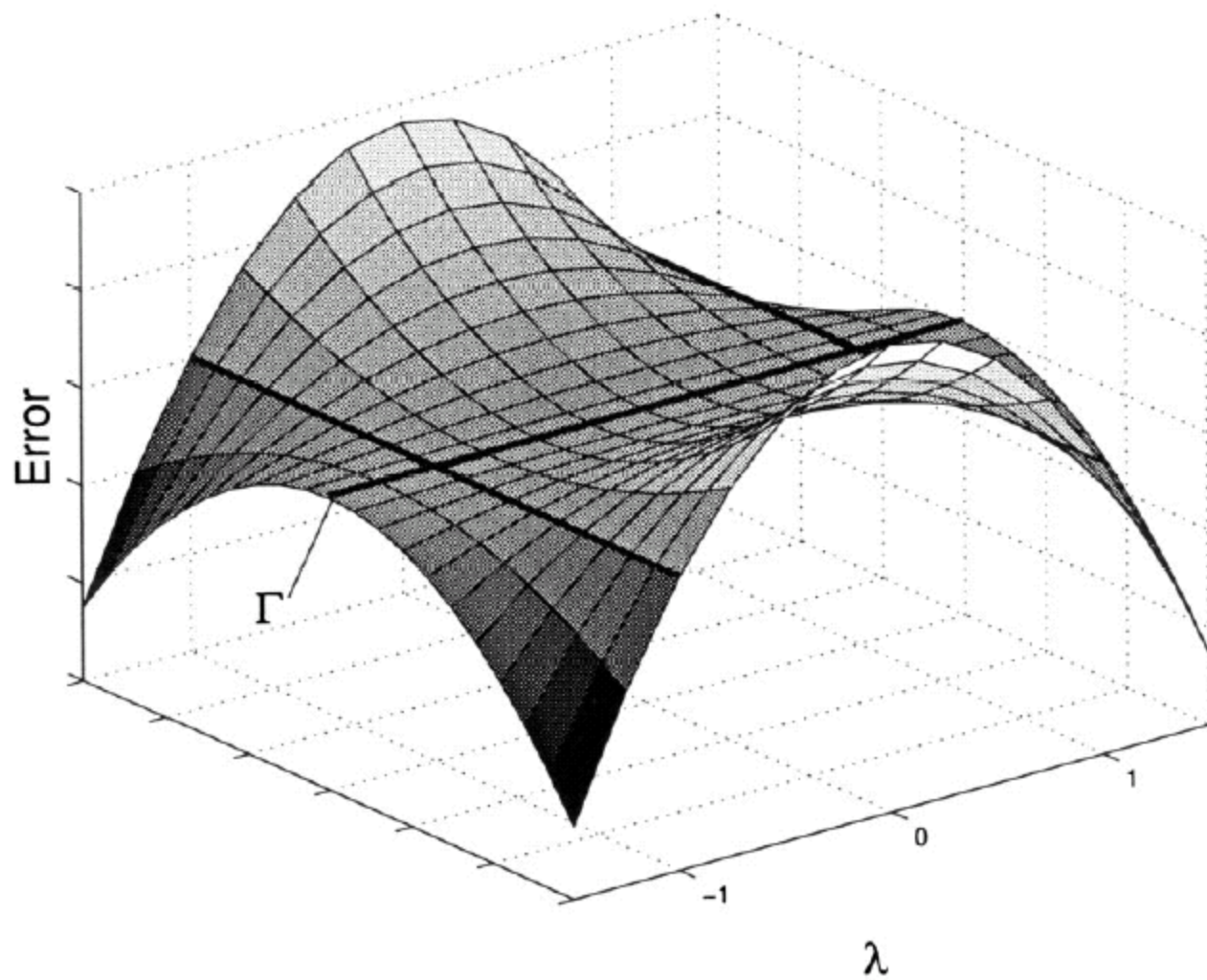
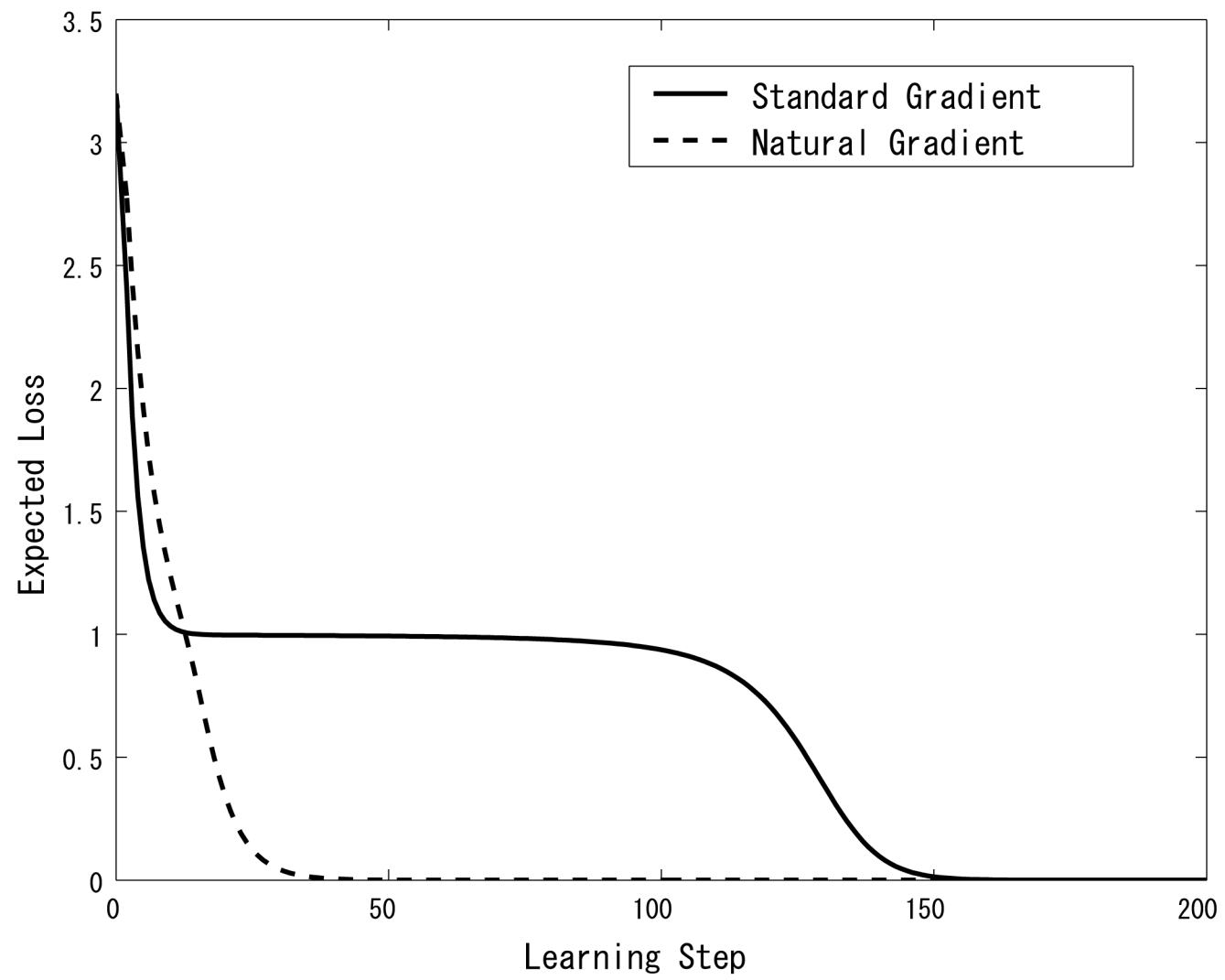
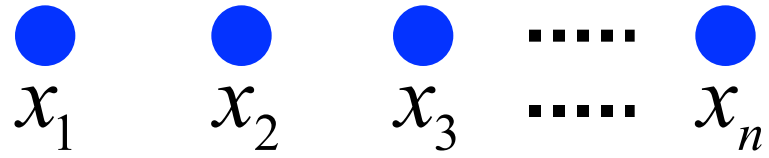


Fig. 5. Critical set with local minima and plateaus.



