

Machine Learning SS: Kyoto U.

**Information Geometry
and Its Applications to
Machine Learning**

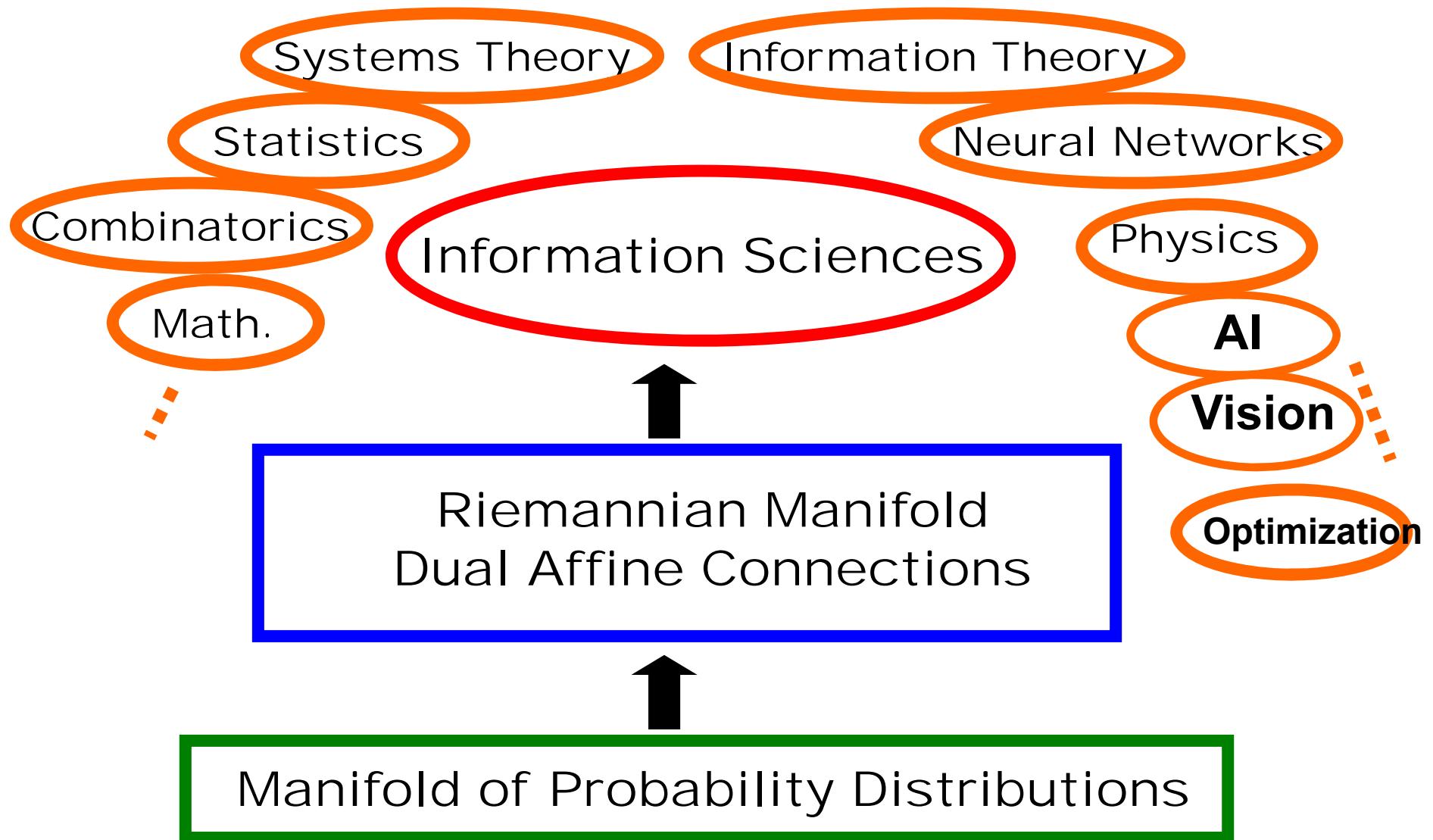
**Shun-ichi Amari
RIKEN Brain Science Institute**

Information Geometry

**-- Manifolds of
Probability Distributions**

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

Information Geometry

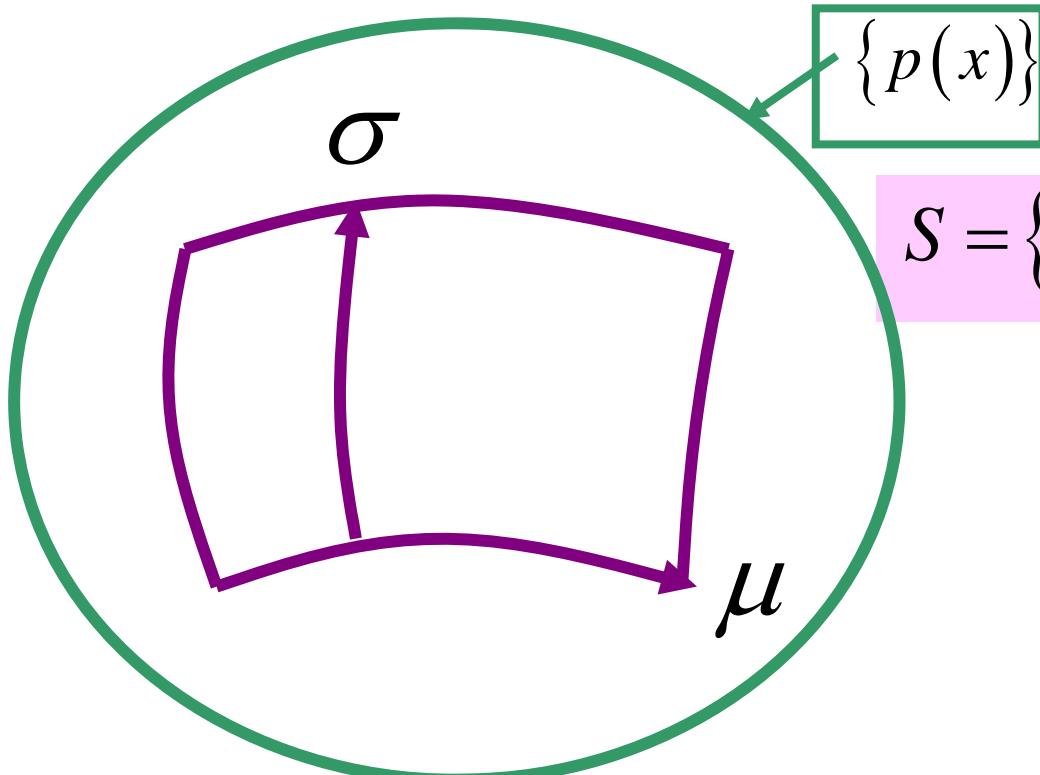


Information Geometry ?

Gaussian distributions

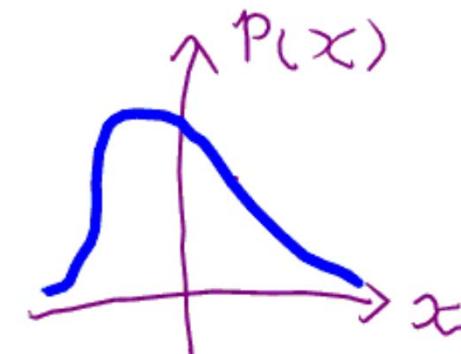
$$S = \{ p(x; \mu, \sigma) \}$$

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



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

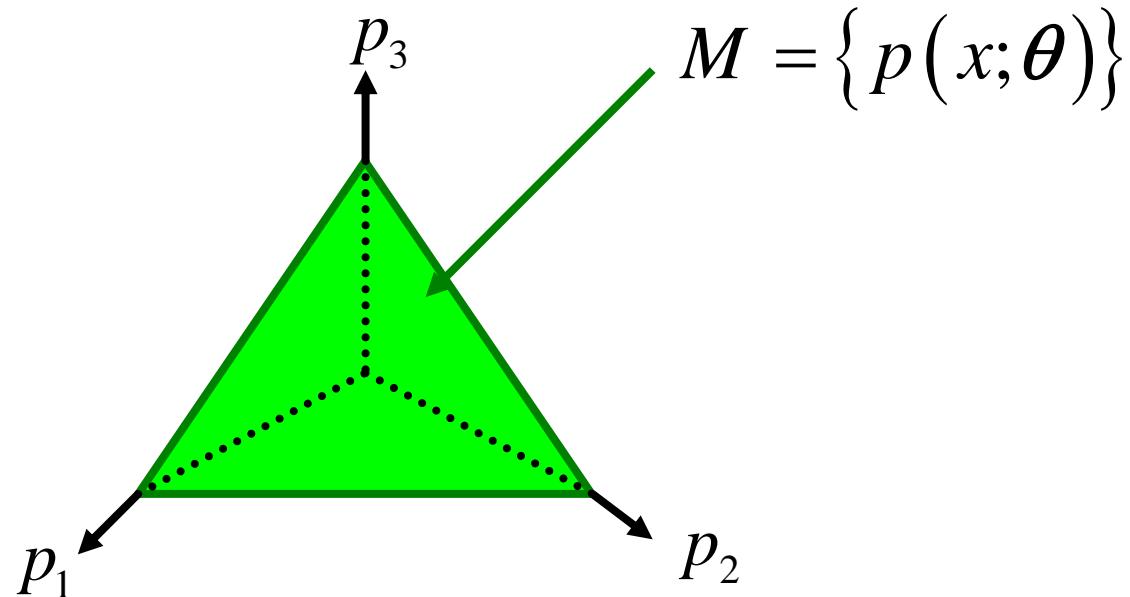
$$\theta = (\mu, \sigma)$$



Manifold of Probability Distributions

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

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



Invariance

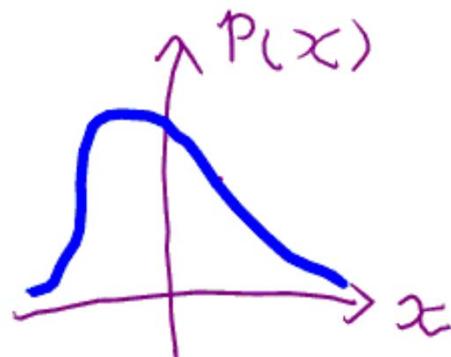
$$S = \{ p(x, \theta) \}$$

Invariant under different representation

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

sufficient statistics

$$\begin{aligned} & \int |p(x, \theta_1) - p(x, \theta_2)|^2 dx \\ & \neq \int |\bar{p}(y, \theta_1) - \bar{p}(y, \theta_2)|^2 dy \end{aligned}$$



Two Geometrical Structures

Riemannian metric

affine connection --- geodesic

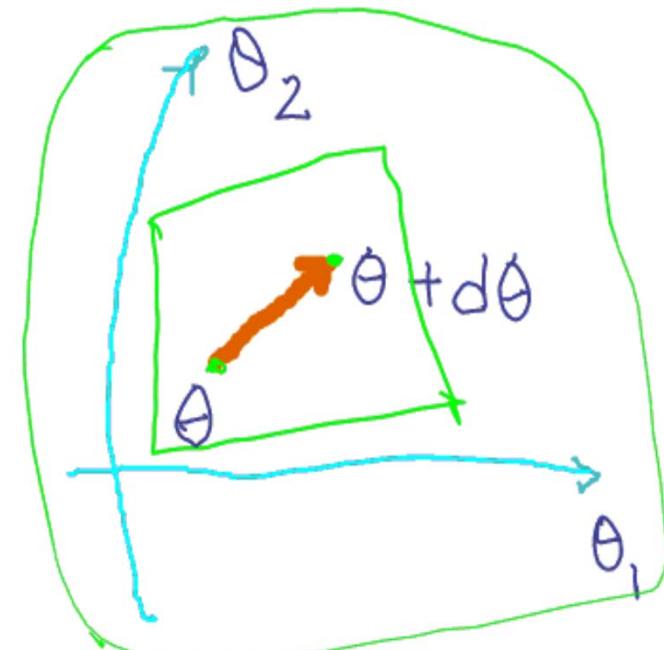
$$ds^2 = \sum g_{ij}(\theta) d\theta_i d\theta_j$$

Fisher information

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

Orthogonality: inner product

$$\langle d_1 \theta, d_2 \theta \rangle = d_1 \theta^T G d_2 \theta$$



Affine Connection

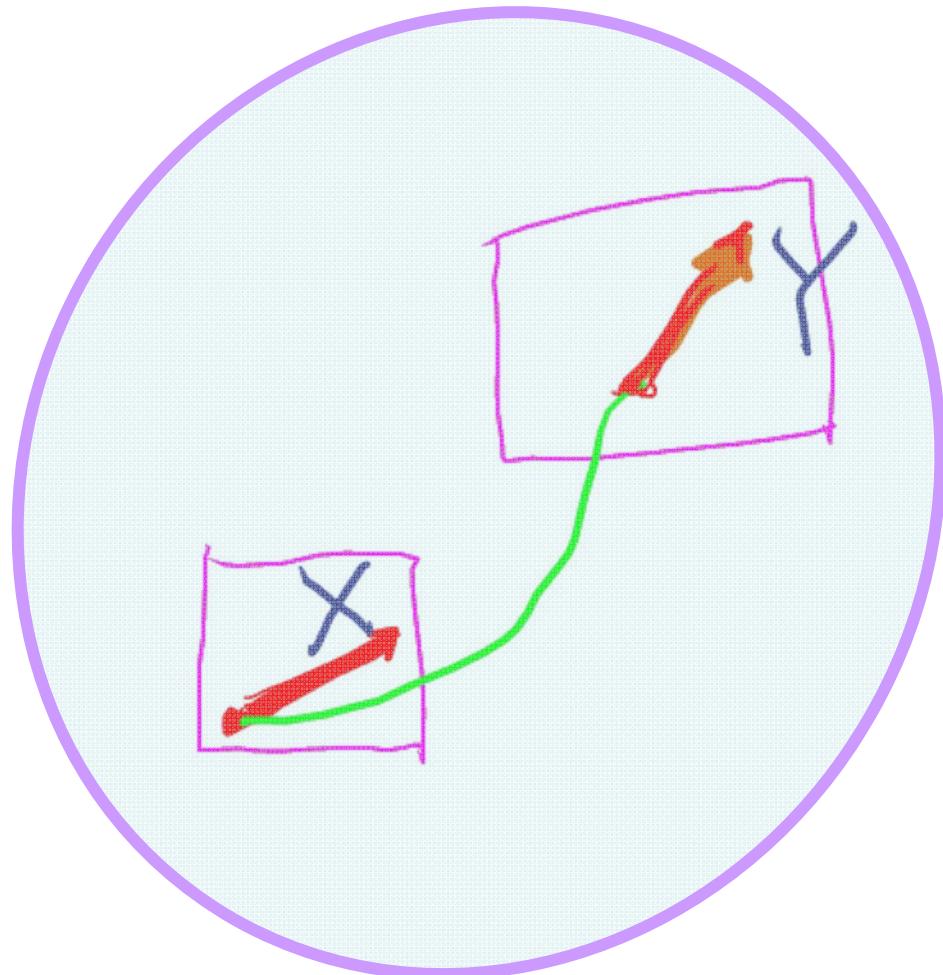
covariant derivative; parallel transport

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

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

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

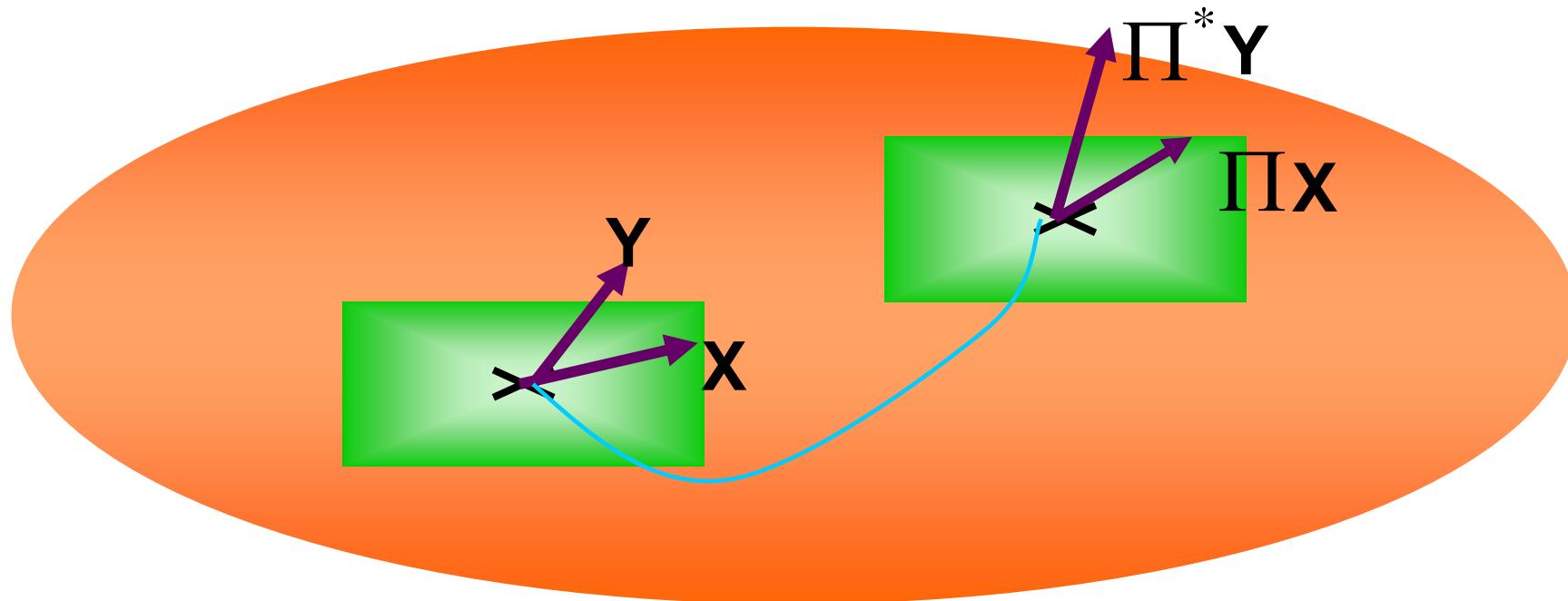
**minimal distance
straight line**



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

(∇, ∇^*)

(Π, Π^*)

e-geodesic

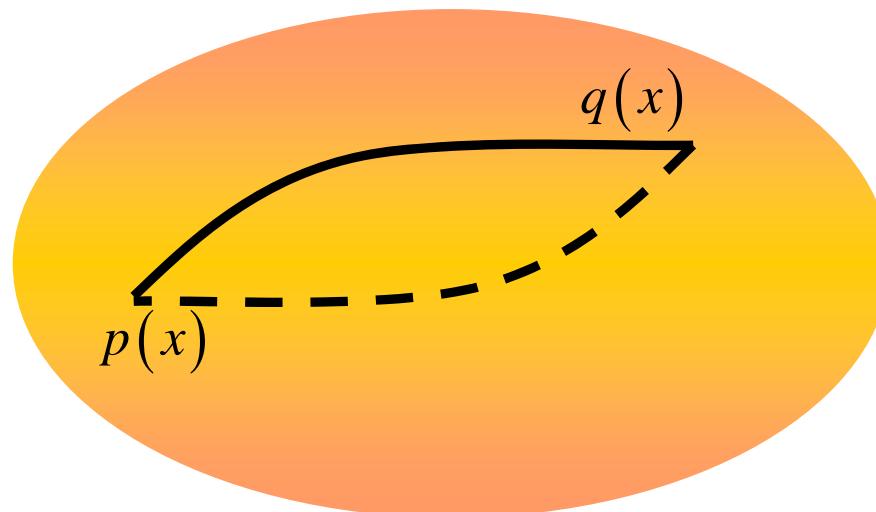
$$\log r(x, t) = t \log p(x) + (1-t) \log q(x) + c(t)$$

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

$$g_{ij}(\xi) = E[\partial_i l \partial_j l]$$

$$T_{ijk}(\xi) = E[\partial_i l \partial_j l \partial_k l]$$

$$\{\mathbf{M}, \mathbf{g}, \mathbf{T}\}$$

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

α -connection

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

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

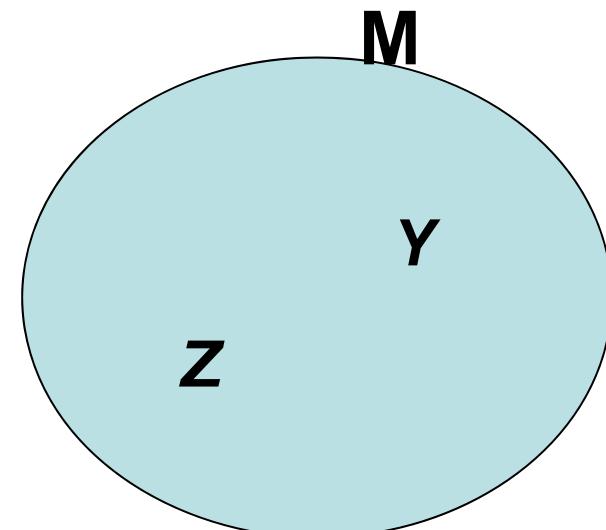
$$X \langle Y, Z \rangle = \langle \nabla_X Y, Z \rangle + \langle Y, \nabla_X^* Z \rangle$$

Divergence: $D[z:y]$

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

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

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



positive-definite

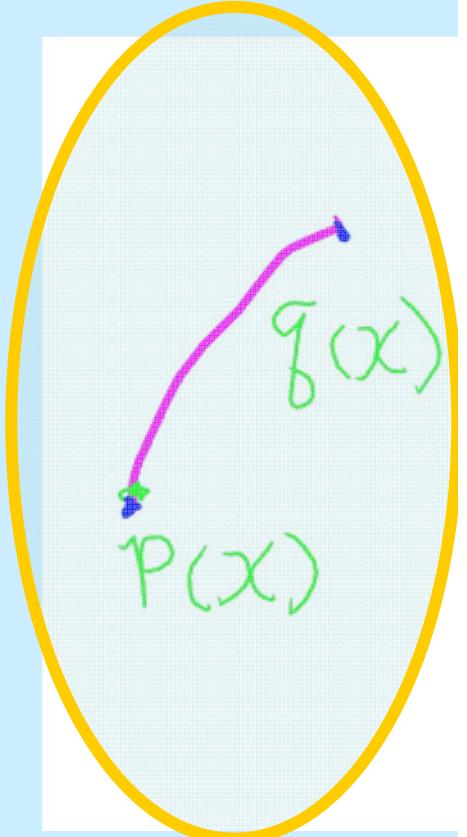
Kullback-Leibler Divergence

quasi-distance

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

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

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



f divergence of \tilde{S}

$$\tilde{S} = \{\tilde{\mathbf{p}}\}, \quad \tilde{p}_i > 0 : (\sum \tilde{p}_i = 1 \text{ nn holds})$$

$$D_f [\tilde{\mathbf{p}} : \tilde{\mathbf{q}}] = \sum \tilde{p}_i f\left(\frac{\tilde{q}_i}{\tilde{p}_i}\right) \geq 0 \quad D_f [\tilde{\mathbf{p}} : \tilde{\mathbf{q}}] = 0 \Leftrightarrow \tilde{\mathbf{p}} = \tilde{\mathbf{q}}$$

not invariant under $\tilde{f}(u) = f(u) - c(u-1)$

α divergence

$$D_\alpha[\tilde{p} : \tilde{q}] = \sum \left\{ \frac{1-\alpha}{2} \tilde{p}_i + \frac{1+\alpha}{2} \tilde{q}_i - \tilde{p}_i^{\frac{1-\alpha}{2}} \tilde{q}_i^{\frac{1+\alpha}{2}} \right\}$$

KL-divergence

$$D[\tilde{p} : \tilde{q}] = \sum \left\{ \tilde{p}_i \log \frac{\tilde{p}_i}{\tilde{q}_i} + \tilde{p}_i - \tilde{q}_i \right\}$$

(α, β) -divergence

$$D_{\alpha,\beta}[p : q] = \sum \left\{ \frac{\alpha}{\alpha + \beta} p_i^{\alpha + \beta} + \frac{\beta}{\alpha + \beta} q_i^{\alpha + \beta} - p^\alpha_i q_i^\beta \right\}$$

$\beta = -\alpha$: α -divergence

$\alpha = 1$: β -divergence

Metric and Connections Induced by Divergence

Riemannian metric

(Eguchi)

$$g_{ij}(z) = \partial_i \partial_j D[z : y]_{|y=z} : D[z : z + dz] = \frac{1}{2} g_{ij}(z) dz_i dz_j$$

affine connections

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

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

$$\Gamma_{ijk}^*(z) = -\partial_i' \partial_j' \partial_k D[z : y]_{|y=z}$$

$$\text{Duality: } X \langle Y, Z \rangle = \langle \nabla_X Y, Z \rangle + \langle Y, \nabla^*_X Z \rangle$$

$$\partial_k g_{ij} = \Gamma_{kij} + \Gamma_{kji}^*$$

$$\Gamma_{ijk} = \Gamma_{ijk}^* - T_{ijk}$$

$$\left\{M\,,g\,,T\,\right\}$$

Dually flat manifold

exponential family; mixture family; $\{p(x); x \text{ discrete}\}$

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

θ -coordinates : affine coordinates, flat, geodesics



η -coordinates: (dual) affine coordinates, flat, geodesics

not Riemannian flat

canonical divergence $D(P: P') : KL - divergence$

Dually flat manifold

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

$$\eta_i = \frac{\partial}{\partial \theta_i} \psi(\theta); \quad \theta_i = \frac{\partial}{\partial \eta_i} \varphi(\eta) : \quad \text{Legendre transformation}$$

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

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

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

ψ : cumulant generating function

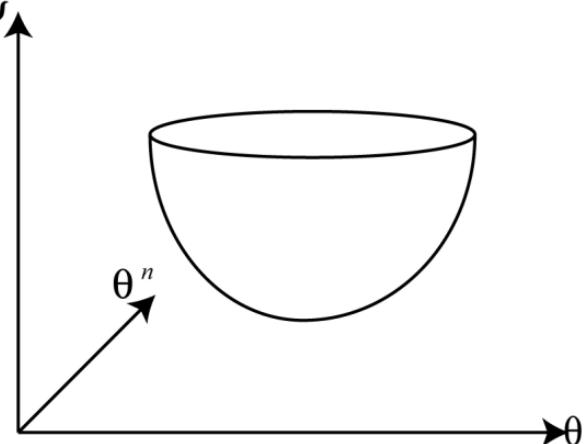
φ : negative entropy

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

Manifold with Convex Function

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

$\psi(\theta)$: convex function



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

negative entropy

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

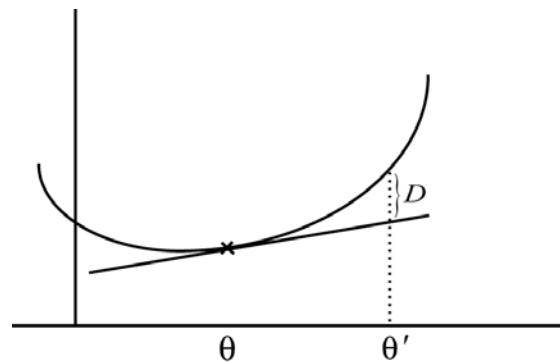
Riemannian metric and flatness

(affine structure)

$$\{S, \psi(\theta), \theta\}$$

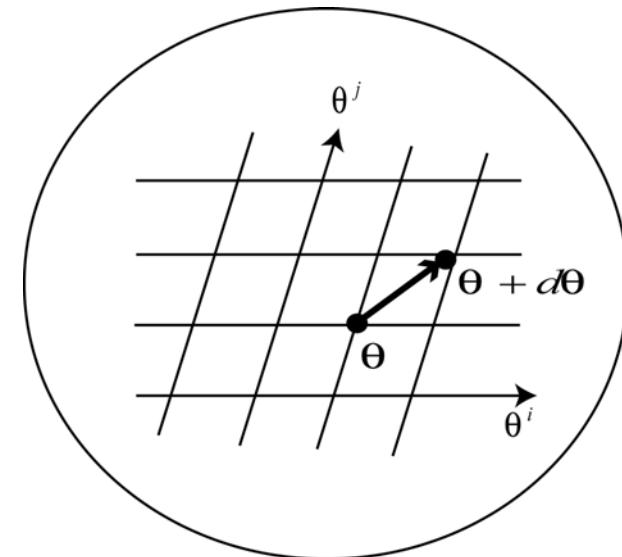
Bregman divergence

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



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

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



Flatness (affine) θ : geodesic (not Levi-Civita)

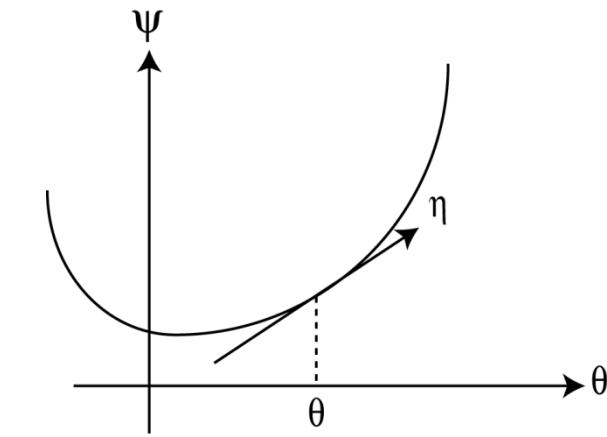
Legendre Transformation

$$\eta_i = \partial_i \psi(\theta)$$

$\theta \leftrightarrow \eta$ one-to-one

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

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



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

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

Two affine coordinate systems (θ, η)

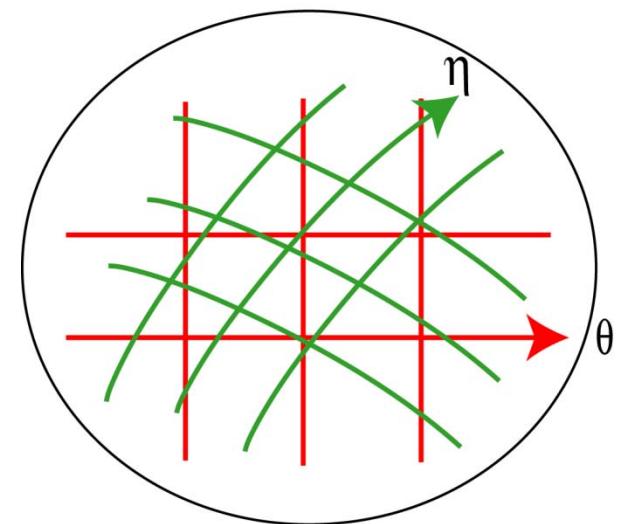
θ : geodesic (e-geodesic)

η : dual geodesic (m-geodesic)

“dually orthogonal”

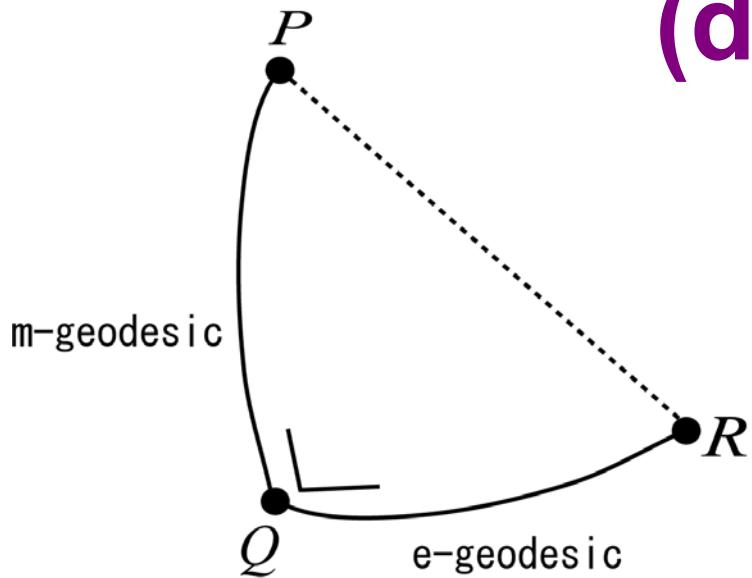
$$\langle \partial_i, \partial^j \rangle = \delta_i^j$$

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



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

Pythagorean Theorem (dually flat manifold)



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

Euclidean space: self-dual $\theta = \eta$

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

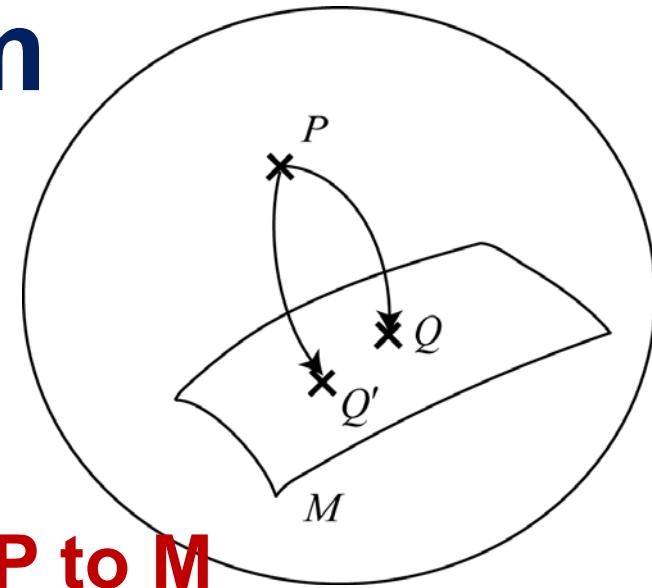
Projection Theorem

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

Q = m-geodesic projection of P to M

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

Q' = e-geodesic projection of P to M



Two Types of Divergence

Invariant divergence (Chentsov, Csiszar)

f-divergence: Fisher- α structure

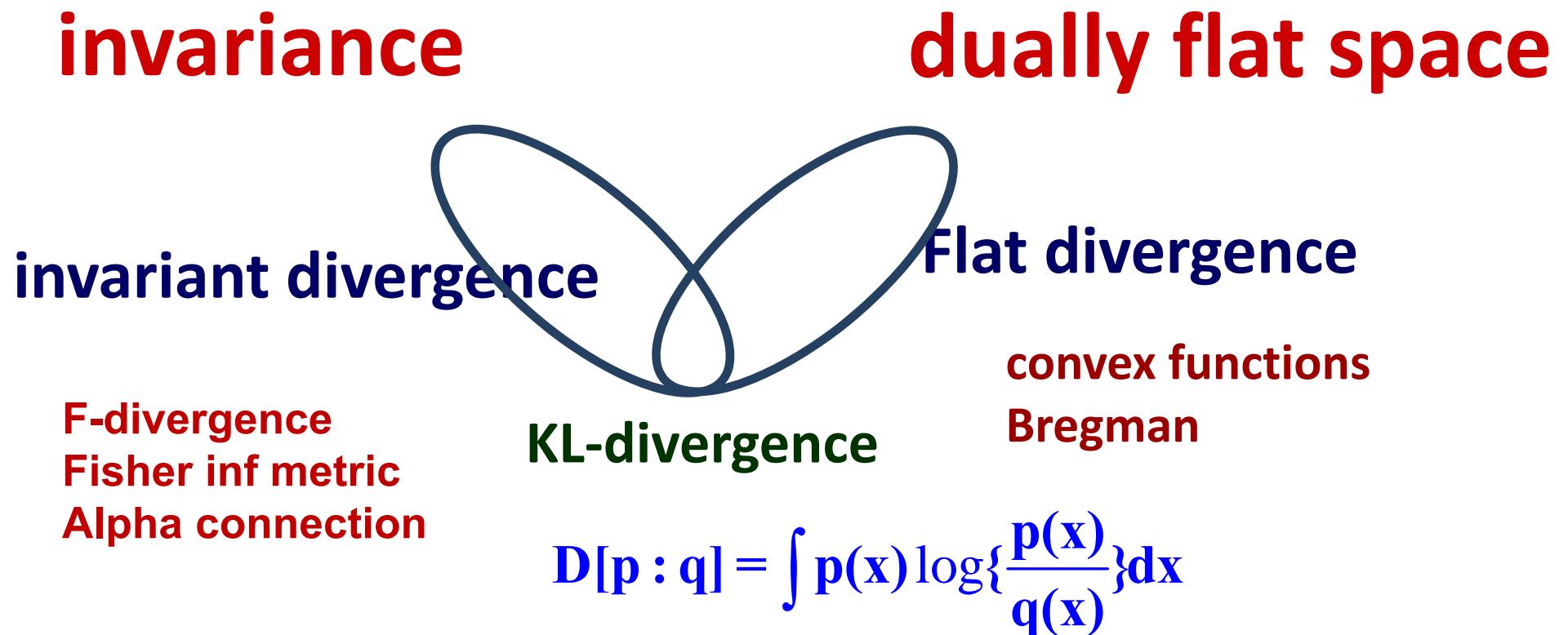
$$D[p : q] = \int p(x) f\left\{ \frac{q(x)}{p(x)} \right\} dx$$