Neural Firing



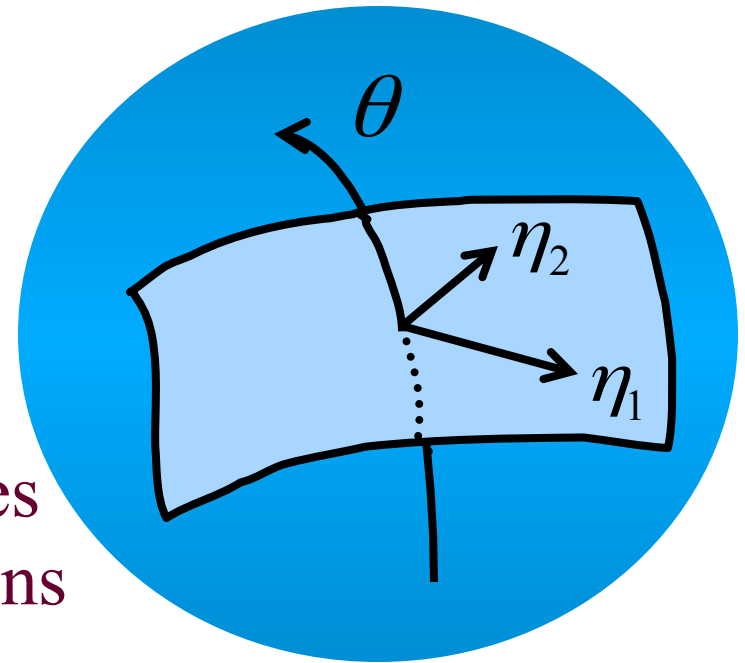
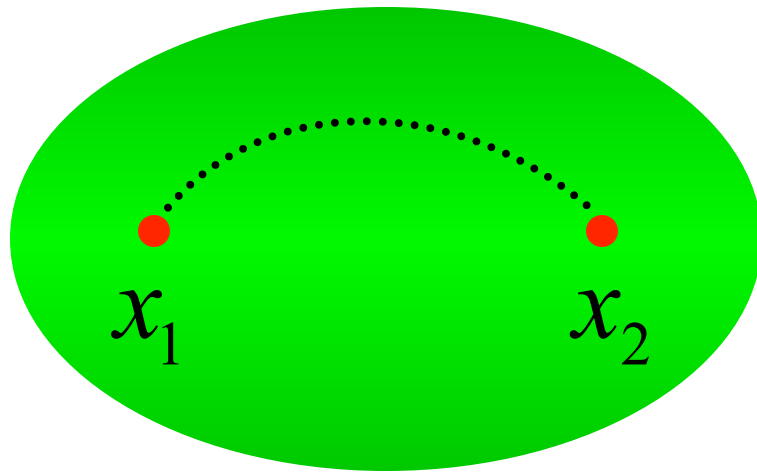
η_i : firing rate

θ_{ij} : correlations

$p(x_1, x_2)$

orthogonal decomposition
higher-order correlations

Correlations of Neural Firing



$$\{p(x_1, x_2)\}$$

firing rates
correlations

$$\{p_{00}, p_{10}, p_{01}, p_{11}\}$$

$$\eta_1 = p_{.1}$$

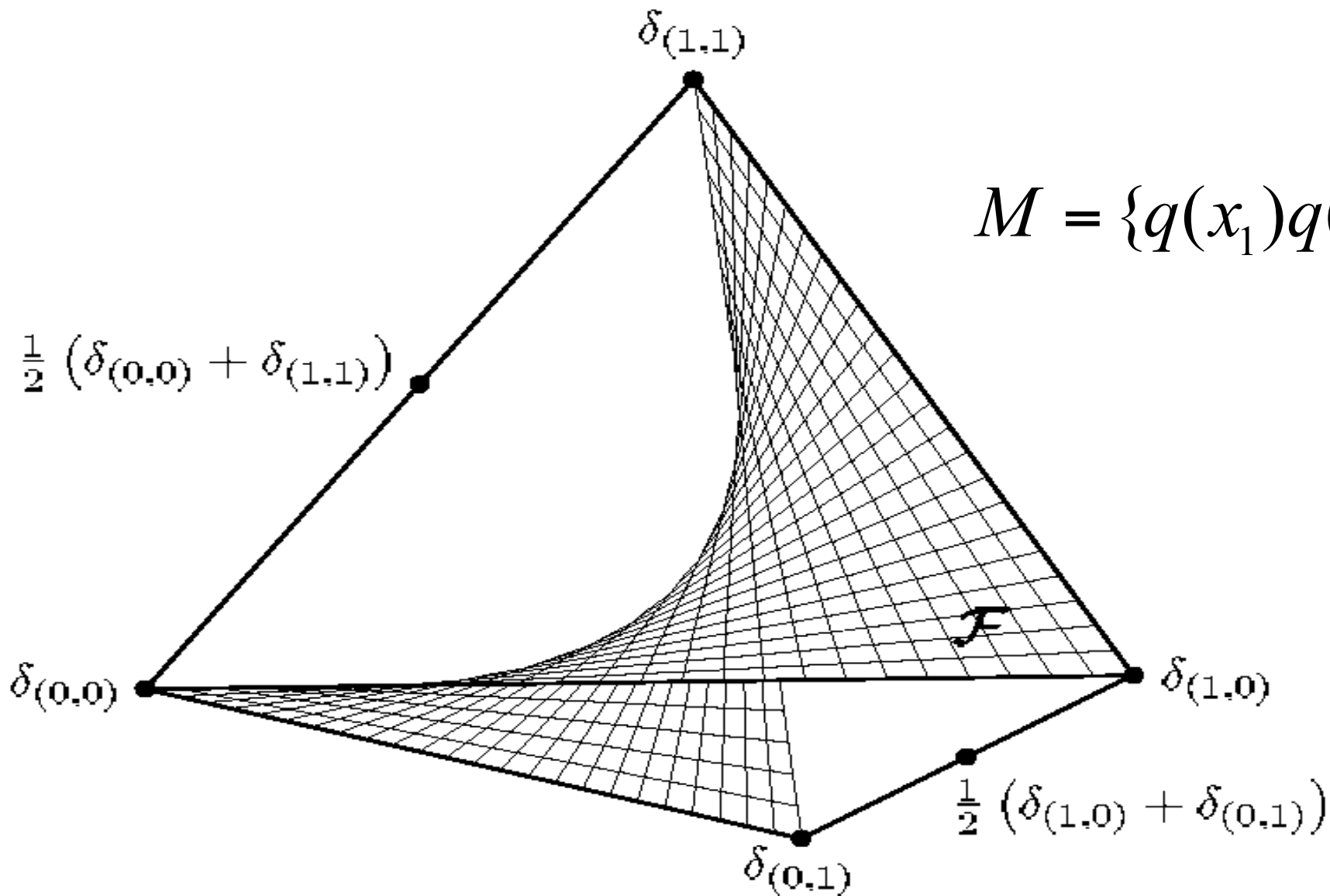
$$\theta = \log \frac{p_{11} p_{00}}{p_{10} p_{01}}$$