Flat divergence (Bregman) – convex function

**KL-divergence belongs to
both classes: flat and invariant**

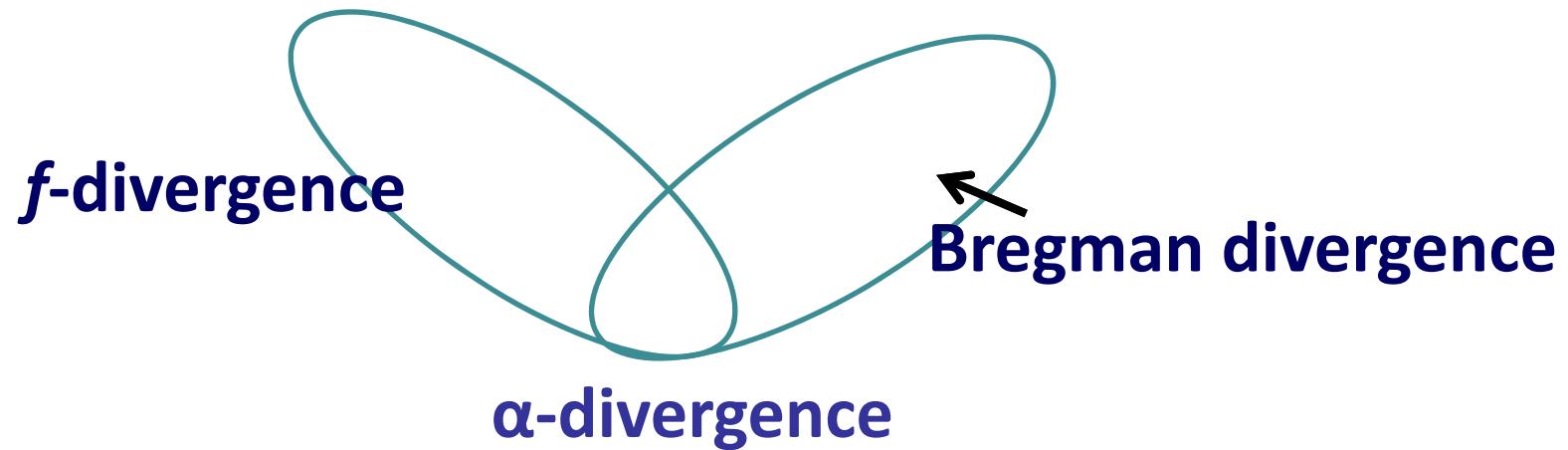
divergence

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



Space of positive measures : vectors, matrices, arrays

$$\tilde{S} = \{\tilde{p}\}, \quad \tilde{p}_i > 0 : (\sum \tilde{p}_i = 1 \text{ } nn \text{ holds})$$



Applications of Information Geometry

Statistical Inference
Machine Learning and AI
Computer Vision
Convex Programming
Signal Processing (ICA; Sparse)
Information Theory, Systems Theory
Quantum Information Geometry

Applications to Statistics

curved exponential family:

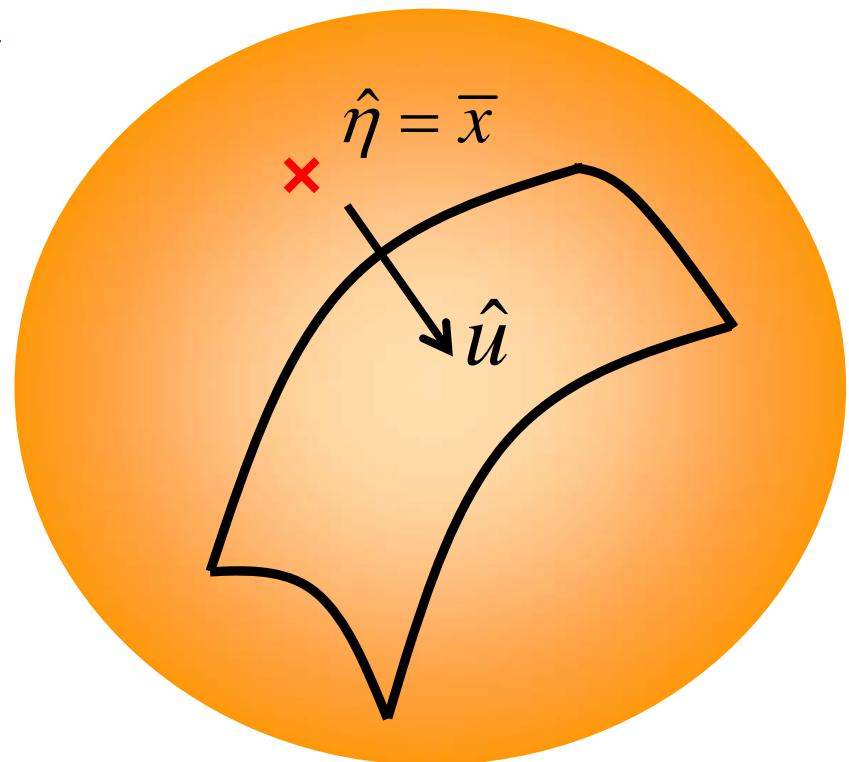
$$p(x, u) \square x_1, x_2, \dots, x_n$$

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

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

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

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



x : discrete X = {0, 1, ..., n}

$S_n = \{p(x) \mid x \in X\}$: **exponential family**

$$p(x) = \sum_{i=0}^n p_i \delta_i(x) = \exp\left[\sum_{i=1}^n \theta^i x_i - \psi(\theta)\right]$$

$$\theta^i = \log p_i - \log p_0; \quad x_i = \delta_i(x); \quad \psi(\theta) = -\log p_0$$

statistical model : $p(x, u)$

High-Order Asymptotics

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

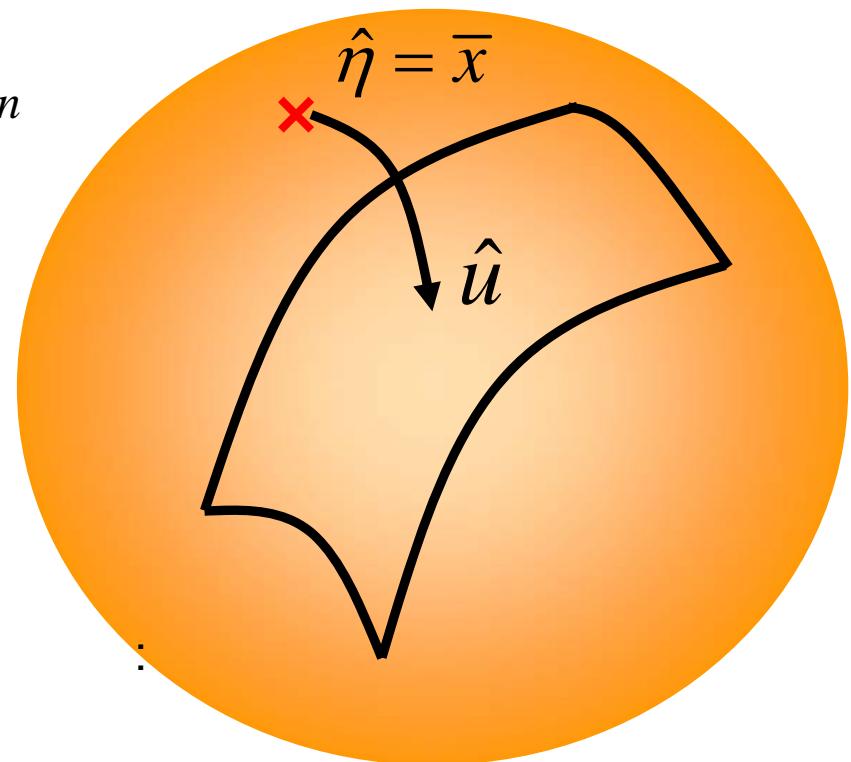
$$\hat{u} = u(x_1, \dots, x_n)$$

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

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

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

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



Information Geometry of Belief Propagation

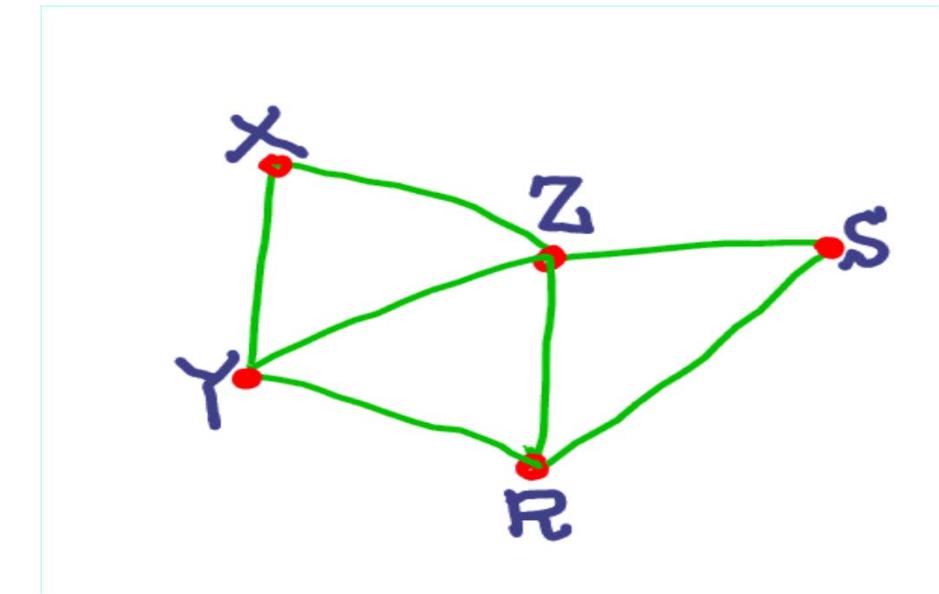
- Shun-ichi Amari (RIKEN BSI)
- Shiro Ikeda (Inst. Statist. Math.)
- Toshiyuki Tanaka (Kyoto U.)

Stochastic Reasoning

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

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

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



Stochastic Reasoning

$q(x_1, x_2, x_3, \dots | \text{observation})$

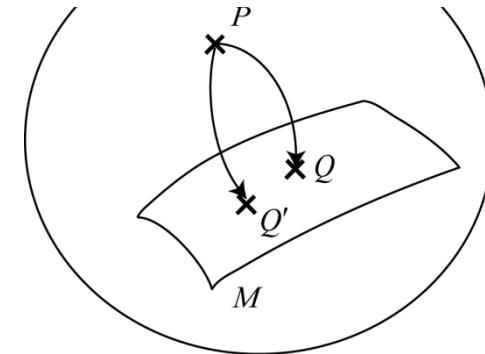
$X = (x_1 \ x_2 \ x_3 \ \dots)$ $x = 1, -1$

$X = \operatorname{argmax} q(x_1, x_2, x_3, \dots)$ maximum likelihood

$X_i = \operatorname{sgn} E[x_i]$ least bit error rate estimator

Mean Value

Marginalization:
projection to independent distributions



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

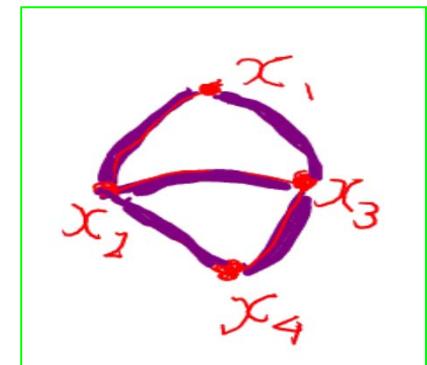
$$q_i(x_i) = \int q(x_1, \dots, x_n) dx_1 \dots d\bar{x}_i \dots dx_n$$

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

$$q(\mathbf{x}) = \exp \left\{ \sum k_i \cdot x_i + \sum_{r=1}^L c_r(\mathbf{x}) - \psi_q \right\}$$

$$c_r(\mathbf{x}) = c_r x_{i_1} \cdots x_{i_s}, \quad r = (i_1 \cdots i_s)$$

$$x_i = \{1, -1\} \quad r = (i_1, i_2)$$



Boltzmann machine, spin glass, neural networks
 Turbo Codes, LDPC Codes

Computationally Difficult

$$q(\mathbf{x}) \rightarrow \eta = E[\mathbf{x}]$$

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

- | **mean-field approximation**
- | **belief propagation**
- | **tree propagation, CCCP (convex-concave)**

Information Geometry of Mean Field Approximation

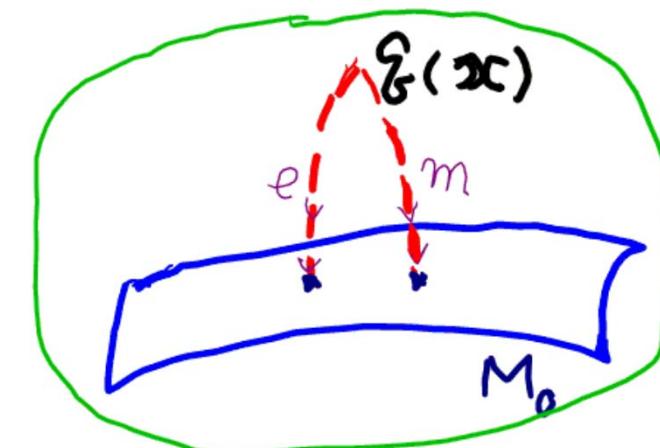
- m-projection
- e-projection

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

$$\Pi_0^m q = \operatorname{argmin} D[q:p]$$

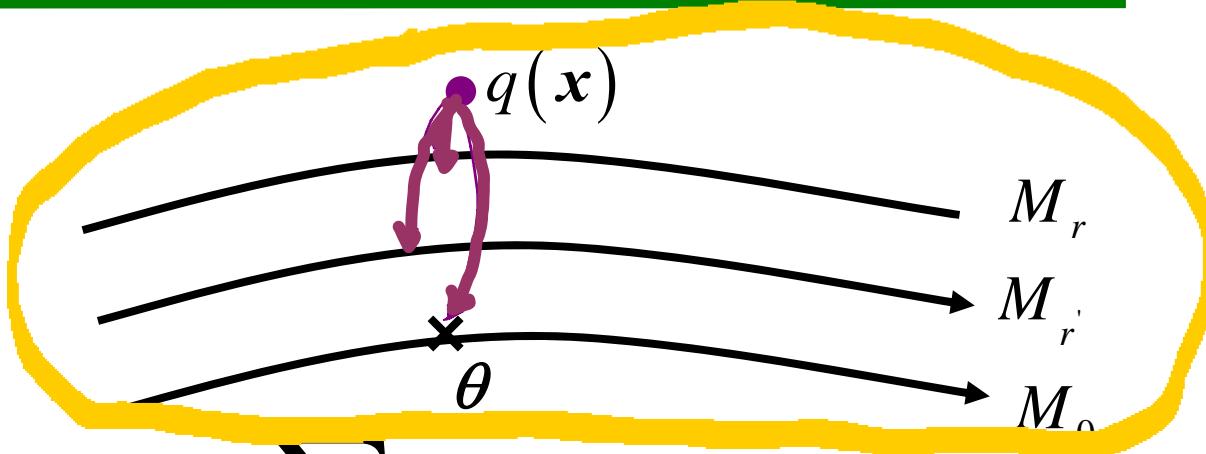
$$\Pi_0^e q = \operatorname{argmin} D[p:q]$$

$$p(x) \in M_0$$



$$M_0 = \{\Pi_i p_i(x_i)\}$$

Information Geometry



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

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

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

$r = 1, \dots, L$

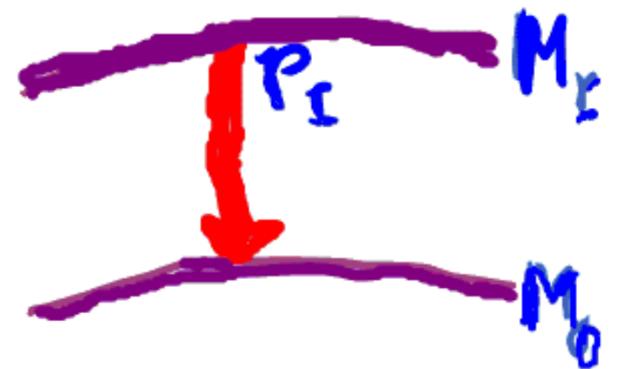
Belief Propagation

$$M_r : p_r(x, \xi_r) = \exp \left\{ c_r(x) + \xi_r \cdot x - \psi_r \right\}$$

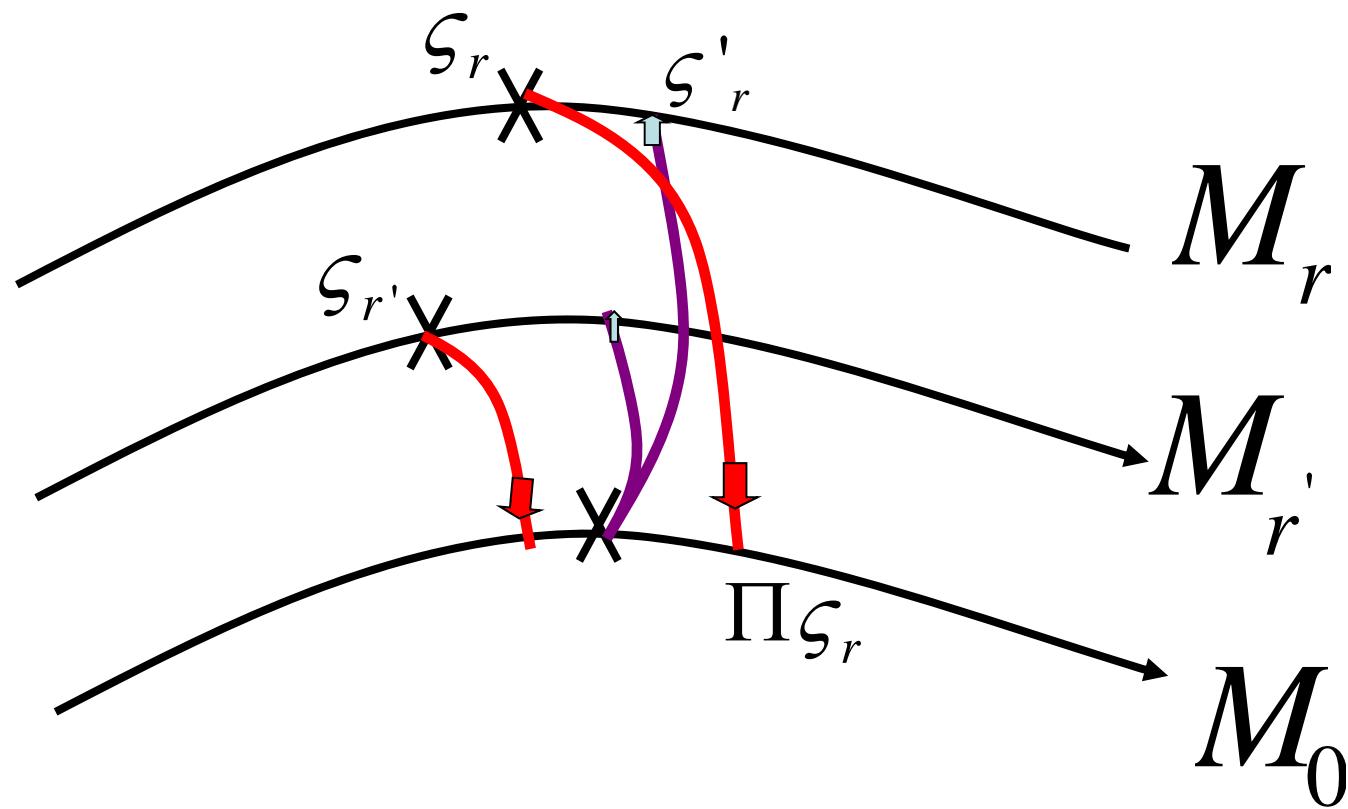
$$\Pi_0 p_r(x, \xi_r^t)$$

$$\theta_r^{t+1} = \Pi_0 p_r(x, \xi_r^t) - \xi_r^t \quad : \text{ belief for } c_r(x)$$

$$\theta^{t+1} = \sum \theta_r^{t+1}$$



Belief Prop Algorithm



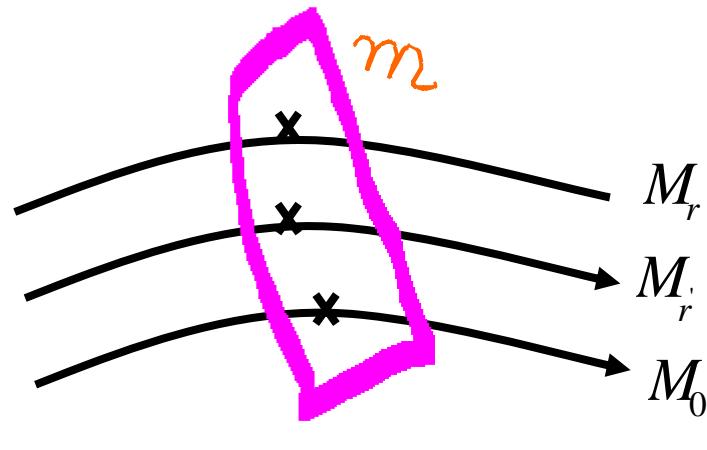
Equilibrium of BP (θ^*, ξ_r^*)