$$\{(\eta_1, \eta_2), \theta\}$$

$$\eta_2 = p_{.1}$$

orthogonal coordinates

Independent Distributions



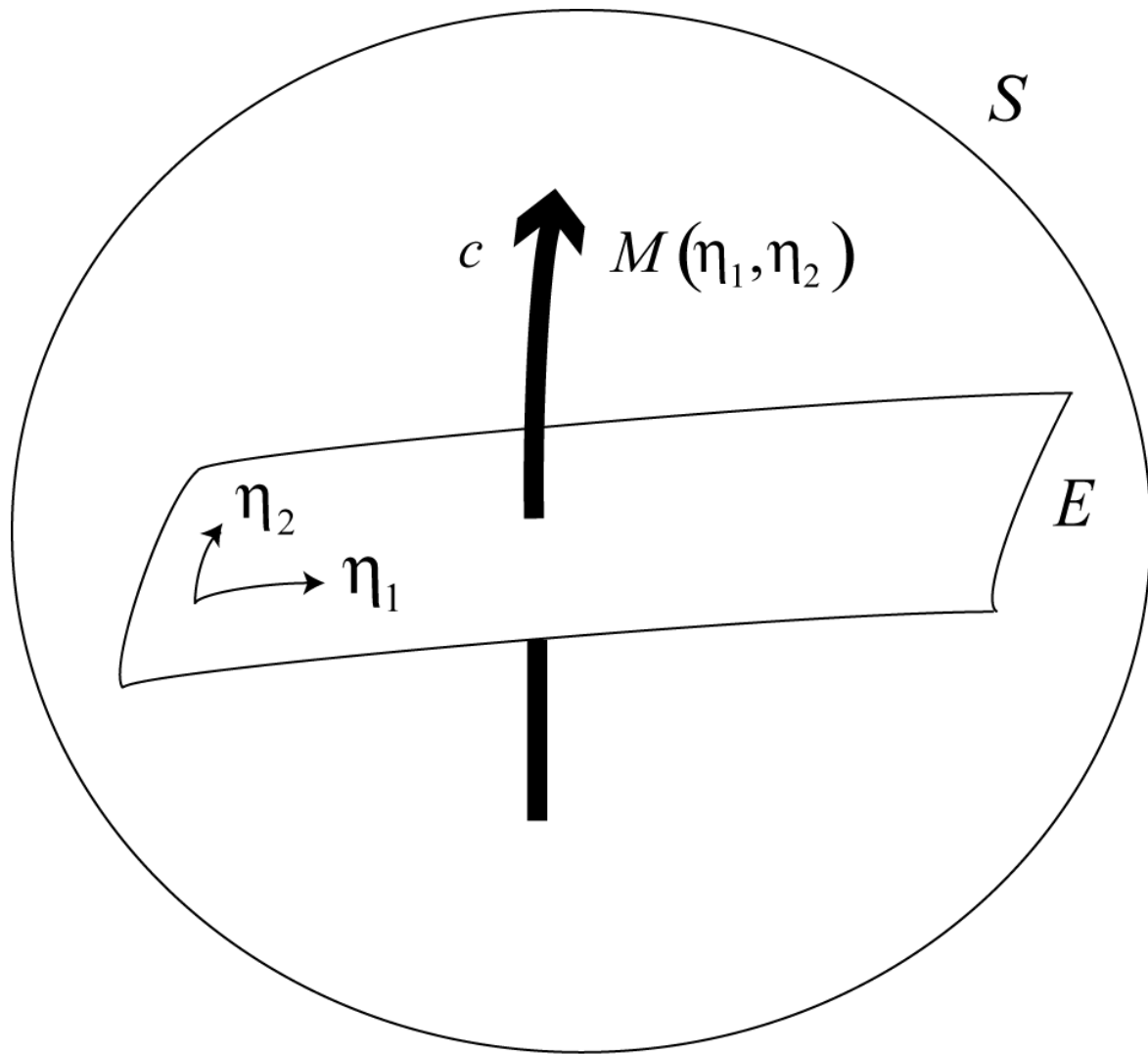
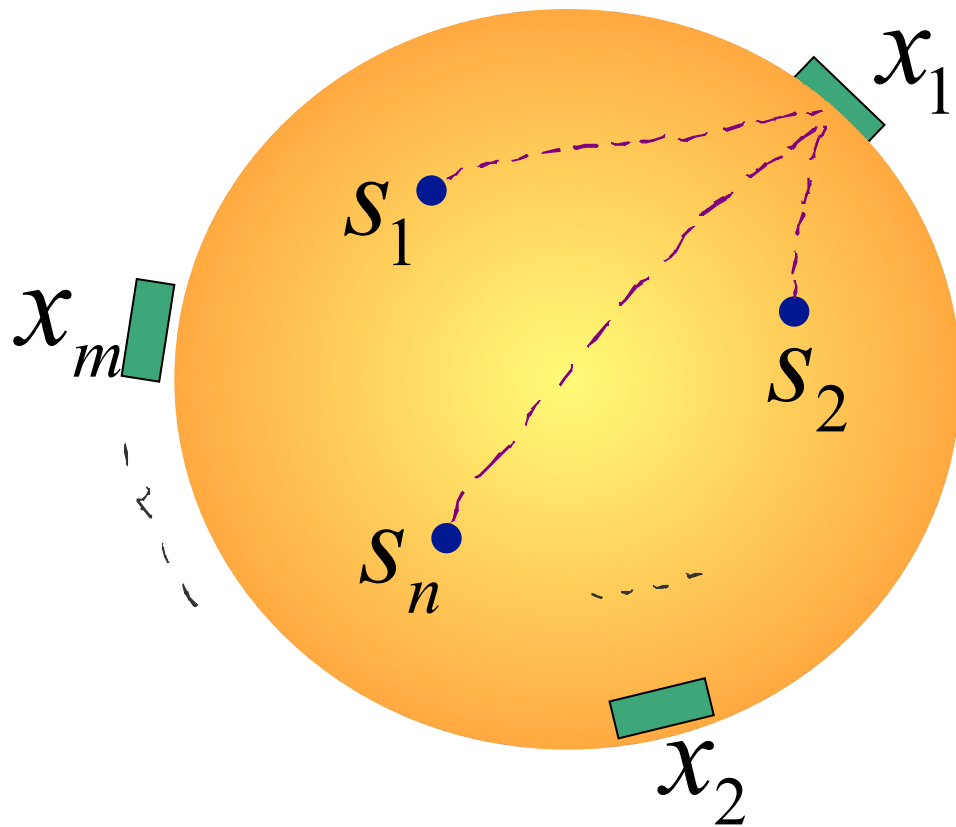


Fig. 1a

mixture and unmixture of independent signals



$$x_i = \sum_{j=1}^n A_{ij} s_j$$

$$\mathbf{x} = \mathbf{A}\mathbf{s}$$

Signal Processing

ICA : Independent Component Analysis

$$\mathbf{x}_t = A\mathbf{s}_t \quad \mathbf{x}_t \rightarrow \mathbf{s}_t$$

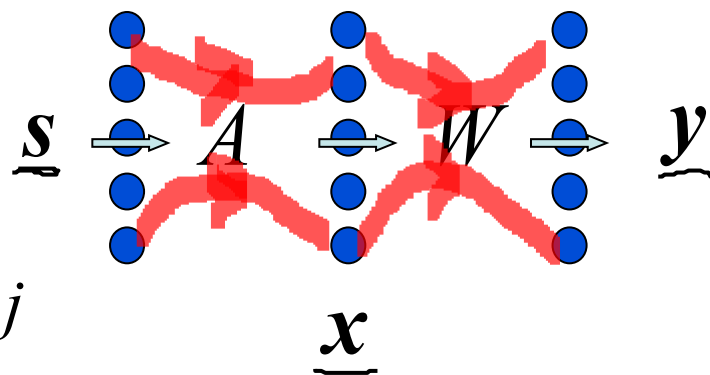
sparse component analysis

positive matrix factorization

Independent Component Analysis

$$\mathbf{x} = \mathbf{A}\mathbf{s}$$

$$x_i = \sum_j A_{ij} s_j$$



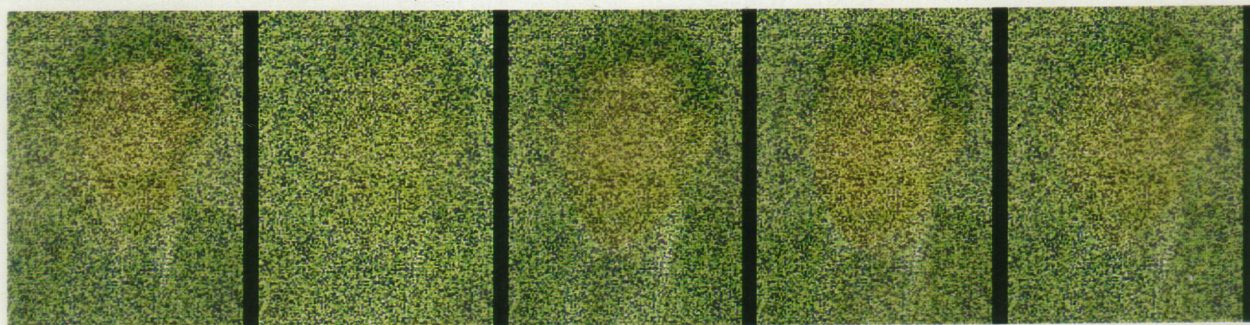
$$\mathbf{y} = \mathbf{W}\mathbf{x} \quad \mathbf{W} = \mathbf{A}^{-1}$$

observations: $\mathbf{x}(1), \mathbf{x}(2), \dots, \mathbf{x}(t)$

recover: $\mathbf{s}(1), \mathbf{s}(2), \dots, \mathbf{s}(t)$

Example of color image separation :

Five original images (but unknown to the neural net)