1) m -condition

$$\theta^* = \Pi_0 p_r(x, \xi_r^*)$$

m -flat submanifold $M(\theta^*)$

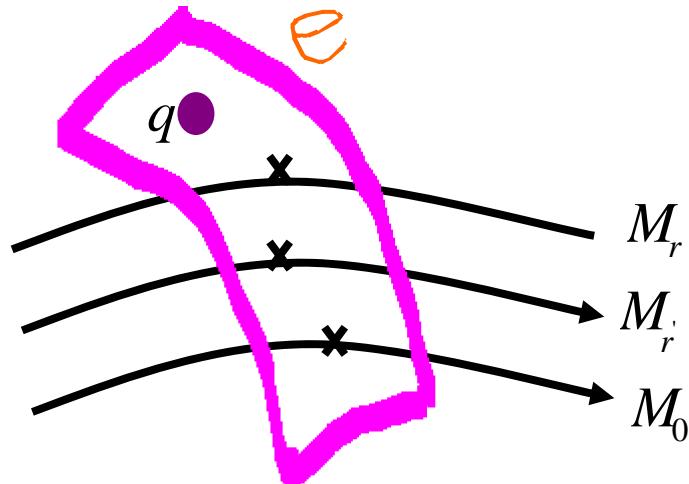
$\xi_1 \oplus \dots$



2) e -condition

$$\theta^* = \frac{1}{L-1} \sum r^* \xi_r^*$$

$q(x) \in e$ -flat submanifold



Free energy:

$$F(\theta, \zeta_1, \dots, \zeta_L) = D[p_0 : q] - \sum D[p_0 : p_r]$$

critical point

$$\frac{\partial F}{\partial \theta} = 0 \quad : e\text{-condition}$$

$$\frac{\partial F}{\partial \zeta_r} = 0 \quad : m\text{-condition}$$

not convex

Belief Propagation

e-condition OK

$$(\theta; \xi_1, \xi_2, \dots, \xi_L), \quad \theta' = \frac{1}{L-1} \sum \xi'_r$$

$$(\xi_1, \xi_2, \dots, \xi_L) \rightarrow (\xi'_1, \xi'_2, \dots, \xi'_L)$$

CCCP

m-condition OK

$$\theta \rightarrow \theta'$$

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

$$\theta' = \Pi_0 p_r(x, \xi'_r); \quad \xi' = \Pi_r p_0(x, \theta')$$

$$\xi_r^{t+1} = \Pi_r\theta^{t+1} = \Pi_r p_0\left(x,\theta^{t+1}\right)$$

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

Convex-Concave Computational Procedure (CCCP) Yuille

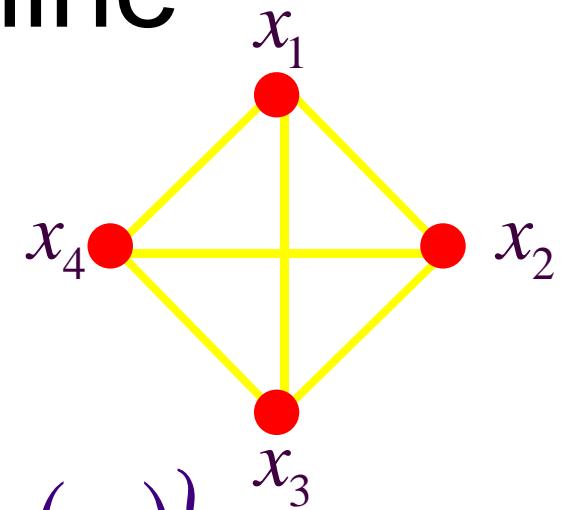
$$F(\theta) = F_1(\theta) - F_2(\theta)$$

$$\nabla F_1(\theta^{t+1}) = \nabla F_2(\theta^t)$$

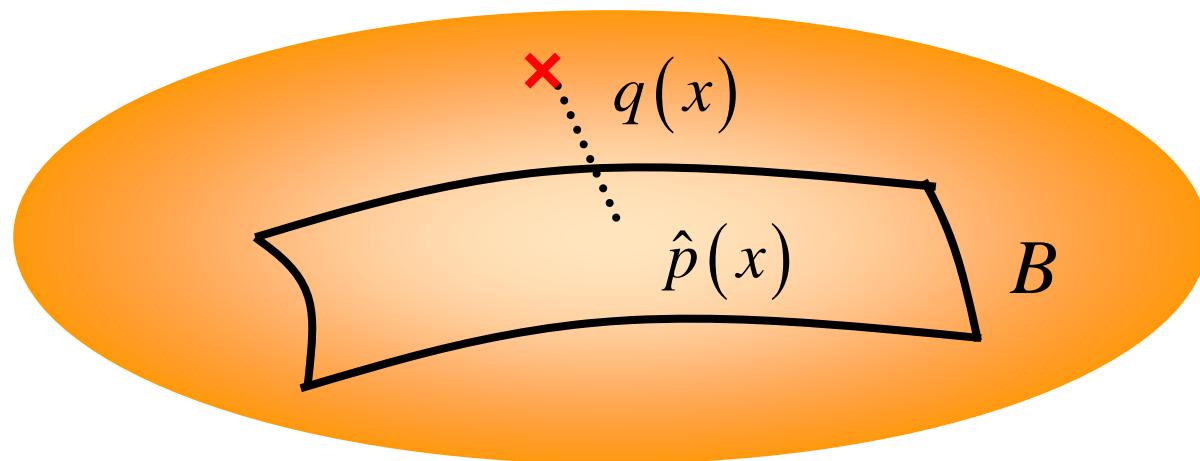
Elimination of double loops

Boltzmann Machine

$$p(x_i = 1) = \varphi\left(\sum w_{ij}x_j - h_i\right)$$

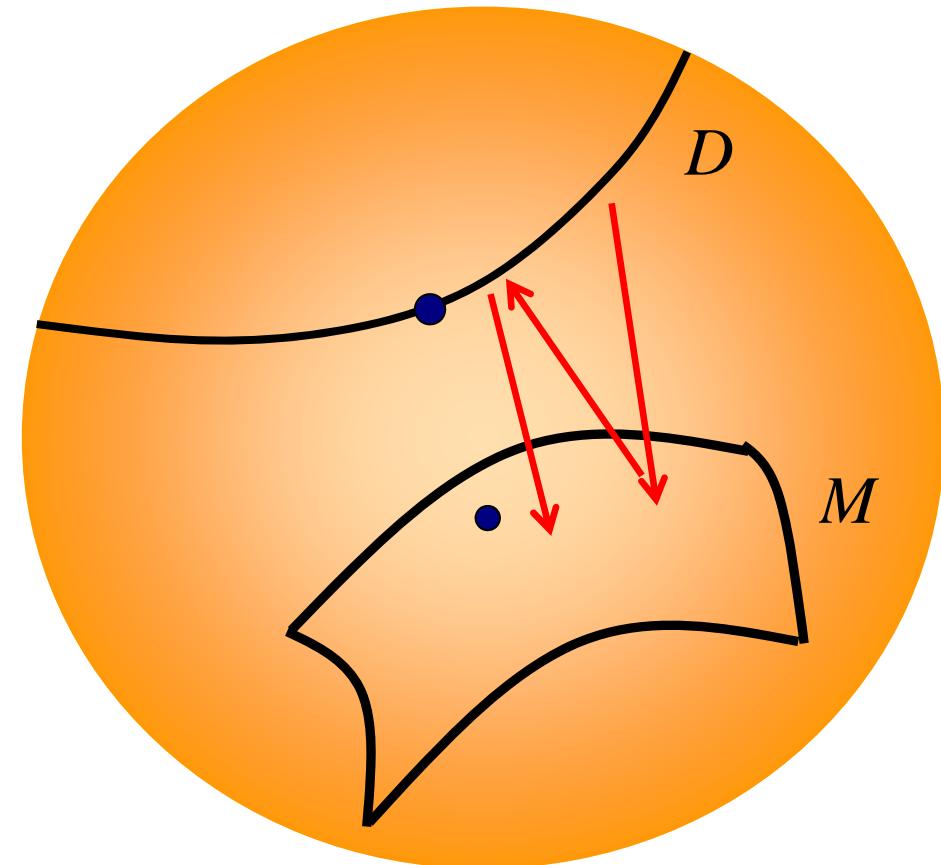


$$p(x) = \exp\left\{\sum w_{ij}x_i x_j - \sum h_i x_i - \psi(w)\right\}$$



Boltzmann machine ---hidden units

- EM algorithm
- e-projection
- m-projection



EM algorithm

hidden variables

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

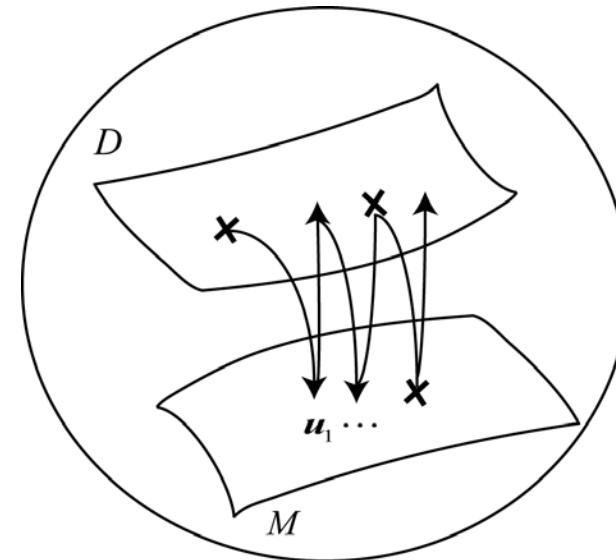
$$D = \{\boldsymbol{x}_1, \dots, \boldsymbol{x}_N\}$$

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

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

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

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



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

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

$$\rho(x) = \exp\{-\kappa |f(x)|^2\}$$

Signal Processing

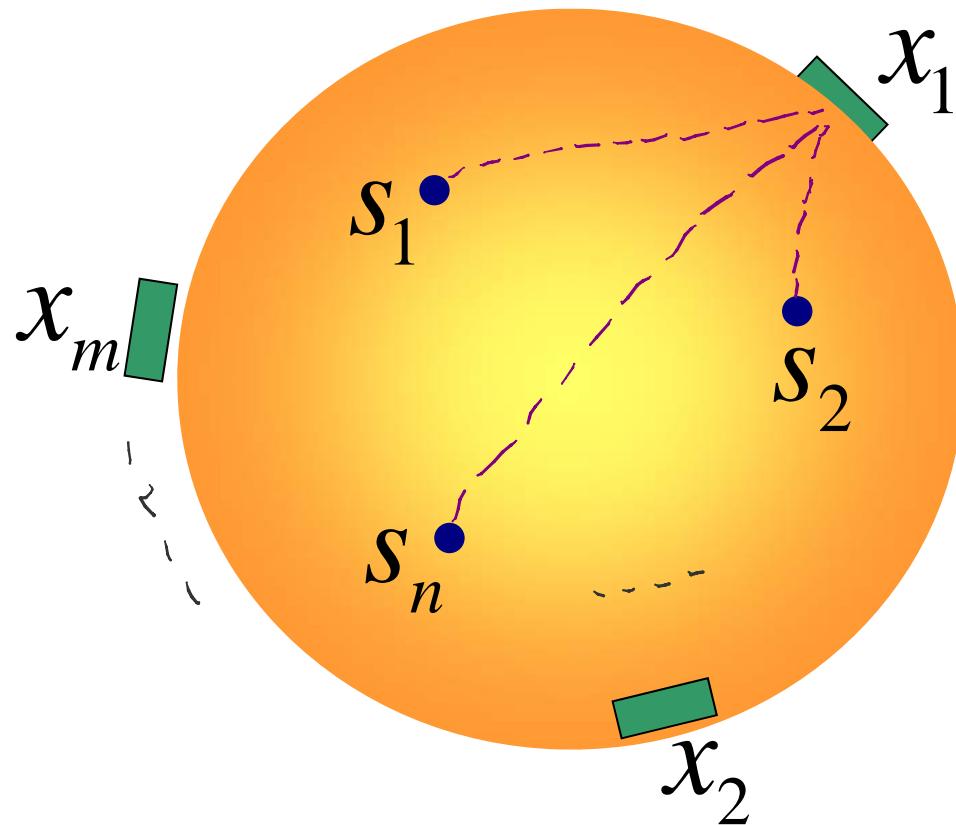
ICA : Independent Component Analysis

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

sparse component analysis

positive matrix factorization

mixture and unmixture of independent signals



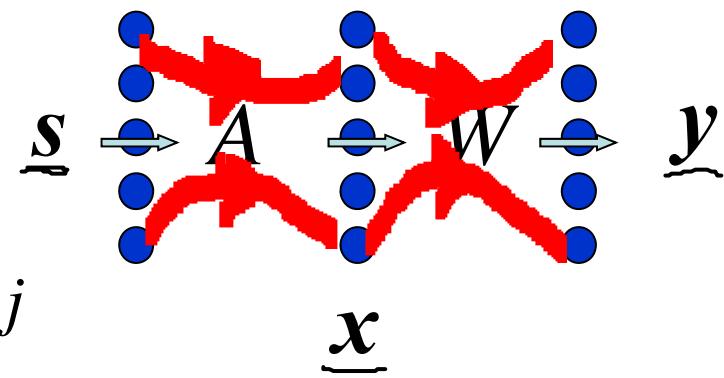
$$x_i = \sum_{j=1}^n A_{ij} s_j$$

$$\mathbf{x} = \mathbf{As}$$

Independent Component Analysis

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

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



observations: $x(1), x(2), \dots, x(t)$
recover: $s(1), s(2), \dots, s(t)$

Example of color image separation :



Five original images (but unknown to the neural net)

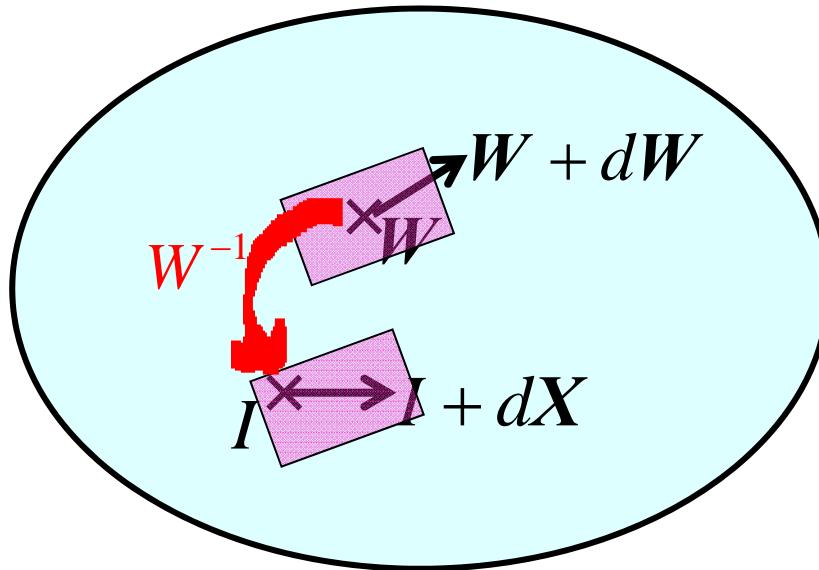


Five mixed images for separation



Final (stable states) of five separated images

Space of Matrices : Lie group



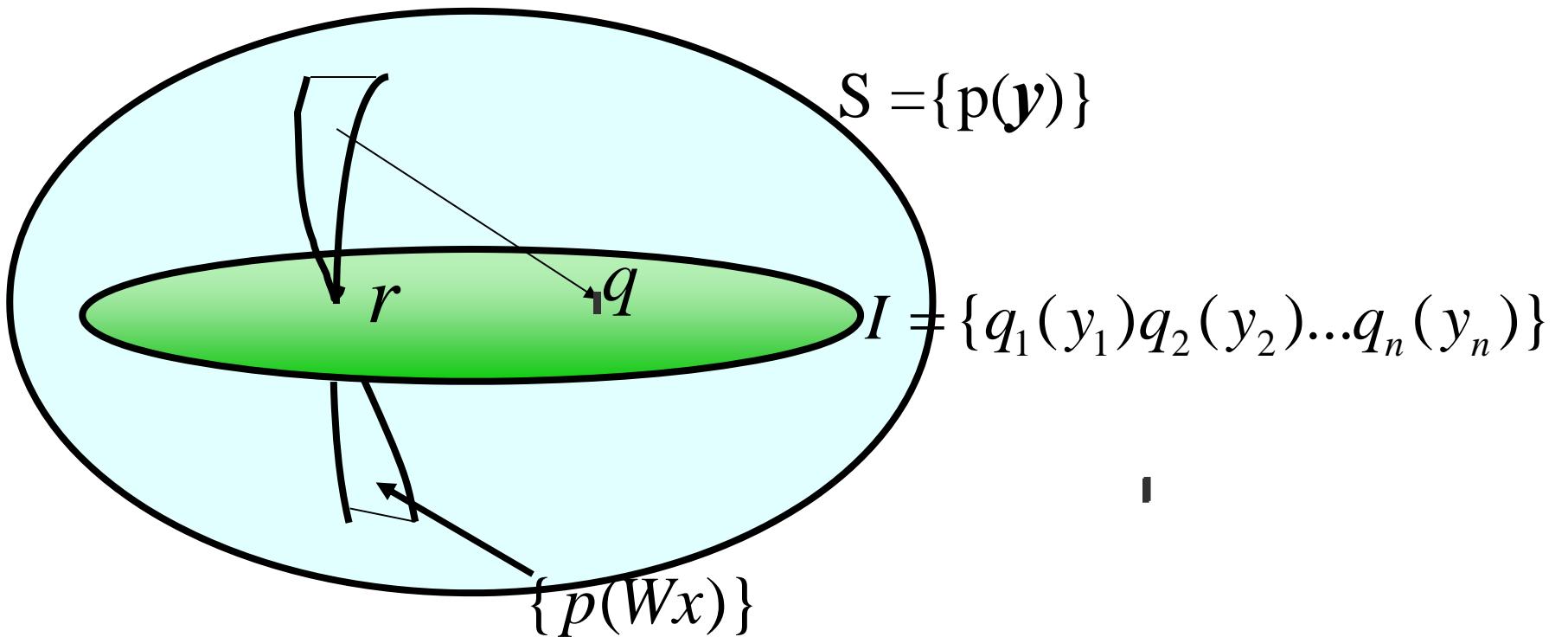
$$dX = dWW^{-1}$$

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

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

dX : **non-holonomic basis**

Information Geometry of ICA



**natural gradient
estimating function
stability, efficiency**

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

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

Semiparametric Statistical Model

$$p(\mathbf{x}; \mathbf{W}, r) = | \mathbf{W} | r(\mathbf{Wx})$$

$$\mathbf{W} = \mathbf{A}^{-1}, \quad r(s): \text{ unknown} \quad \quad r = \prod r_i$$