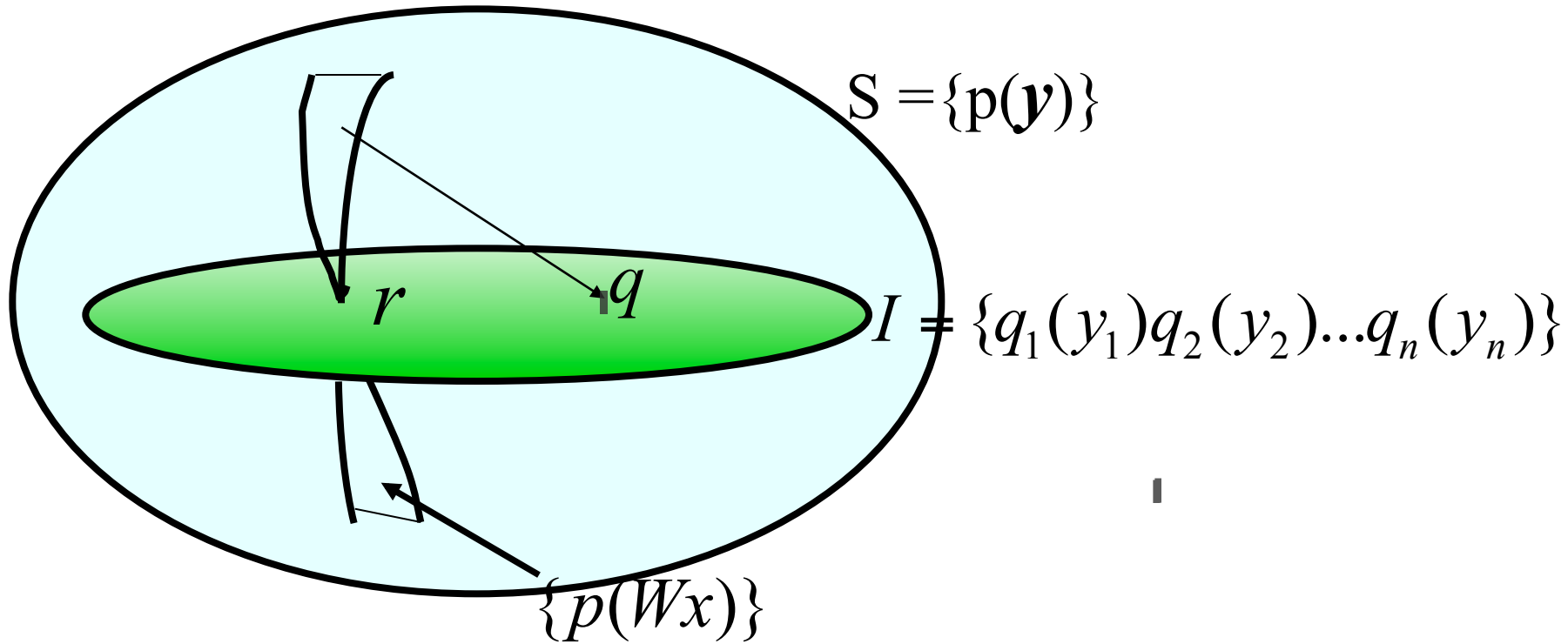


Five mixed images for separation



Final (stable states) of five separated images

Information Geometry of ICA

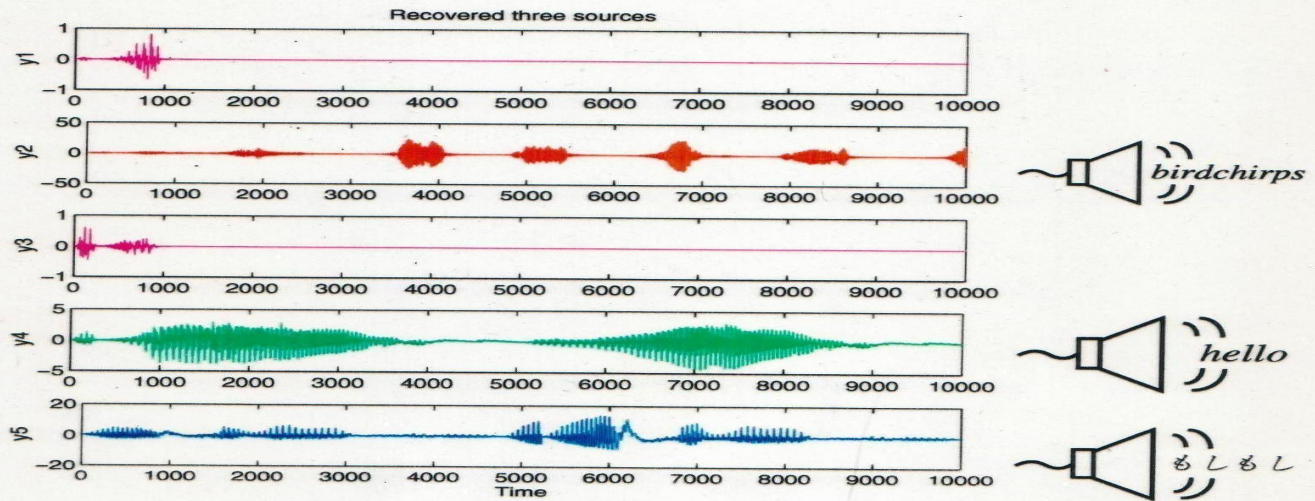
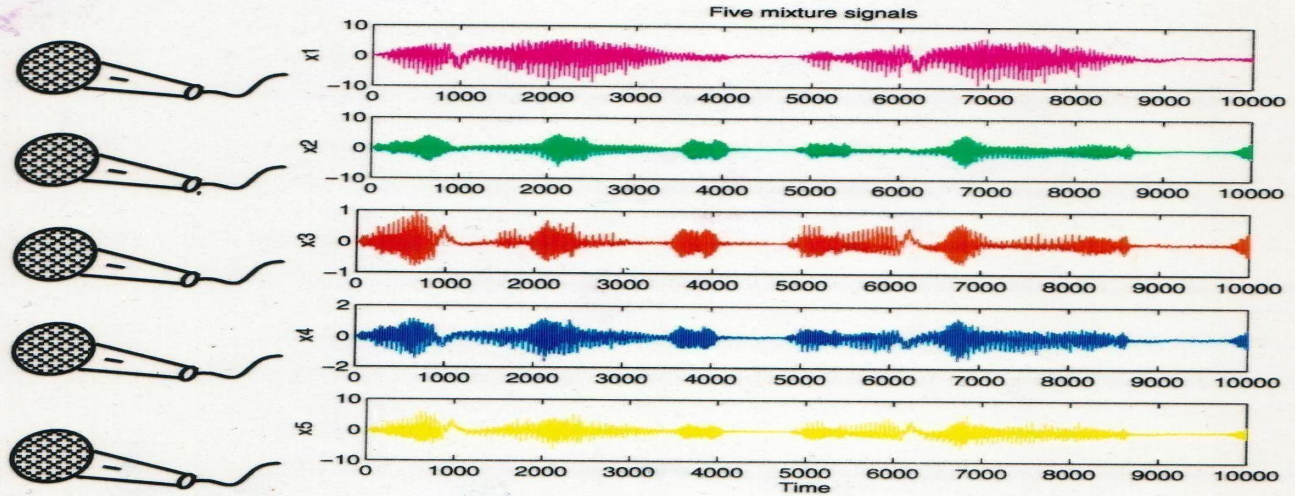


natural gradient
estimating function
stability, efficiency

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

$r(\mathbf{y})$

Cocktail party experiment



- 5 microphones (sensors) and only 3 speakers

Natural Gradient

$$\Delta W = -\eta \frac{\partial l(y, W)}{\partial W} W^T W$$

Basis Given: overcomplete case

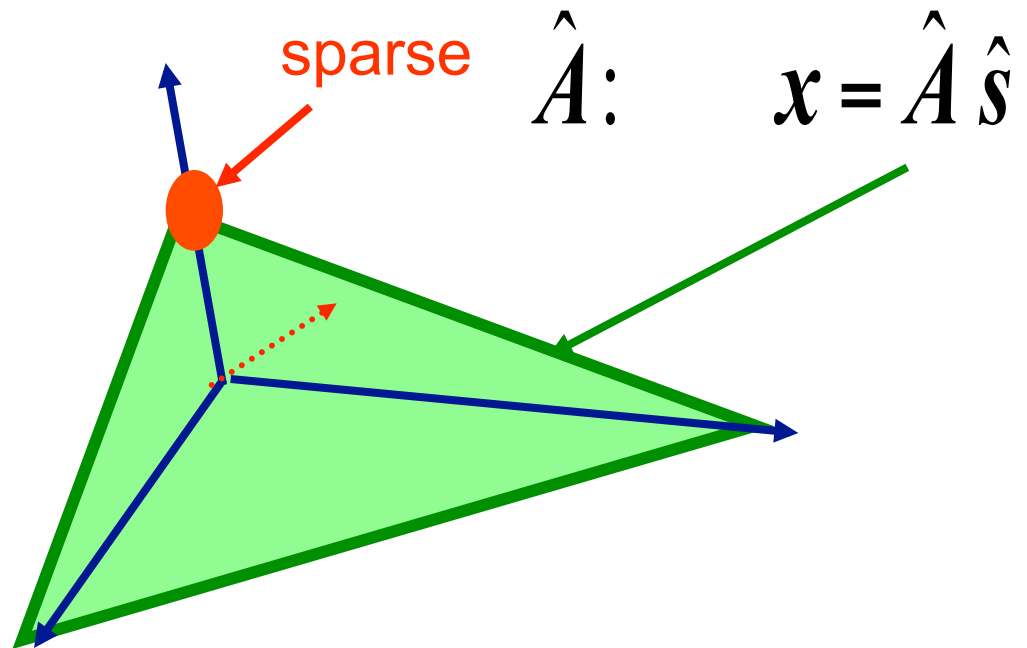
Sparse Solution

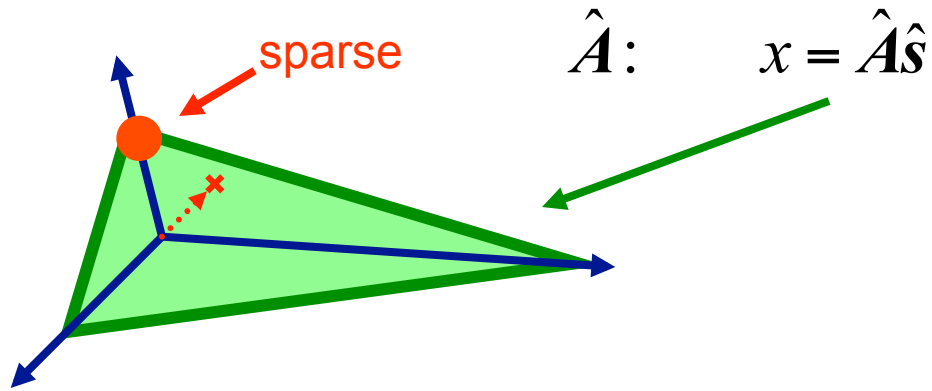
$$\mathbf{x} = A\mathbf{s} = \sum s_i \mathbf{a}_i$$