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

Natural Gradient

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

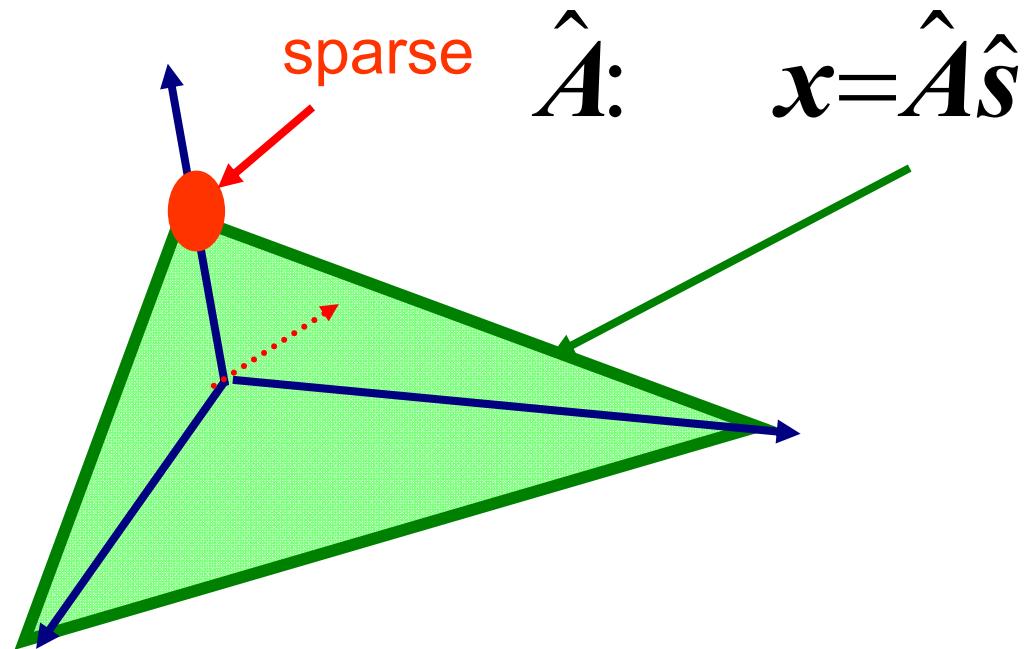
Basis Given: overcomplete case Sparse Solution

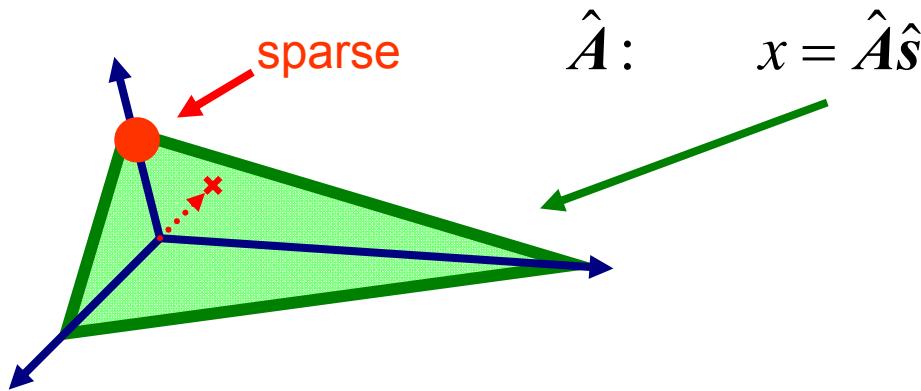
$$\mathbf{x} = \mathbf{A}\mathbf{s} = \sum s_i \mathbf{a}_i$$

many solutions

many $s_i \rightarrow 0$

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



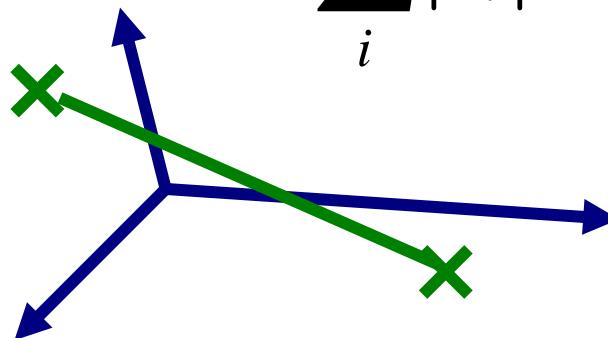


generalized inverse

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

sparse solution

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



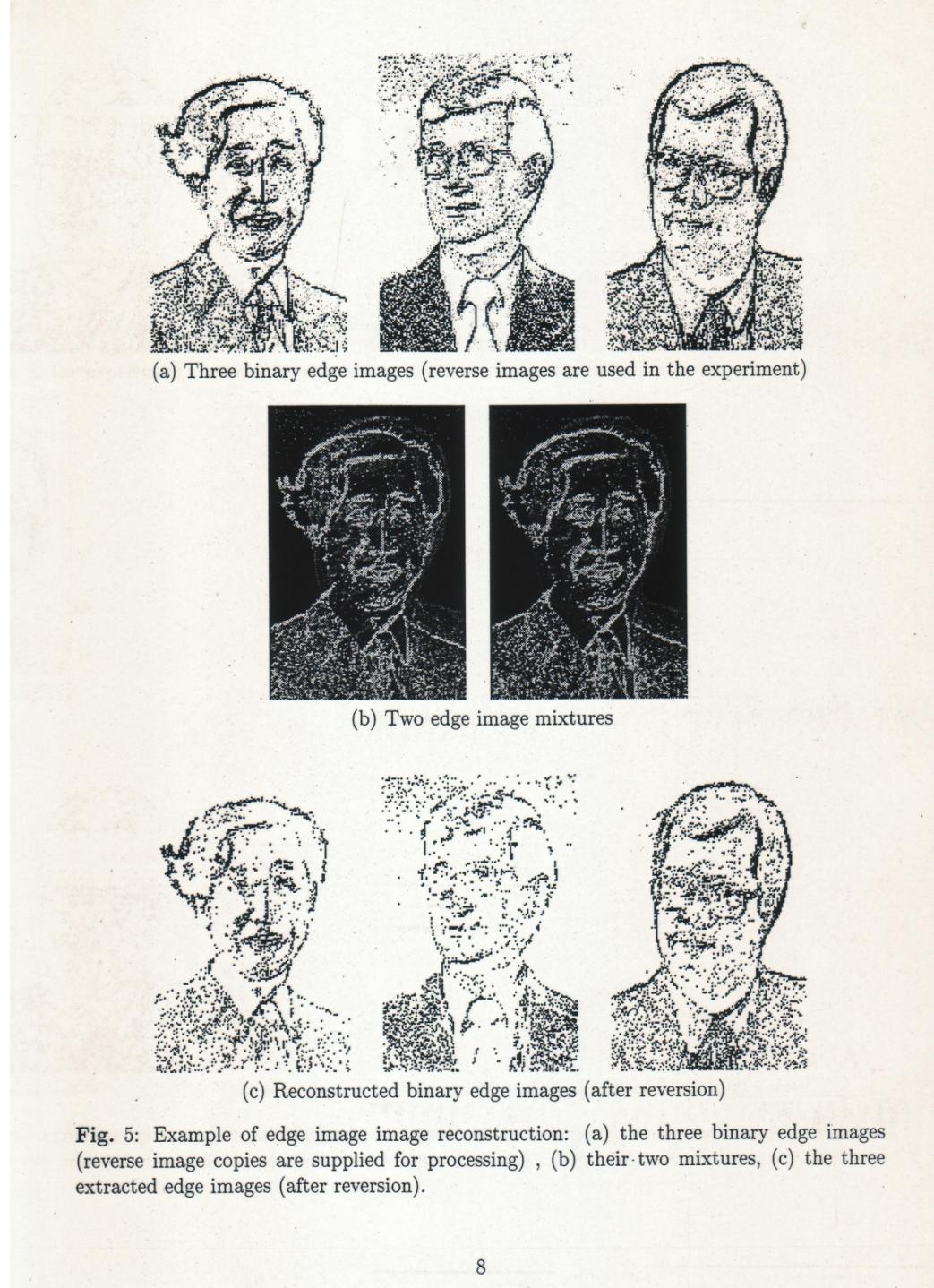


Fig. 5: Example of edge image image reconstruction: (a) the three binary edge images (reverse image copies are supplied for processing) , (b) their two mixtures, (c) the three extracted edge images (after reversion).

Overcomplete Basis and Sparse Solution

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

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

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

non-linear denoising

Sparse Solution

$$\min \varphi(\beta)$$

penalty $F_p(\beta) = \sum |\beta_i|^p$: Bayes prior

$F_0(\beta) = \#\{i : \beta_i \neq 0\}$: **sparsest solution**

$F_1(\beta) = \sum |\beta_i|$: **L_1 solution**

$F_p(\beta)$: $0 \leq p \leq 1$ **Sparse solution: overcomplete case**

$F_2(\beta) = \sum |\beta_i|^2$: **generalized inverse solution**

Optimization under Sparsity Condition

$$\begin{cases} \min \varphi(\beta) & : \text{convex function} \\ \text{constraint} & F(\beta) \leq c \end{cases}$$

typical case: $\varphi(\beta) = \frac{1}{2} |y - X\beta|^2 = \frac{1}{2} (\beta - \beta^*)^T G (\beta - \beta^*)$

$$F(\beta) = \frac{1}{p} \sum |\beta_i|^p ; \quad p = 2, \quad p = 1, \quad p = 1/2$$

L1-constrained optimization

P_c Problem

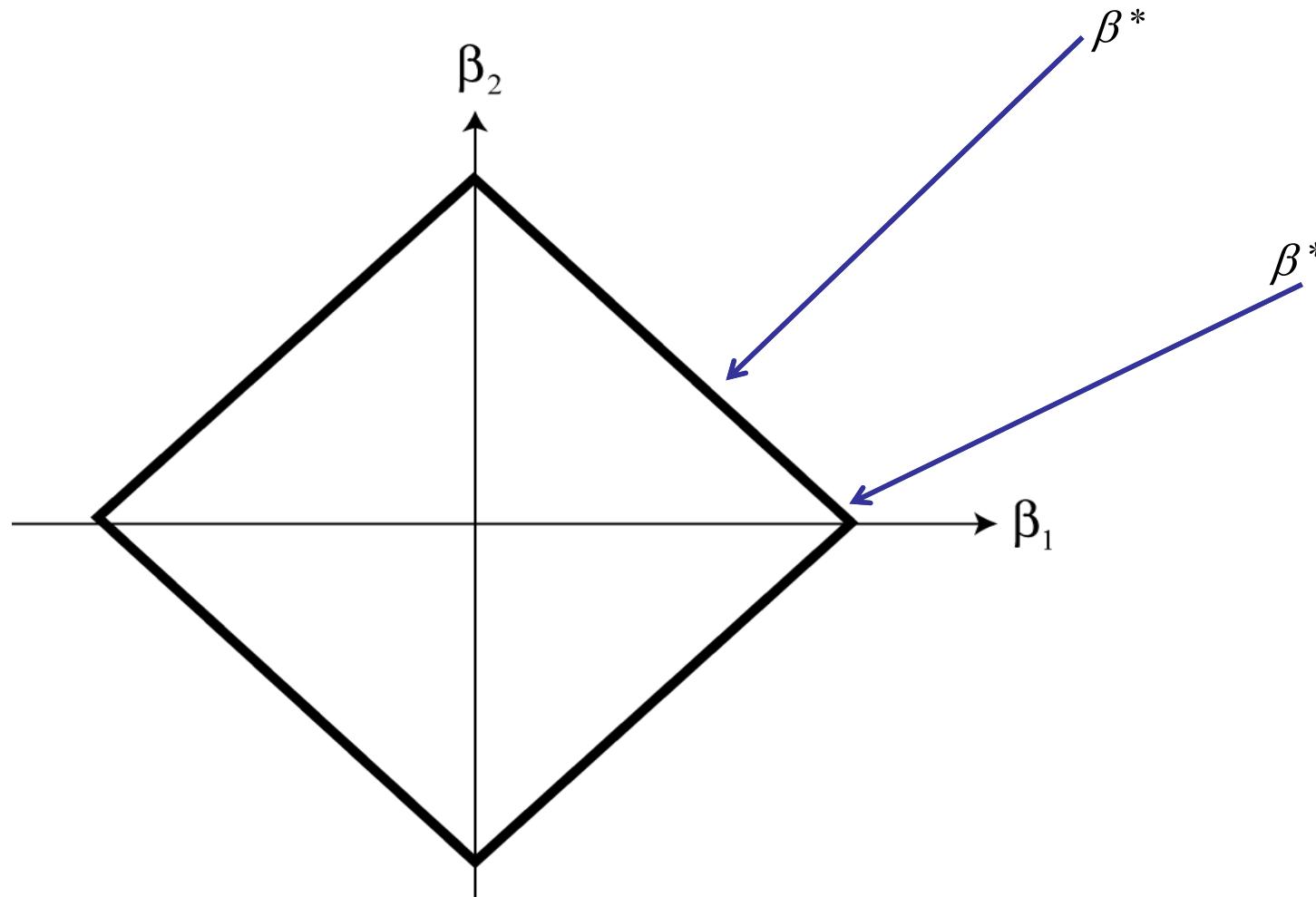
$$\begin{aligned} \min \varphi(\beta) & \quad \text{under } F(\beta) \leq c && \textbf{LASSO} \\ \text{solution } \beta^*(c) & : c = 0 \rightarrow \infty \\ & \beta_c^* = 0 \rightarrow \beta^* \end{aligned}$$

P_λ Problem

$$\begin{aligned} \min \varphi(\beta) + \lambda F(\beta) & && \textbf{LARS} \\ \text{solution } \beta^*(\lambda) & : \lambda = \infty \rightarrow 0 \\ & \beta_\lambda^* = 0 \rightarrow \beta^* \end{aligned}$$

solutions β_c^* and β_λ^* : coincide, $\lambda = \lambda(c)$, $p \geq 1$
 $p < 1$: $\lambda = \lambda(c)$ multiple, noncontinuous
stability different

Projection from β^* to $F = c$ (information geometry)



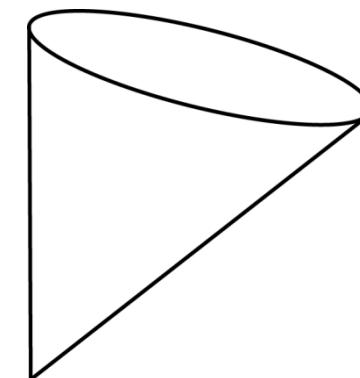
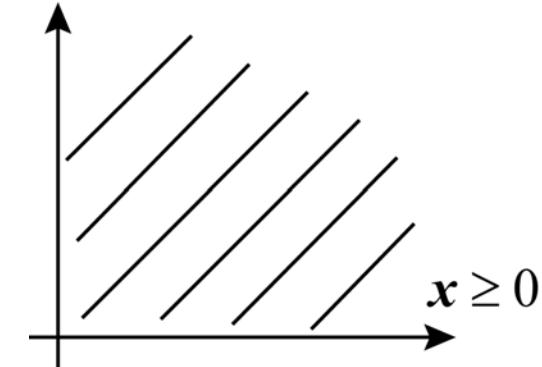
Convex Cone Programming

P : positive semi-definite matrix

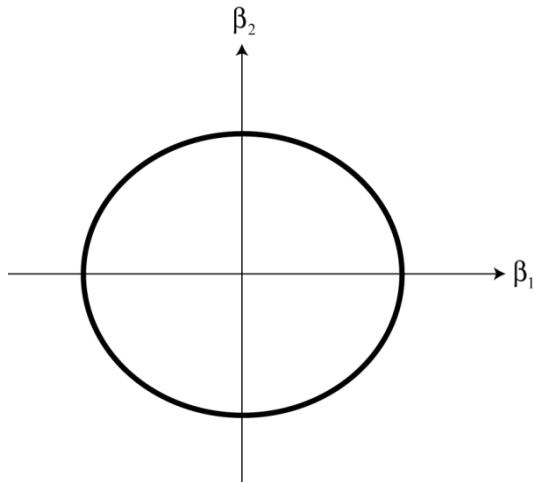
convex potential function

dual geodesic approach

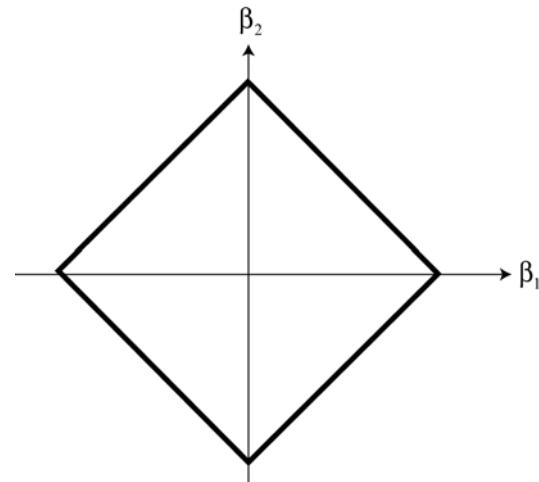
$$Ax = b, \quad \min c \cdot x$$



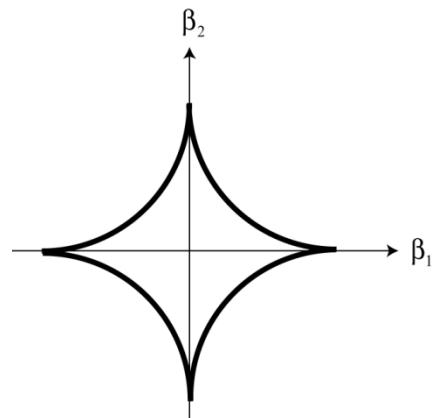
Support vector machine



a) $R_c : n = 2, p > 1$



b) $R_c : n = 2, p = 1$



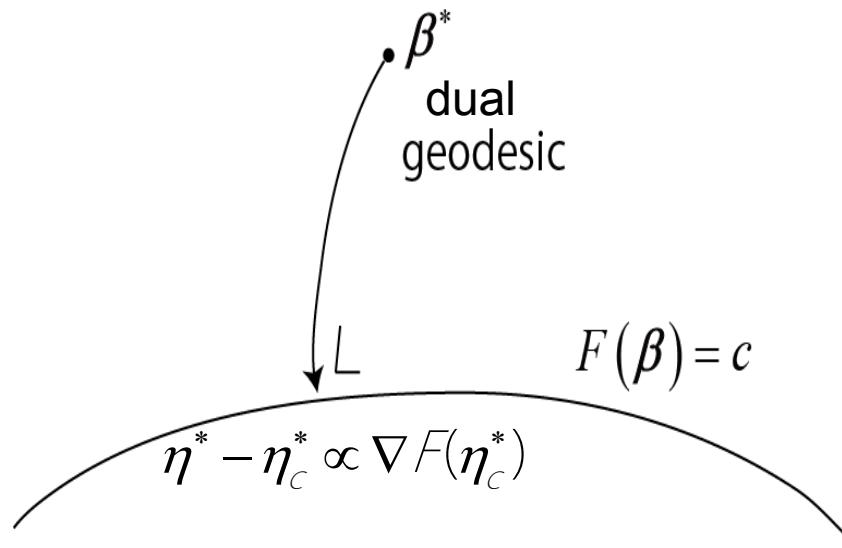
c) $R_c : n = 2, p < 1$

non-convex

Fig. 1

orthogonal projection, dual projection

$\min \varphi(\beta) = \mathbf{D}[\beta : \beta^*], \quad F(\beta) = c : \text{dual geodesic projection}$



$$\eta_c^* \propto \nabla F(\eta_c^*)$$

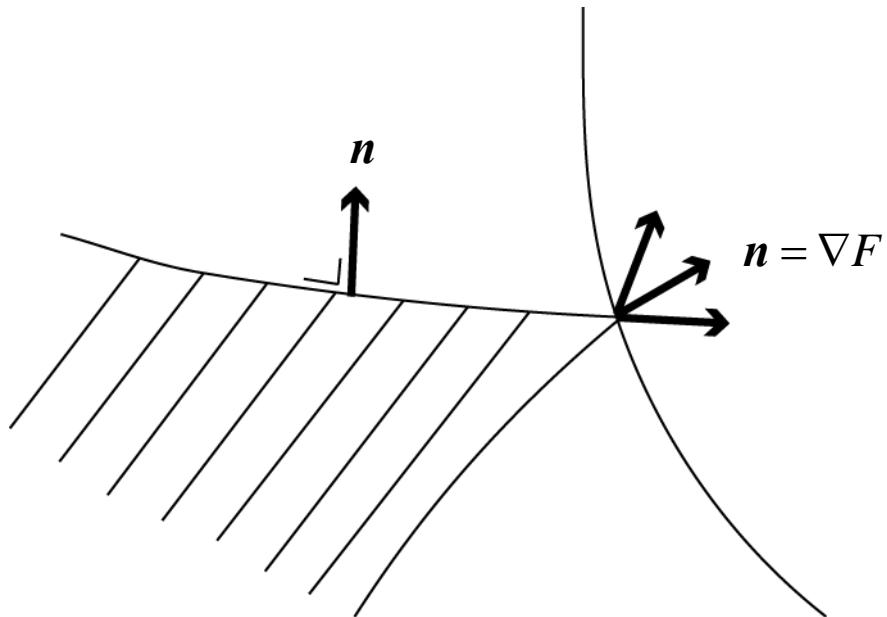


Fig. 5 subgradient

LASSO path and LARS path (stagewise solution)

$$\min \varphi(\beta) : F(\beta) = c$$

$$\min \varphi(\beta) + \lambda F(\beta)$$

$$\beta^*(c), \beta^*(\lambda) \quad \text{c} \Leftrightarrow \lambda \text{ correspondence}$$

Active set and gradient

$$A(\beta) = \{i \mid \beta_i \neq 0\}$$

$$\nabla F_p(\beta) = \begin{cases} \text{sgn}(\beta_i) |\beta_i|^{-(1-p)}, & i \in A \\ (-\infty, \infty), & i \notin A \\ [-1, 1] \end{cases}$$

Solution path

$$\nabla_A \varphi(\boldsymbol{\beta}_c^*) + \lambda_c \nabla_A F(\boldsymbol{\beta}_c^*) = 0, \quad \boldsymbol{\beta}_c^*$$

$$\left\{ \nabla_A \nabla_A \varphi(\boldsymbol{\beta}_c^*) + \lambda_c \nabla_A \nabla_A F(\boldsymbol{\beta}_c^*) \right\} \cdot \dot{\boldsymbol{\beta}}_c = -\dot{\lambda}_c \nabla_A F(\boldsymbol{\beta}_c)$$

$$\dot{\boldsymbol{\beta}}_c = -\dot{\lambda}_c K^{-1} \nabla_A F(\boldsymbol{\beta}_c^*) \quad ; \quad \dot{\boldsymbol{\beta}}_c = \frac{d}{dc} \boldsymbol{\beta}_c$$

$$K = G(\boldsymbol{\beta}_c^*) + \lambda_c \nabla \nabla F(\boldsymbol{\beta}_c^*)$$

$$(\nabla \nabla F_1 = 0; \quad \nabla F_1 = (\operatorname{sgn} \beta_i) : L_1)$$

Solution path in the subspace of the active set

$$\nabla_A \phi(\beta_\lambda^*) + \lambda \nabla_A F(\beta_\lambda^*) = 0 \quad \nabla_A : \text{active direction}$$

$$\dot{\beta}_\lambda^* = -K_A^{-1} \nabla_A F(\beta_\lambda^*)$$

turning point $A \rightarrow A'$

Gradient Descent Method

$$\min L(x+a): \quad g_{ij}a^i a^j = \varepsilon^2$$

$$\nabla L = \left\{ \frac{\partial}{\partial x_i} L(x) \right\}: \quad \text{covariant}$$

$$\tilde{\nabla} L = \left\{ \sum g^{ji} \frac{\partial}{\partial x_i} L(x) \right\}: \quad \text{contravariant}$$

$$x_{t+1} = x_t - c \nabla L(x_t)$$

Extended LARS ($p = 1$) and Minkovskian grad

$$\text{norm } \|\mathbf{a}\|_p = \sum |a_i|^p$$

$$\max \psi(\boldsymbol{\beta} + \varepsilon \mathbf{a}) \quad \text{under } \|\mathbf{a}\|_p = 1$$

$$\psi(\boldsymbol{\beta} + \varepsilon \mathbf{a}) - \lambda \|\mathbf{a}\|_p$$

$$p = 1^+$$

$$\nabla_1 \psi(\boldsymbol{\beta}) = \mathbf{1}_A \begin{cases} \operatorname{sgn} \eta_i, & |\eta_i| = \max \{|\eta_1|, \dots, |\eta_N|\} \\ 0, & \text{otherwise} \end{cases}$$

$$\boldsymbol{\eta} = \nabla \psi(\boldsymbol{\beta})$$

$$i^* = \arg \max |f_i|$$

$$\max |f_i| = |f_{i^*}| = |f_{j^*}|$$

$$\left(\tilde{\nabla}F\right)_i=\begin{cases} 1,&\text{for } i=i^*\text{ and } j^*,\\ 0&\text{otherwise.}\end{cases}$$

$$\beta_{t+1} = \beta_t - \eta \tilde{\nabla} F \quad \quad \quad \textcolor{red}{\textbf{LARS}}$$

$$\tilde{\nabla}F = \nabla f \quad \text{Euclidean case}$$

$$\tilde{\nabla}F = c(\operatorname{sgn} f_i) |f_i|^{\frac{1}{p-1}}$$

$$\tilde{\nabla}F = c(\operatorname{sgn} f_{i^*}) \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \alpha \rightarrow 1$$

L1/2 constraint: non-convex optimization

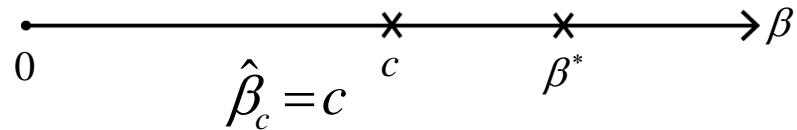
λ -trajectory and -trajectory

Ex. **1-dim**

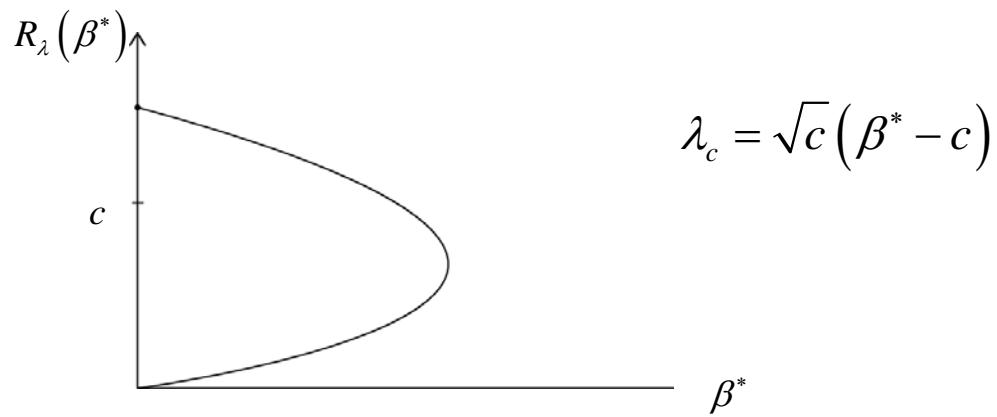
$$\varphi(\beta) = \frac{1}{2}(\beta - \beta^*)^2$$

$$f_\lambda(\beta) = \phi + \lambda F = \frac{1}{2}(\beta - 2)^2 + 2\lambda\sqrt{\beta}$$

$$P_c : \min(\beta - \beta^*)^2, \quad |\beta| \leq c$$



$$P_\lambda : \nabla f_\lambda = 0 \quad \beta - \beta^* + \frac{\lambda}{\sqrt{\beta}} = 0 \quad \hat{\beta} = R_\lambda(\beta^*) \quad : \text{ Xu Zongben's operator}$$



λ

ICCN-Huangshan(黃山)

Sparse Signal Analysis

Shun-ichi Ammari (甘利俊一)

RIKEN Brain Science Institute

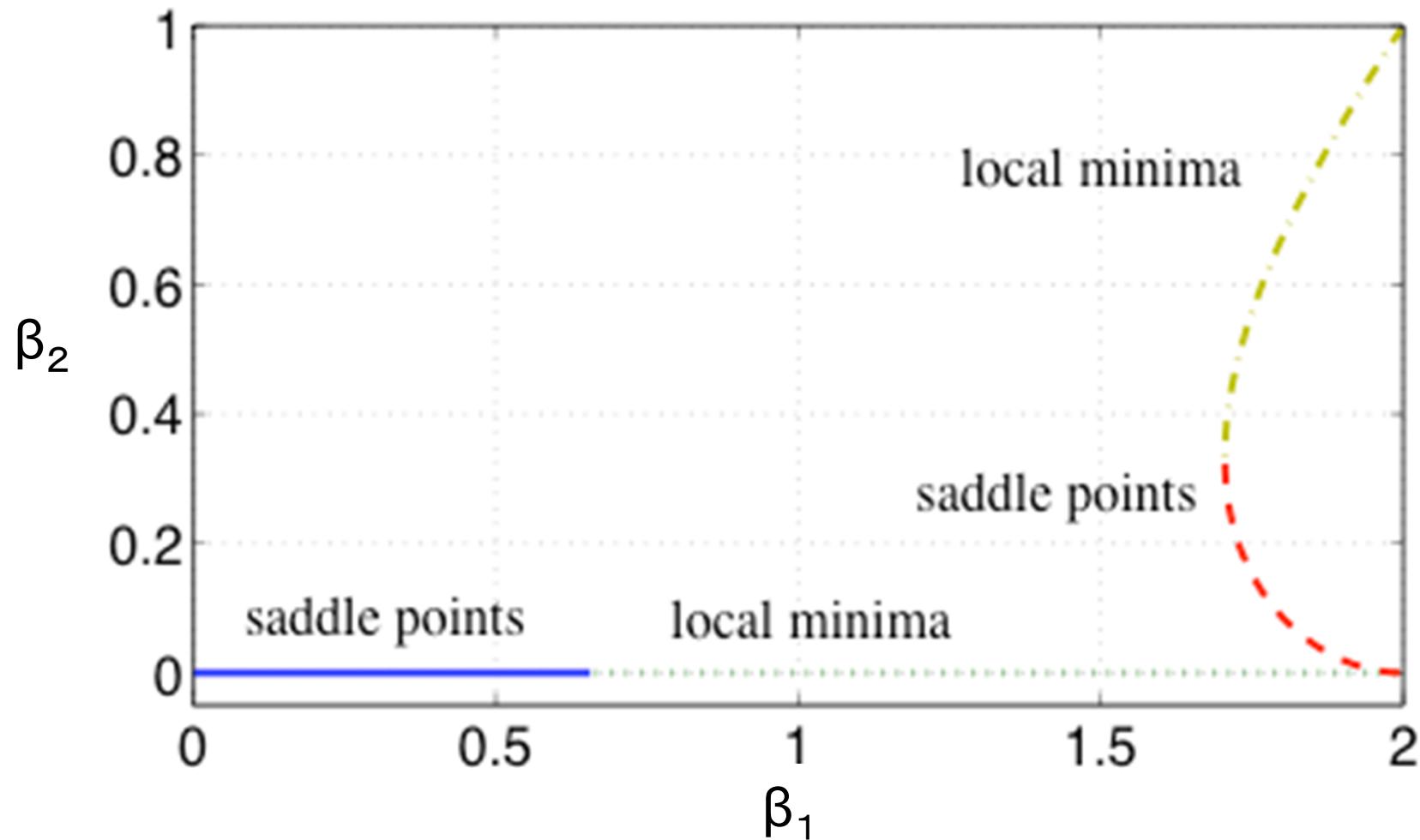
(Collaborator: Masahiro Yukawa, Niigata University)

Solution Path : $\lambda \leftrightarrow c$

**not continuous, not-monotone
jump**

$$\beta_\lambda \Leftrightarrow \beta_c$$

An Example of the greedy path

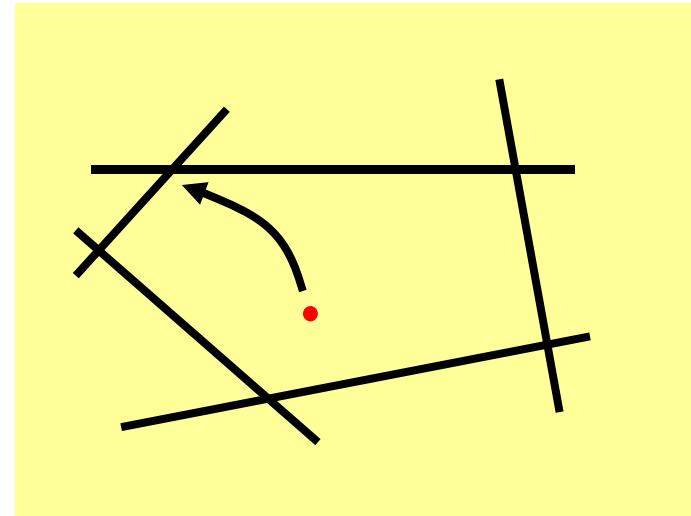


Linear Programming

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

$$\max \sum c_i x_i$$

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



Convex Programming

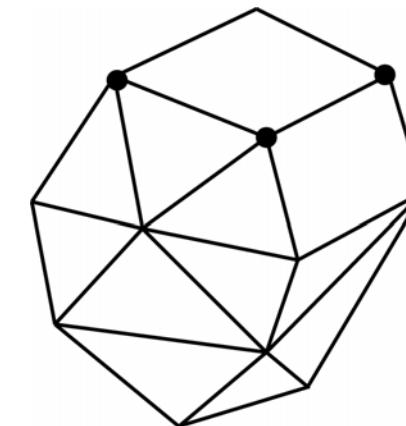
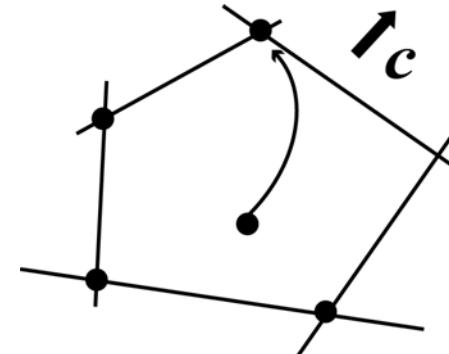
— Inner Method

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

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

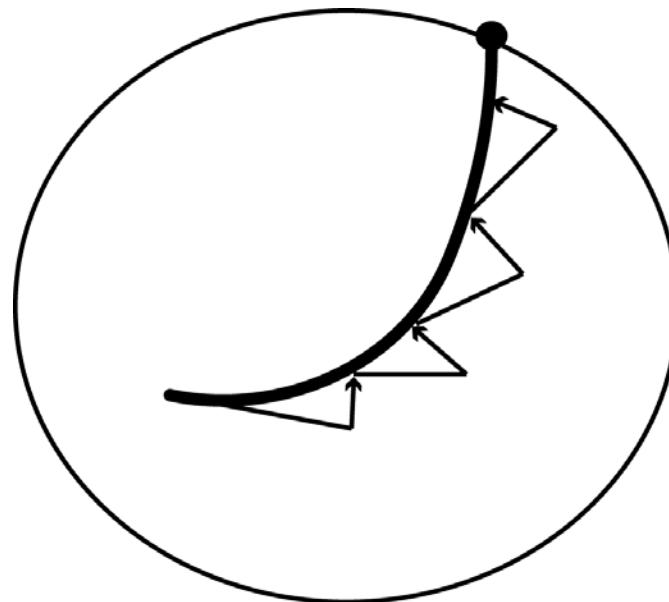
$$\begin{aligned}\psi(x) = & \sum \log \left(\sum A_{ij} x_j - b_i \right) \\ & + \sum \log x_i\end{aligned}$$