many solutions

many $s_i \rightarrow 0$

$$\mathbf{x}_t = \hat{A}\mathbf{s}_t$$



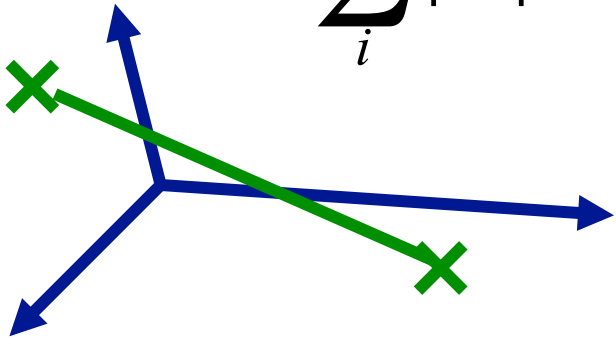


generalized inverse

L_2 -norm: $\min \sum |\hat{s}_i|^2$

sparse solution

L_1 -norm: $\min \sum |\hat{s}_i|$



Overcomplete Basis and Sparse Solution

$$\mathbf{x} = \sum s_i \mathbf{a}_i = A\mathbf{s}$$

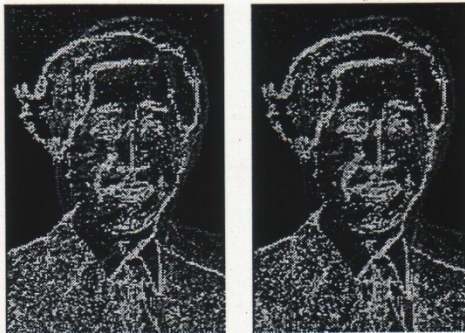
$$\min \|\mathbf{s}\|_1 = \sum |s_i|$$

$$\min \|\mathbf{A}\mathbf{s} - \mathbf{x}\|_p + \alpha \|\mathbf{s}\|_p,$$

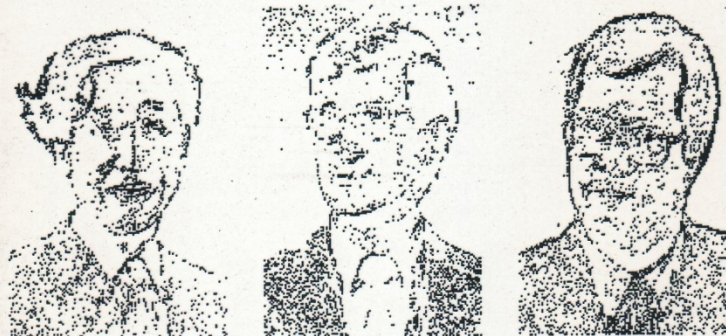
non-linear denoising



(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).

Linear Systems



ARMA

$$x_{t+1} = \frac{1 + b_1 z^{-1} + \dots + b_q z^{-q}}{1 + a_1 z^{-1} + \dots + a_p z^{-p}} u_t$$

$$\theta = (a_1, \dots, a_p : b_1, \dots, b_q)$$

$$x_{t+1} = f(\theta, z^{-1}, u_t)$$

AR---e-flat

MA---m-flat

Information Geometry of Belief Propagation

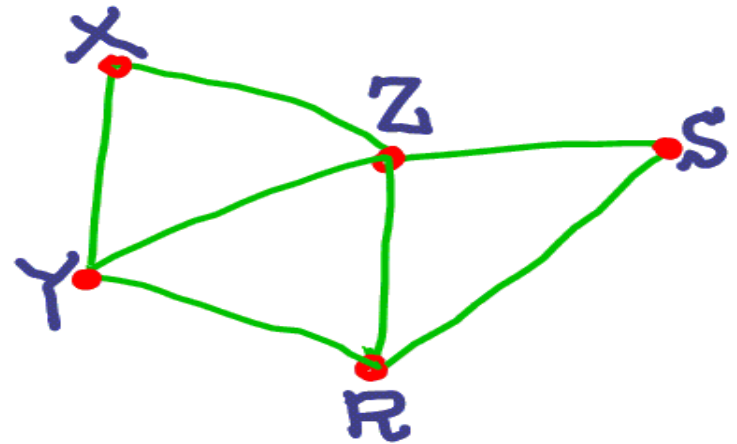
- Shun-ichi Amari (RIKEN BSI)
- Shiro Ikeda (Inst. Statist. Math.)
- Toshiyuki Tanaka (Tokyo Metropolitan U.)

Stochastic Reasoning

$$p(x, y, z, r, s)$$

$$p(x, y, z \mid r, s)$$

$$x, y, z, \dots = 1, -1$$



Mean Value

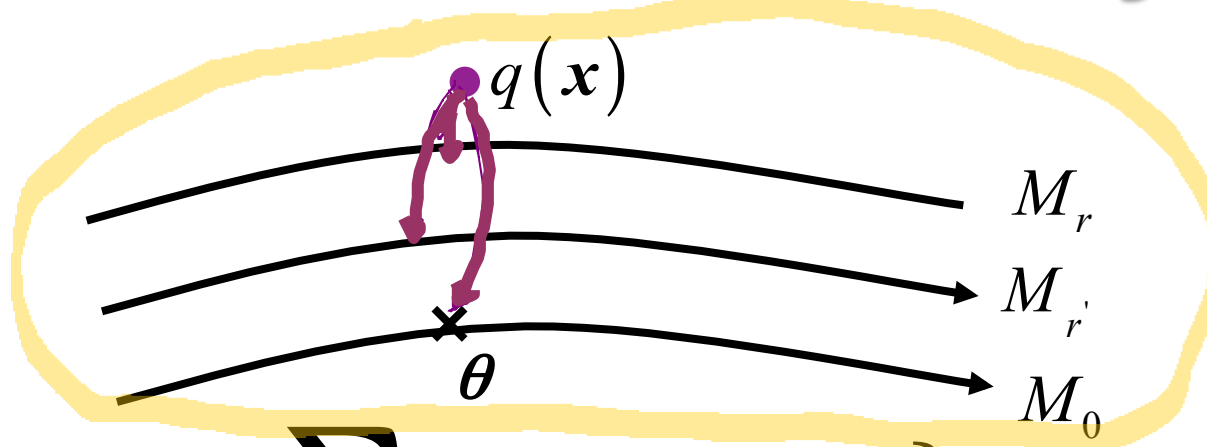
Marginlization

$$\Pi_0 q(\mathbf{x}) = q_1(x_1)q_2(x_2)\dots q_n(x_n) = q_0(\mathbf{x})$$

$$q_i(x_i) = \int q(x_1, \dots, x_n) dx_1 \dots dx_i \dots dx_n$$

$$\boldsymbol{\eta} = \mathbf{E}_q[\mathbf{x}] = \mathbf{E}_{q_0}[\mathbf{x}]$$

Information Geometry



$$q(x) = \exp \left\{ \sum c_r(x) - \phi \right\}$$

$$M_0 = \left\{ p_0(\mathbf{x}, \boldsymbol{\theta}) \right\} = \exp \left\{ \boldsymbol{\theta} \cdot \mathbf{x} - \psi_0 \right\}$$

$$M_r = \left\{ p_r(\mathbf{x}, \boldsymbol{\zeta}_r) = \exp \left\{ c_r(\mathbf{x}) + \boldsymbol{\zeta}_r \cdot \mathbf{x} - \psi_r \right\} \right\}_{r=1, L, L}$$

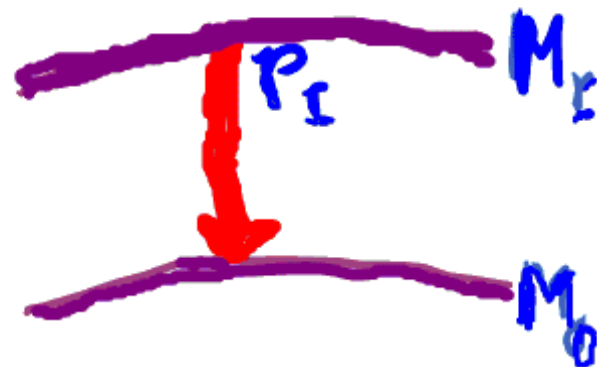
Belief Propagation

$$\prod_0 p_r(x, \xi_r) \quad p(x, \xi_r) = \exp\{c_r(x) + \xi_r \cdot x - \psi_r\}$$