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



Simplex method ; inner method

Polynomial-Time Algorithm



curvature : step-size

$$|H^{(m)}|^2$$

$$\min : tc \cdot x + \psi(x) \quad x = \delta(t) \quad \nabla^* - \text{geodesic}$$

Neural Networks

Multilayer Perceptron

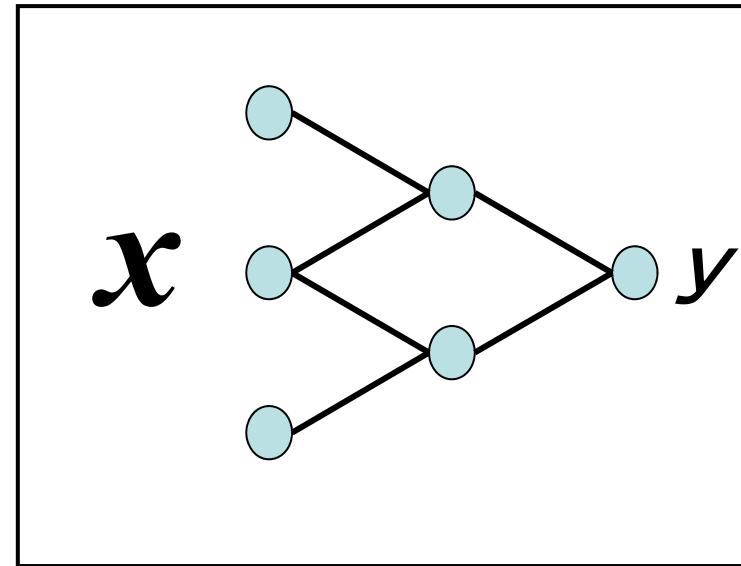
Higher-order correlations

Synchronous firing

Multilayer Perceptrons

$$y = \sum v_i \varphi(w_i \cdot x) + n$$

$$x = (x_1, x_2, \dots, x_n)$$

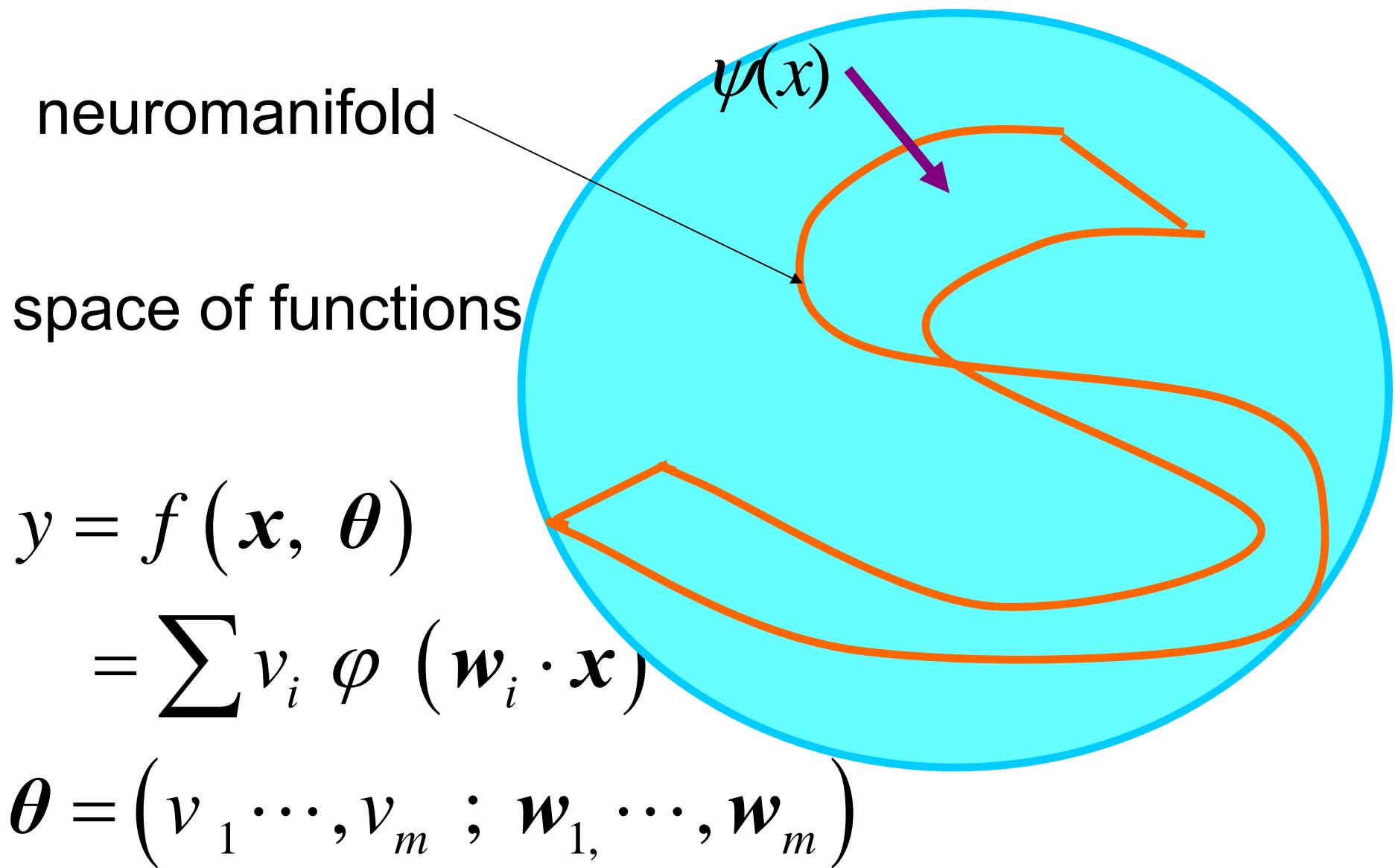


$$p(y|x;\theta) = c \exp\left\{-\frac{1}{2}(y - f(x,\theta))^2\right\}$$

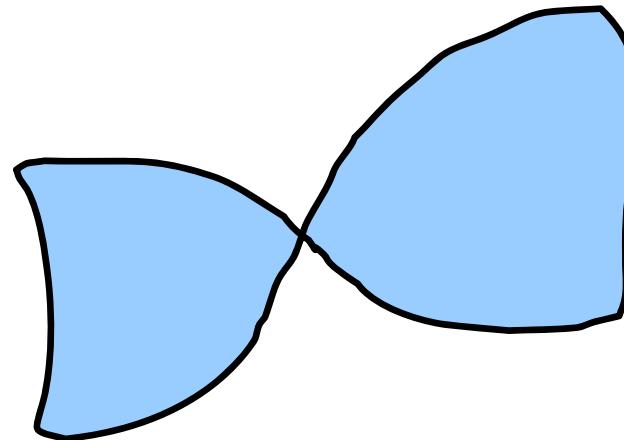
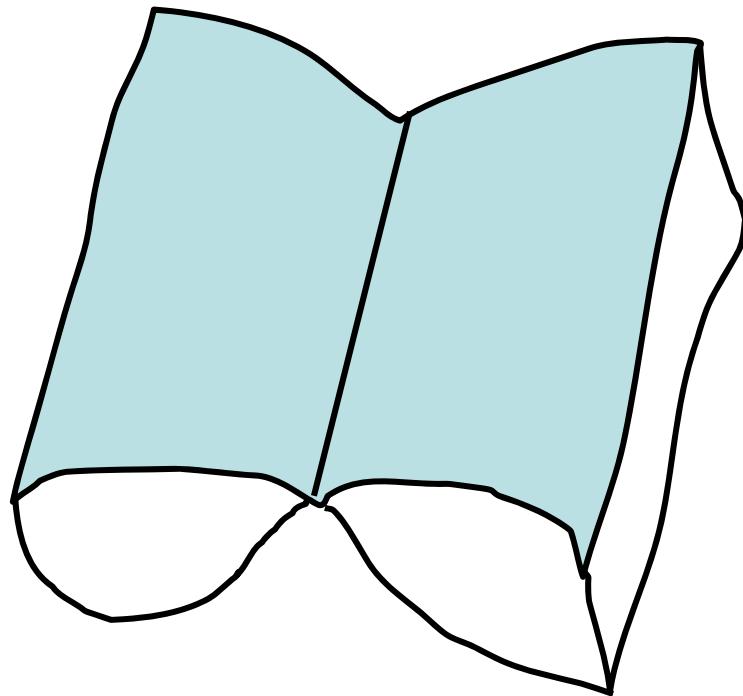
$$f(x,\theta) = \sum v_i \varphi(w_i \cdot x)$$

$$\theta = (w_1, \dots, w_m; v_1, \dots, v_m)$$

Multilayer Perceptron

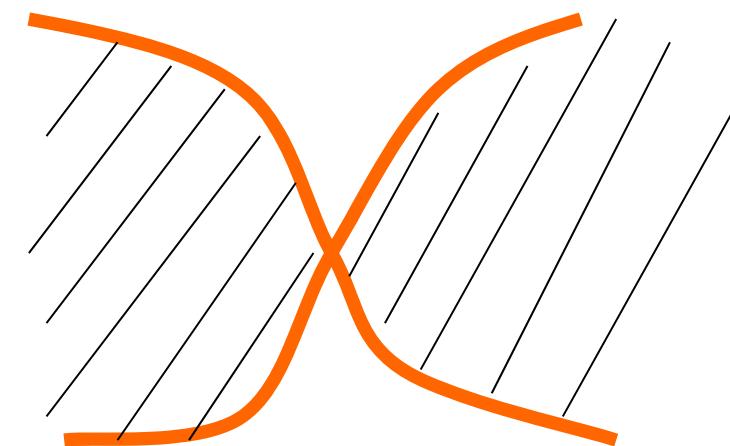
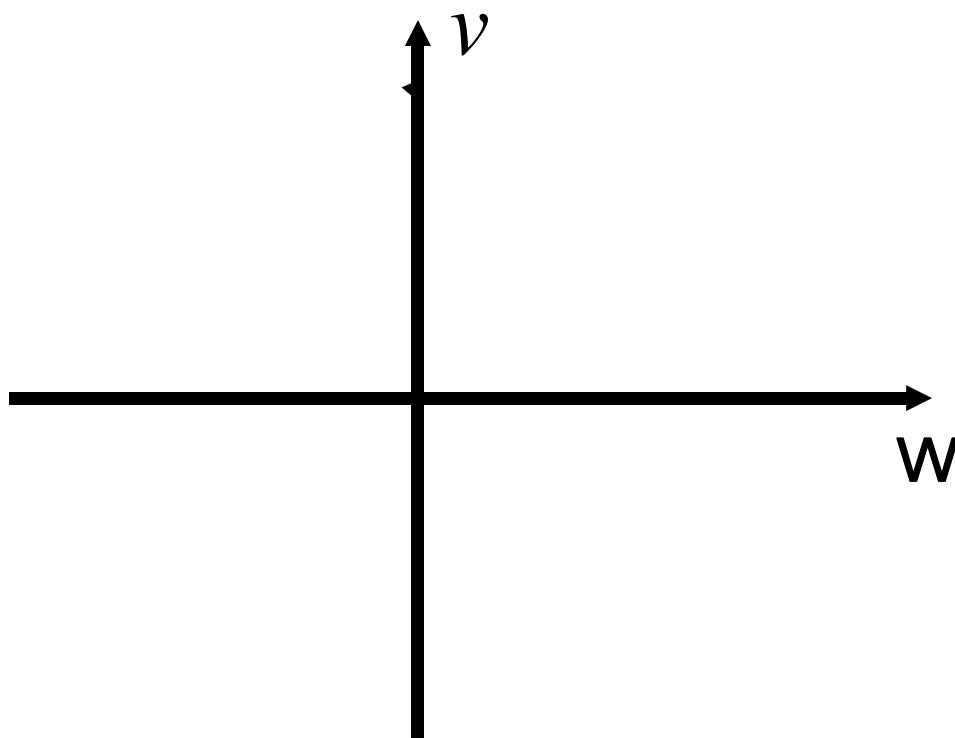


singularities



Geometry of singular model

$$y = v\varphi(w \cdot x) + n$$
$$v | w | = 0$$



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)

$$\Delta \boldsymbol{\theta}_t = -\eta_t \frac{\partial E}{\partial \boldsymbol{\theta}} \quad \tilde{\nabla} E = G^{-1} \nabla E \text{ --steepest descent}$$

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

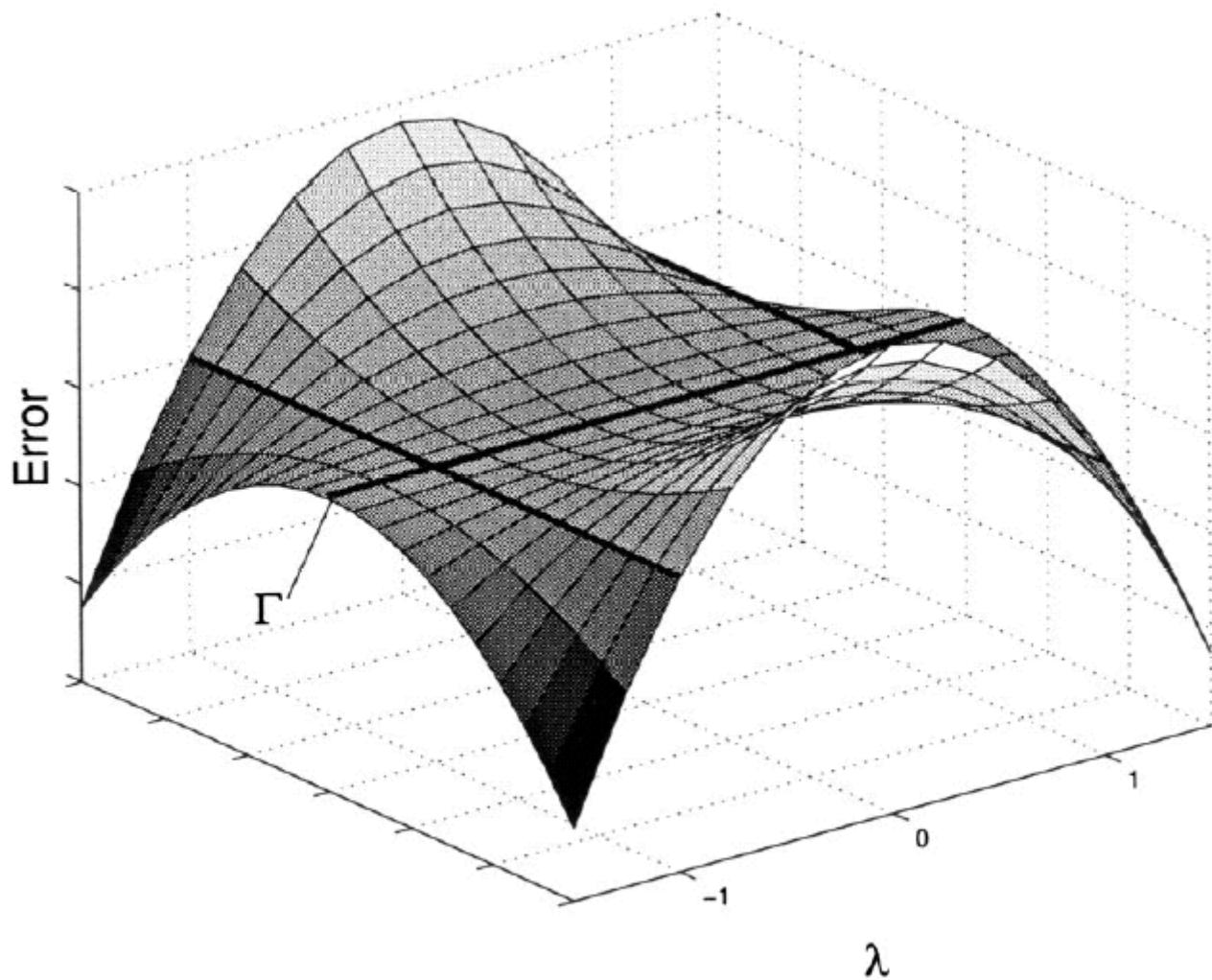
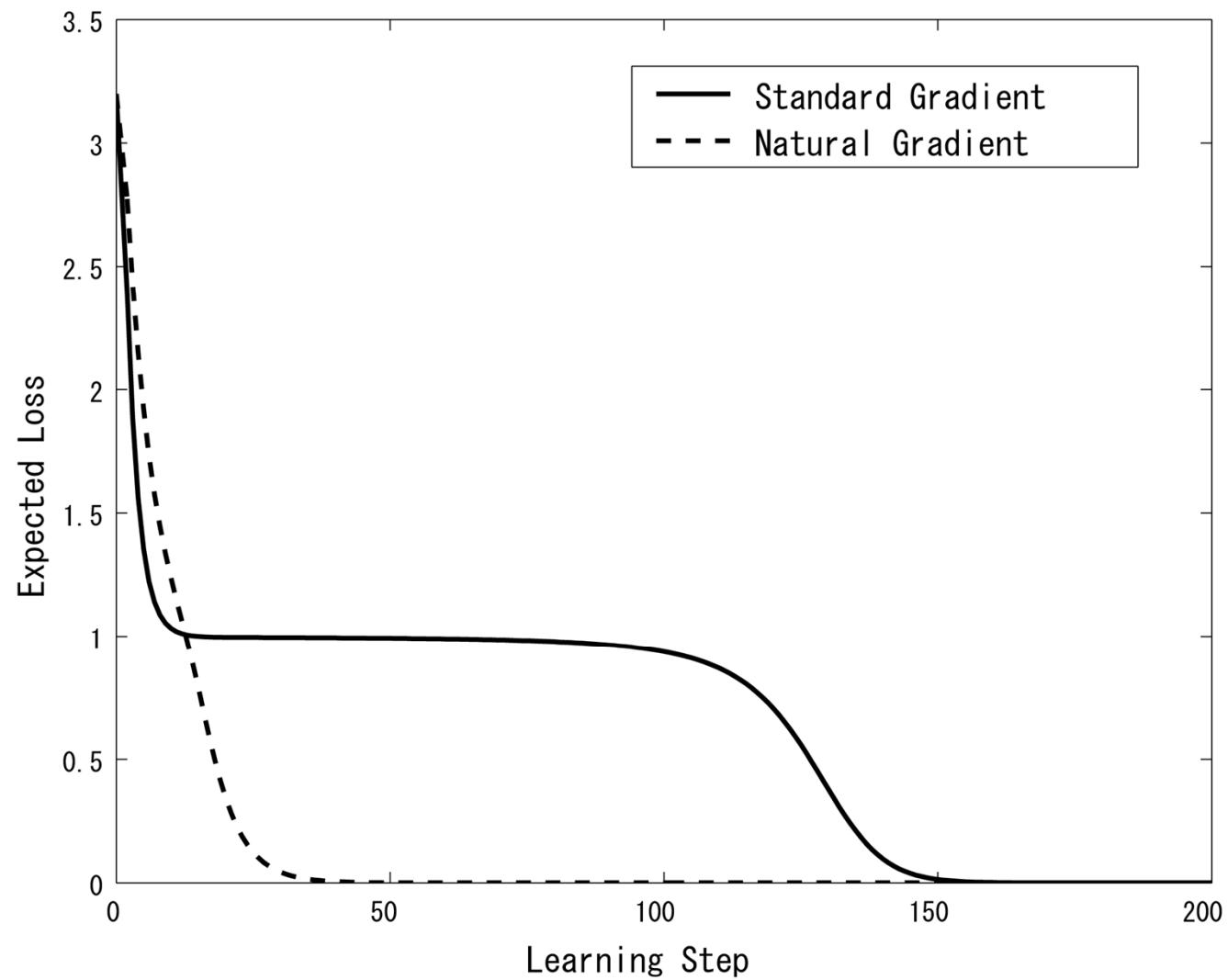