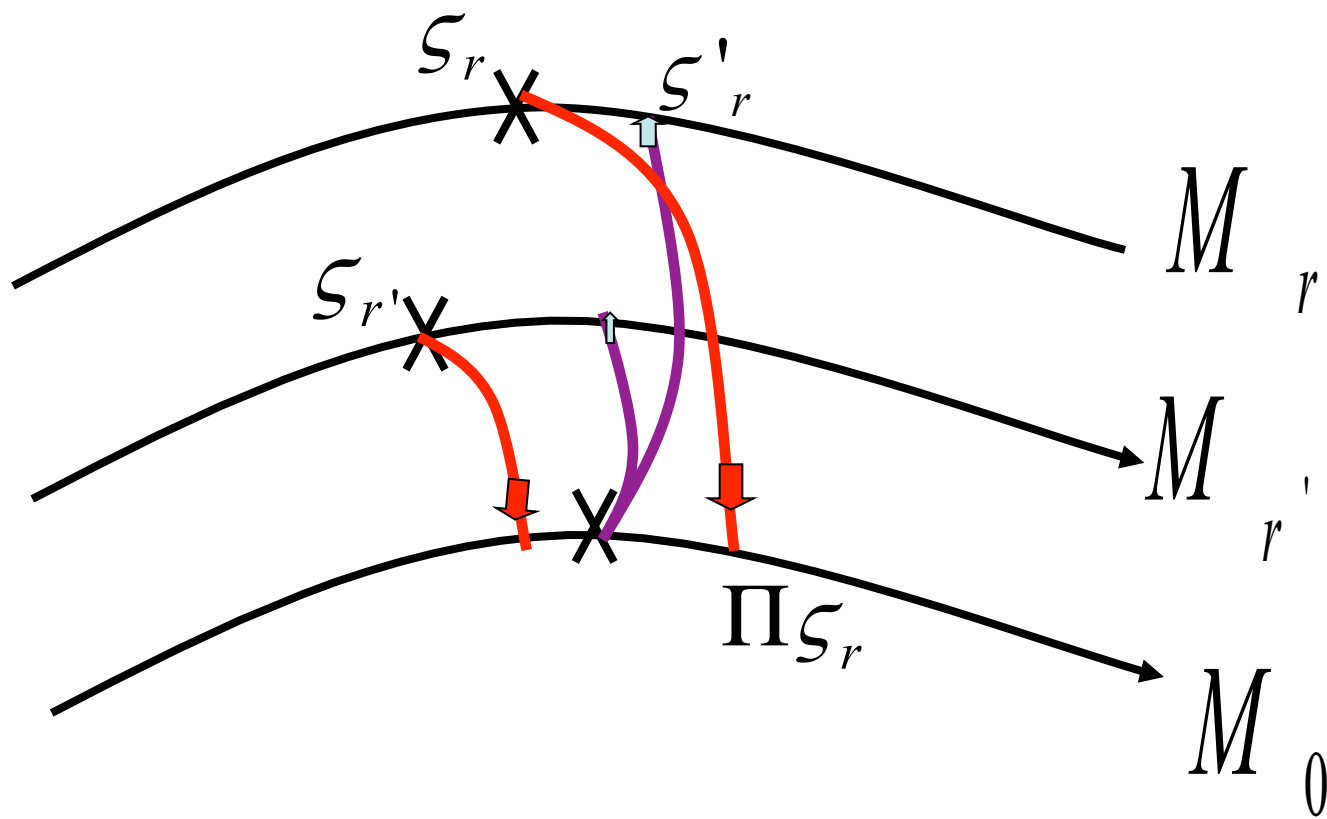
$$\xi_r^{t+1} = \prod_0 p_r(x, \xi_r^t) - \xi_r^t : \text{belief for } c_r(x)$$

$$\xi_r^{t+1} = \sum_{r' \neq r} \prod_0 \left\{ p(\xi_{r'}^t) - \xi_{r'}^t \right\} = \sum \xi_{r'}^{t+1}$$

$$\theta^{t+1} = \frac{1}{L-1} \sum_r \xi_r^{t+1} = \sum \xi_r^{t+1}$$



Belief Prop Algorithm



Belief Propagation

e-condition OK

$$F(\theta; \xi_1, \xi_2, \dots, \xi_L), \quad \theta = \theta(\xi_1, \dots, \xi_L)$$

$$(\xi_1, \dots, \xi_L) \rightarrow (\xi'_1, \dots, \xi'_L)$$

CCCP m-condition OK

$$\xi_1(\theta), \xi_2(\theta), \xi_L(\theta)$$

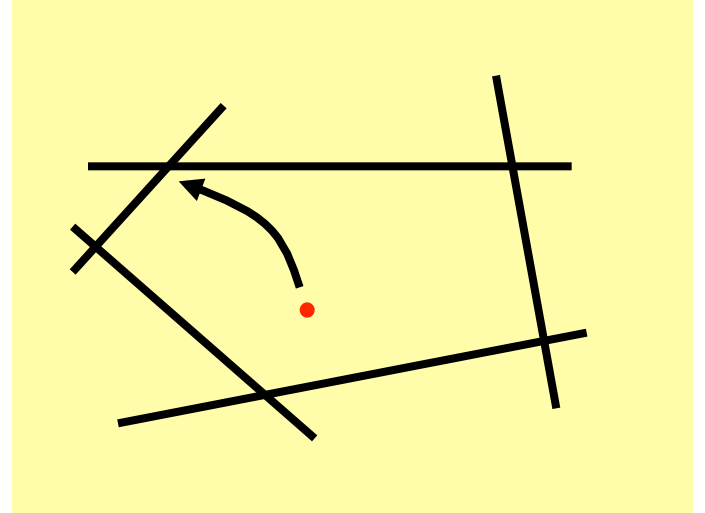
$$\theta \rightarrow \theta'$$

Linear Programming

$$\sum A_{ij}x_j \geq b_i$$

$$\max \sum c_i x_i$$

$$\psi(\mathbf{x}) = \sum_i \log\left(\sum A_{ij}x_j - b_i\right)$$



Convex Programming

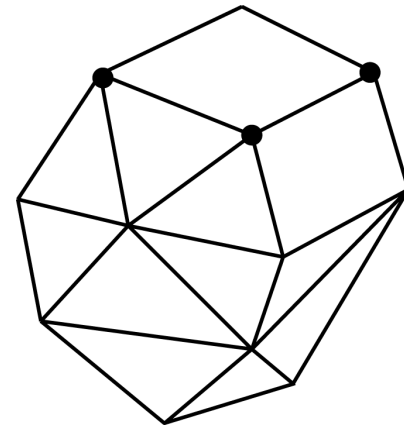
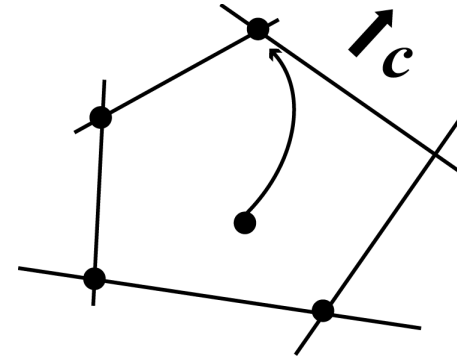
— Inner Method

$$LP: Ax \geq b, \quad c \cdot x \geq 0$$

$$\min \quad c \cdot x$$

$$\psi(x) = \sum \log\left(\sum A_{ij}x_j - b_i\right) + \sum \log x_i$$

$$\eta = \partial_i \psi(x)$$



Simplex method ; inner method

Convex Cone Programming

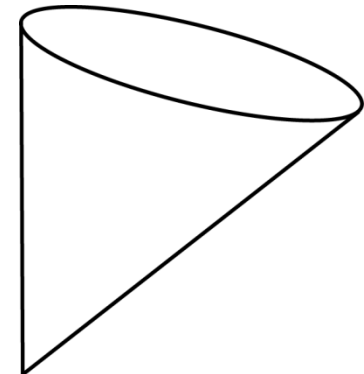
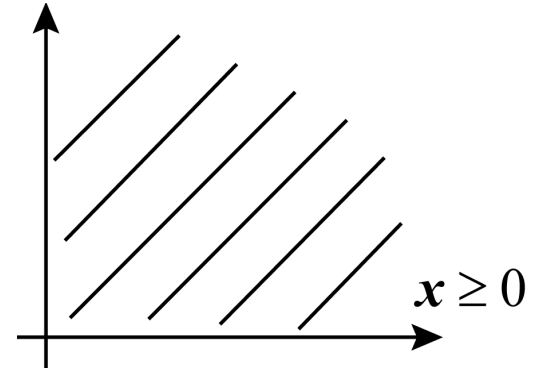
P : positive semi-definite matrix

convex potential function

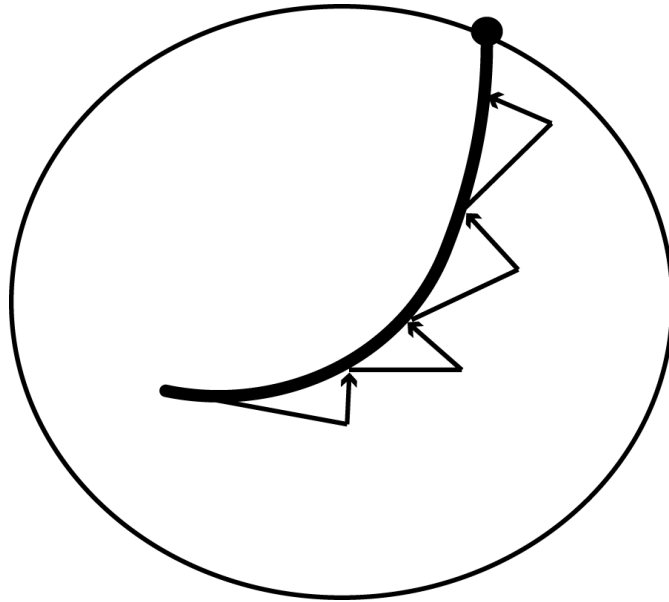
dual geodesic approach

$$Ax = b, \quad \min c \cdot x$$

Support vector machine



Polynomial-Time Algorithm



curvature : step-size

$$\left| H^{(m)} \right|^2$$

$$\min : tc \cdot x + \psi(x)$$

$$x = \delta(t)$$

∇^* – geodesic

Machine Learning

Boosting : combination of weak learners

$$D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$$

$$y_i = \pm 1$$

$$f(\mathbf{x}, \mathbf{u}) : y = h(\mathbf{x}, \mathbf{u}) = \text{sgn } f(\mathbf{x}, \mathbf{u})$$

Weak Learners

$$H(\mathbf{x}) = \text{sgn}\left(\sum \alpha_t h_t(\mathbf{x})\right)$$

$$\varepsilon_t : \text{Prob} \{h_t(\mathbf{x}_i) \neq y_i\} \quad |W_t$$

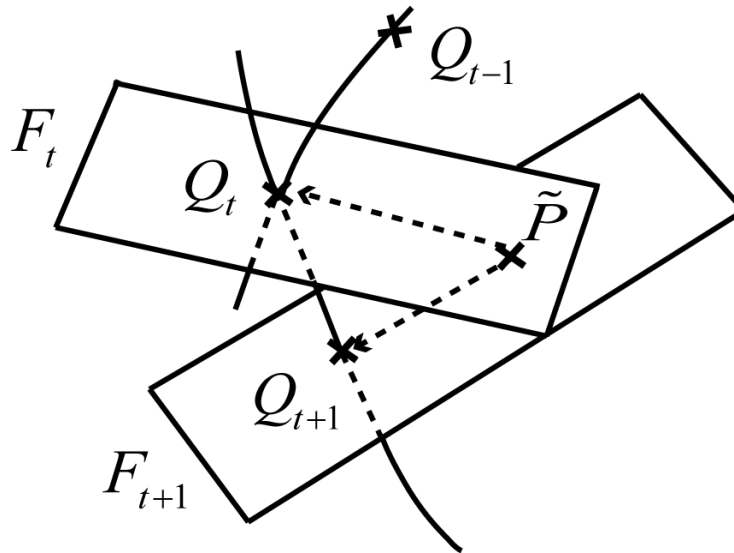
$$W_{t+1}(i) = cW_t(i) \exp\{-\alpha_t y_i h_t(x_i)\}$$

weight distribution

Boosting — generalization

$$Q_t = \left\{ Q_t(y|x) = Q_{t-1}(y|x) \exp\left\{ \alpha_t y h_t(x) - \beta_t \right\} \right\}$$

$$F_t = \left\{ P(y, x) E y h_t(x) = \text{const} \right\}$$



$$\alpha_t : \min D \left[\tilde{P} \parallel Q_t \right]$$

$$D(\tilde{P} \parallel Q_{t+1}) < D(\tilde{P} \parallel Q_t)$$

SVM : support vector machine

Embedding $z_i = \phi_i(x)$

$$f(x) = \sum w_i \phi_i(x) = \sum \alpha_i y_i K(x_i, x)$$

Kernel $K(x, x') = \sum \phi_i(x) \phi_i(x')$

Conformal change of kernel

$$K(x, x') \longrightarrow \rho(x) \rho(x') K(x, x')$$

$$\rho(x) = \exp\{-\kappa f(x)\}$$

KL-divergence, α -divergence

$$D_\alpha [p(x):q(x)] = \frac{4}{1-\alpha^2} \left\{ 1 - \int p(x)^{\frac{1-\alpha}{2}} q(x)^{\frac{1+\alpha}{2}} dx \right\}$$

$$\alpha \rightarrow -1: D_{-1}[p:q] = KL[p:q] = \int p(x) \log \frac{p(x)}{q(x)} dx$$

$$\alpha \rightarrow +1: D_1[p:q] = KL[q:p]$$

α -representation

$$f_\alpha(p) = \frac{2}{1-\alpha} \{p(x)\}^{\frac{2}{1-\alpha}}$$

$$\alpha\text{-family} : p(x, \theta)^{\frac{1-\alpha}{2}} = c \sum \theta_i q_i(x)^{\frac{1-\alpha}{2}}$$

{	exponential family	$(\alpha = -1)$	dually flat
	mixture family	$(\alpha = 1)$	

Csiszar f -divergence, U -divergence

$$D_f [p(x) : q(x)] = \int f \left(\frac{q(x)}{p(x)} \right) p(x) dx$$

$$D_U [p(x) : q(x)]$$

α -divergence from convex function

$$D_{\alpha}(P, P') = \frac{4}{1 - \alpha^2} \left[1 - \exp \left\{ \psi \left(\frac{1 - \alpha}{2} \theta + \frac{1 + \alpha}{2} \theta' \right) \right. \right. \\ \left. \left. + \psi \left(\frac{1 - \alpha}{2} \theta \right) + \psi \left(\frac{1 + \alpha}{2} \theta' \right) \right\} \right]$$

α -geodesic, duality, projection

Integration of evidences:

$$x_1, x_2, \dots, x_m$$

arithmetic mean

geometric mean

harmonic mean

α -mean