Fig. 5. Critical set with local minima and plateaus.



conformal transformation

q-Fisher information

$$g_{ij}^{(q)}(p) = \frac{q}{h_q(p)} g_{ij}^F(p)$$

q-divergence

$$D_q[p(x):r(x)] = \frac{1}{(1-q)h_q(p)} \left(1 - \int p(x)^q r(x)^{1-q} dx \right)$$

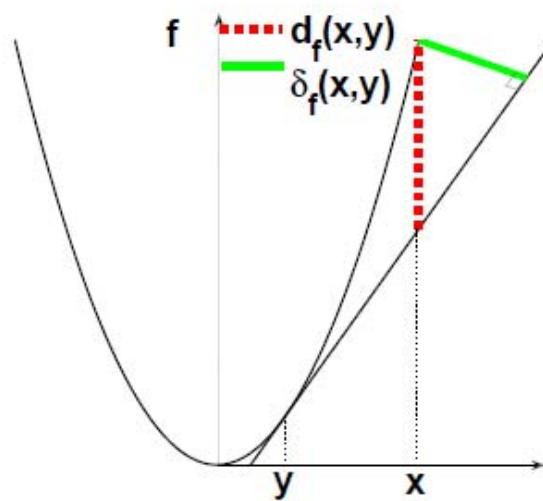
Total Bregman Divergence and its Applications to Shape Retrieval

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

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

Total Bregman Divergence

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

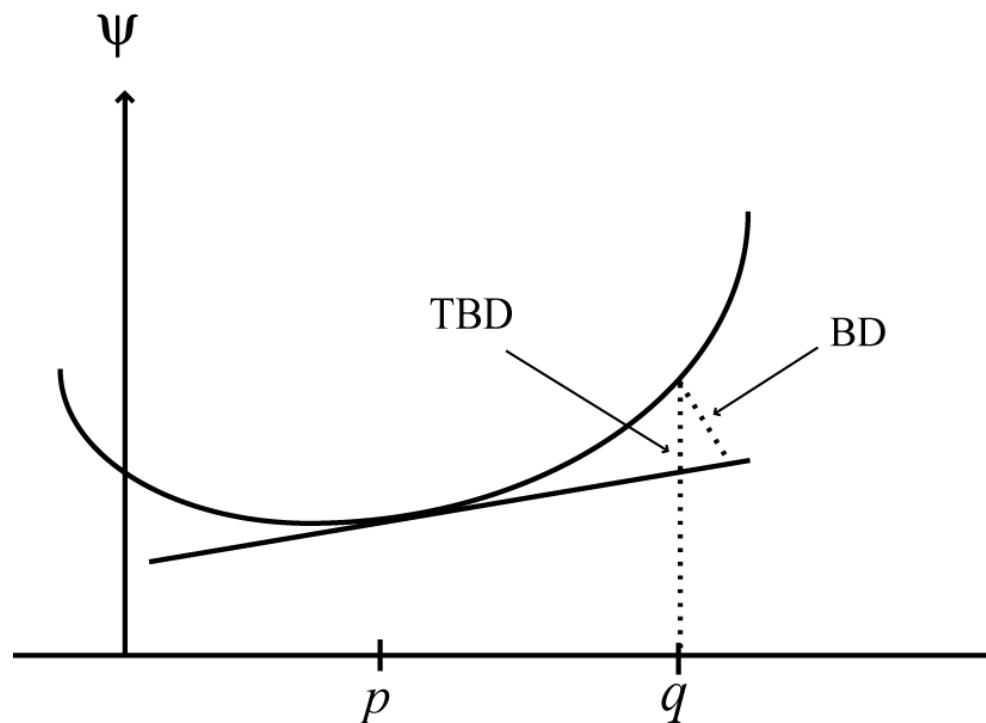


- rotational invariance
- conformal geometry

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.

Total Bregman divergence (Vemuri)

$$\text{TBD}(p : q) = \frac{\varphi(p) - \varphi(q) - \nabla \varphi(q) \cdot (p - q)}{\sqrt{1 + |\nabla \varphi(q)|^2}}$$

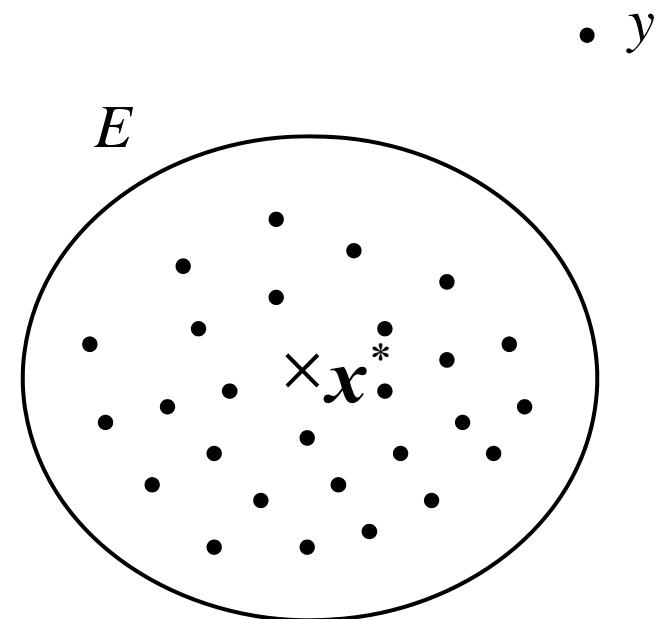


Clustering : t -center

$$E = \{x_1, \dots, x_m\}$$

T-center of
E

$$\mathbf{x}^* = \arg \min \sum_i TD[\mathbf{x}, \mathbf{x}_i]$$



$$\textbf{\textit{t-center}}\boldsymbol{x}^*$$

$$\nabla f\left(\boldsymbol{x}^*\right)\!=\!\frac{\sum w_i \nabla f\left(\boldsymbol{x}_i\right)}{\sum w_i}$$

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

q -super-robust estimator (Eguchi)

$$\max \hat{p}(x, \xi) \rightarrow \max \frac{p(x, \xi)}{h_{q+1}}$$

bias-corrected q -estimating function

$$s_q(x, \xi) = \hat{p}(x, \xi) \left\{ \partial_i \log p - c_{q+1}(\xi) \right\}$$

$$c_{q+1} = \frac{1}{q+1} \partial \log h_{q+1}(\xi)$$

$$\sum_{i=0}^N s_q(x_i, \xi) = 0 \quad \Leftrightarrow \max \frac{1}{h_{q+1}} \sum \hat{p}(x_i, \xi)$$

Conformal change of divergence

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

$$\tilde{g}_{ij} = \sigma(p) g_{ij}$$

$$\tilde{T}_{ijk} = \sigma(T_{ijk} + s_k g_{ij} + s_j g_{ik} + s_i g_{jk})$$

$$s_i = \partial_i \log \sigma$$

t-center is robust

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

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

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

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

$$z(x^*, y) = G^{-1} \frac{\nabla f(y) - \nabla f(x^*)}{w(y)}$$

$$G = \frac{1}{n} \sum w(x_i) \nabla \nabla f(x_i)$$

Robust: z is bounded

$$\frac{\nabla f(y)}{w(y)} = \frac{\nabla f(y)}{\sqrt{1 + \|\nabla f(y)\|^2}} < \infty$$

$$z(x^*, y) = G^{-1} \frac{y}{\sqrt{1 + \|y\|^2}}$$

$z(x^*, y) = y$ Euclidean case $f = \frac{1}{2} |x|^2$

$$w(y) > 1$$

MPEG7 database

- Great intraclass variability, and small interclass dissimilarity.



Other TBD applications

Diffusion tensor imaging (DTI) analysis [Vemuri]

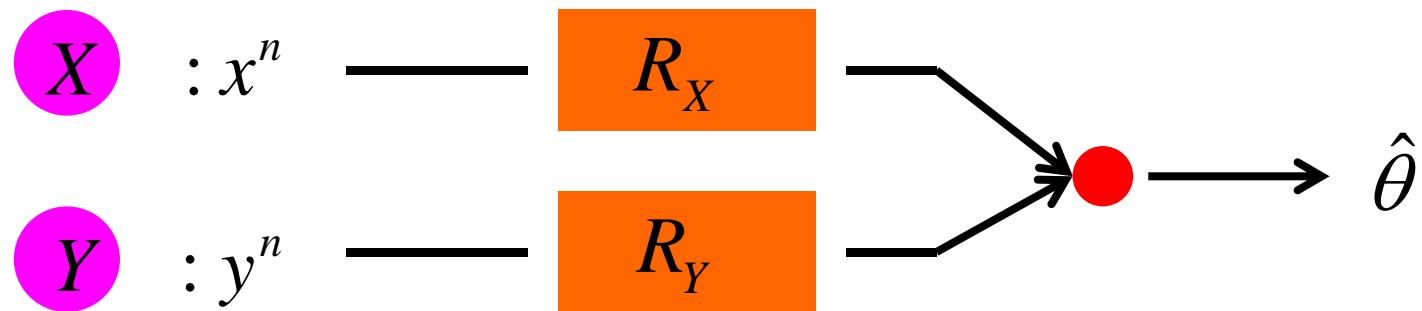
- **Interpolation**
- **Segmentation**

Baba C. Vemuri, Meizhu Liu, Shun-ichi Amari and Frank Nielsen, *Total Bregman Divergence and its Applications to DTI Analysis*, IEEE TMI, to appear

TBD application-shape retrieval

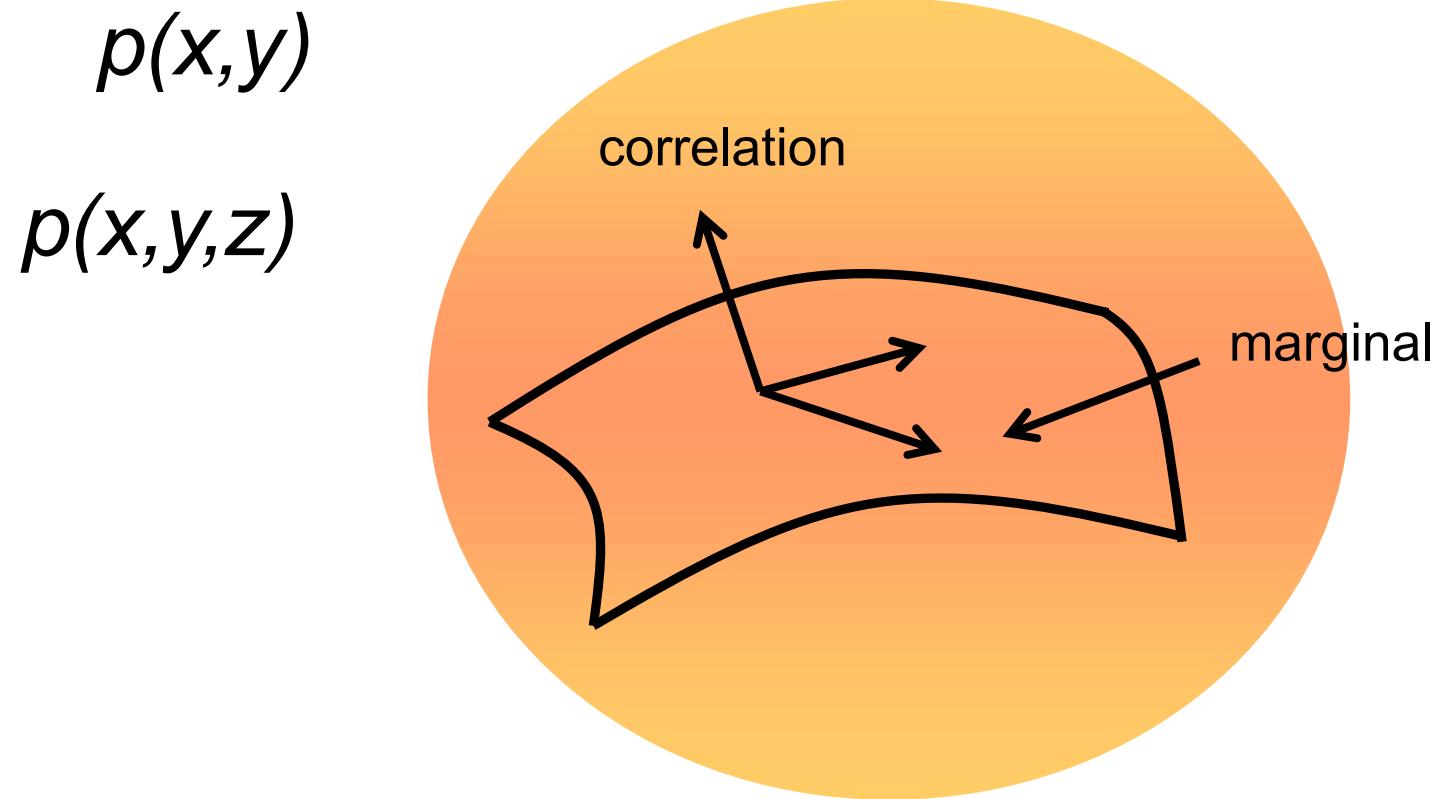
- Using MPEG7 database;
- 70 classes, with 20 shapes each class
(Meizhu Liu)

Multiterminal Information & Statistical Inference



$$p(x, y; \theta)$$

$$|M_X| = 2^{R_X n} \quad |M_r| = 2^{R_Y n}$$



$$G = G_M + G_C$$

0-rate Slepian-Wolf

Linear Systems



ARMA

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

AR---e-flat

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

MA---m-flat

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

Machine Learning

Boosting : combination of weak learners

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

$$y_i = \pm -1$$

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

Weak Learners

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

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

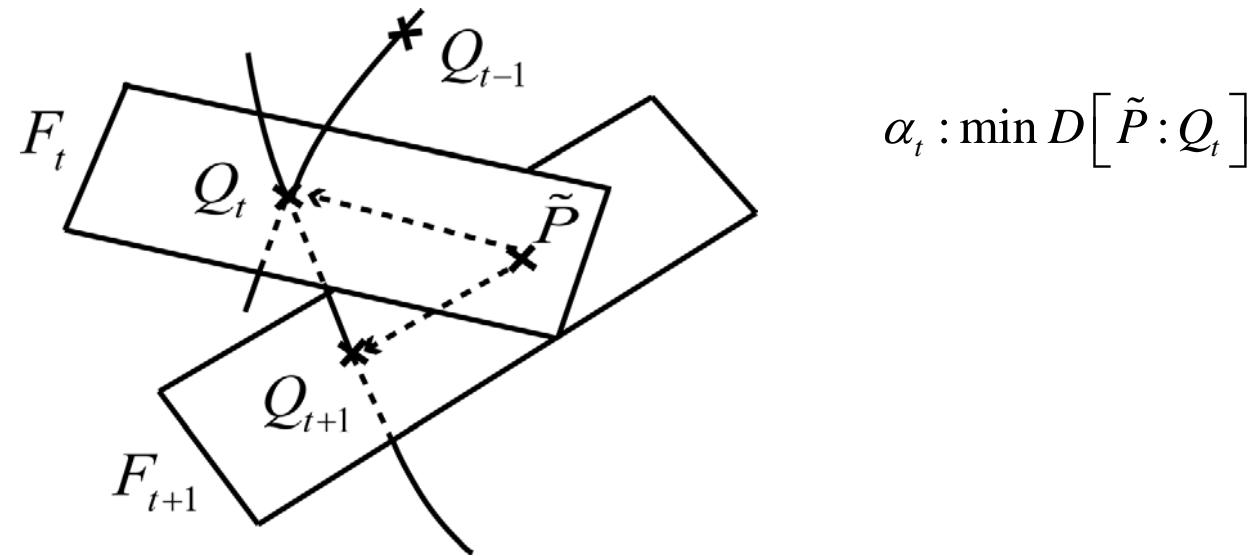
$$W_{t+1}(i) = c W_t(i) \exp\{-\alpha_t y_i h_t(x_i)\}$$

weight distribution

Boosting —generalization

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

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



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

Integration of evidences:

$$x_1, x_2, \dots x_m$$

arithmetic mean

geometric mean

harmonic mean

α -mean

Various Means

$$\frac{a+b}{2} : \sqrt{ab} : \frac{2}{\frac{1}{a} + \frac{1}{b}}$$

arithmetic geometric harmonic

Any other mean?

Generalized mean: f-mean

f(u): monotone; f-representation of u

$$m_f(a,b) = f^{-1}\left\{\frac{f(a)+f(b)}{2}\right\}$$

scale free

$$m_f(ca,cb) = cm_f(a,b)$$

α -representation

$$f_\alpha(u) = u^{\frac{1-\alpha}{2}}, \quad \alpha \neq 1$$

$$\log u, \quad \alpha=1$$

α -mean : $m_\alpha(p_1(s), p_2(s))$

$$\alpha = 1 : \sqrt{ab}$$

$$\alpha = -1 : \frac{a+b}{2}$$

$$\alpha = 0 : (\sqrt{a} + \sqrt{b})^2 = \frac{a+b}{4} + \frac{1}{2}\sqrt{ab}$$

$$\alpha = \infty \qquad m_\alpha = \min(a, b)$$

$$\alpha = -\infty \qquad m_\alpha = \max(a, b)$$

α – Family of Distributions

$$\{p_1(s), \dots, p_k(s)\} \quad p(x; \theta) = f_\alpha^{-1}\left\{\sum \theta_i f_\alpha(p_i(x))\right\}$$

mixture family :

$$p_{mix}(s) = \sum_{i=1}^k t_i p_i(s), \quad \sum t_i = 1$$

exponential family :

$$\log p_{\text{exp}}(s) = \sum t_i \log p_i(s) - \psi$$

