

# **Applications of Information Geometry**

**Clustering based on Divergence: Cluster Center**

**Stochastic Reasoning: Belief Propagation in Graphical Model**

**Support Vector Machine**

**Bayesian Framework and Restricted Boltzmann Machine**

**Natural Gradient in Multilayer Perceptron Learning**

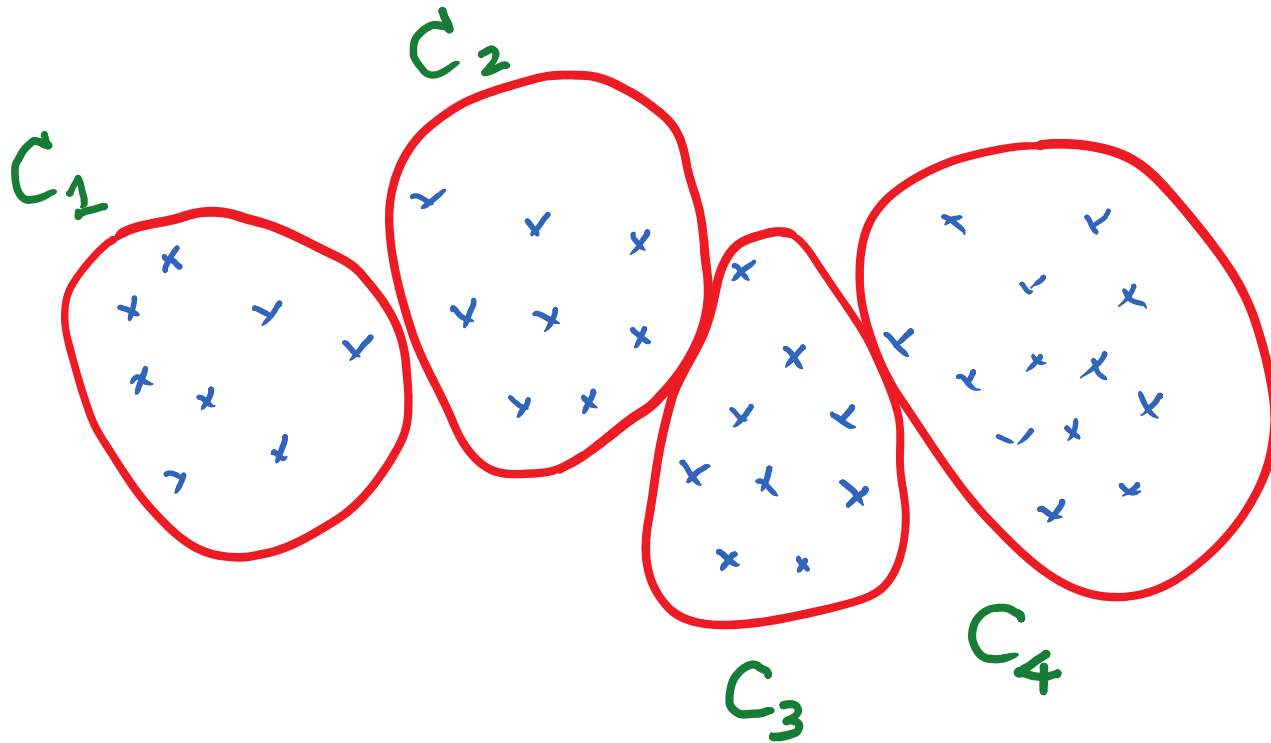
**Independent Component Analysis**

**Sparse Signal Analysis and Minkovskian Gradient**

**Convex Optimization**

# Clustering of Points

$$D = \{\theta_1, \theta_2, \theta_3, \dots, \theta_N\}$$

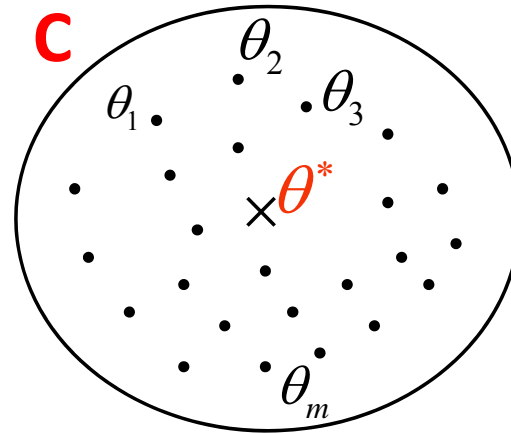


# Clustering : center of a cluster of points

Center of C

$$C = \{\theta_1, \dots, \theta_m\}$$

$$\theta^* = \arg \min_{\theta} \sum_i D[\theta, \theta_i]$$



$\eta^*$  is simply the arithmetic average in the dual coordinate system

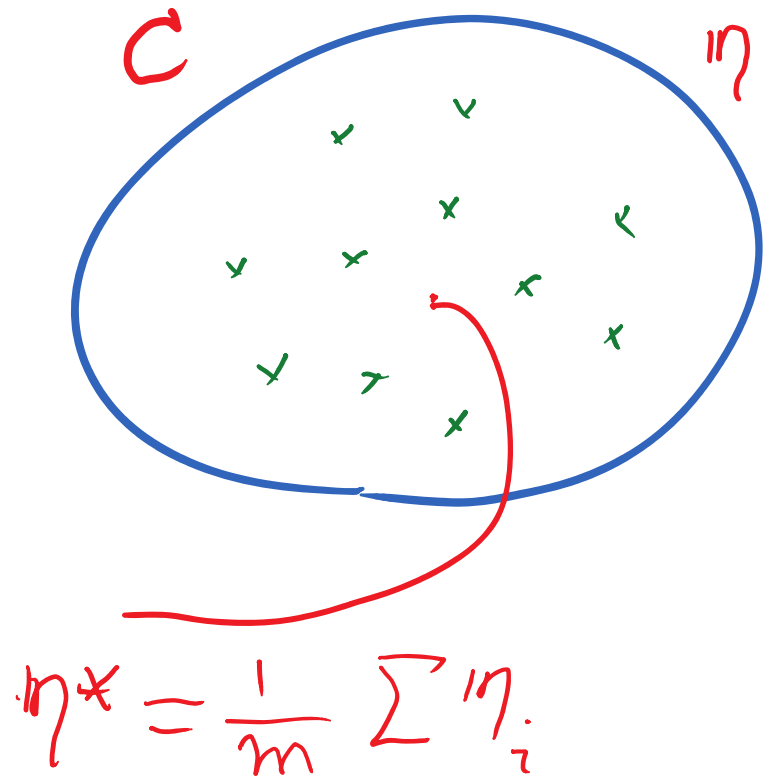
$$D[\theta : \theta_i] = \psi(\theta) + \varphi(\eta_i) - \theta \cdot \eta_i$$

$$\partial_{\theta} D[\theta : \theta_i] = \eta - \eta_i$$

$$\sum (\eta - \eta_i) = 0$$

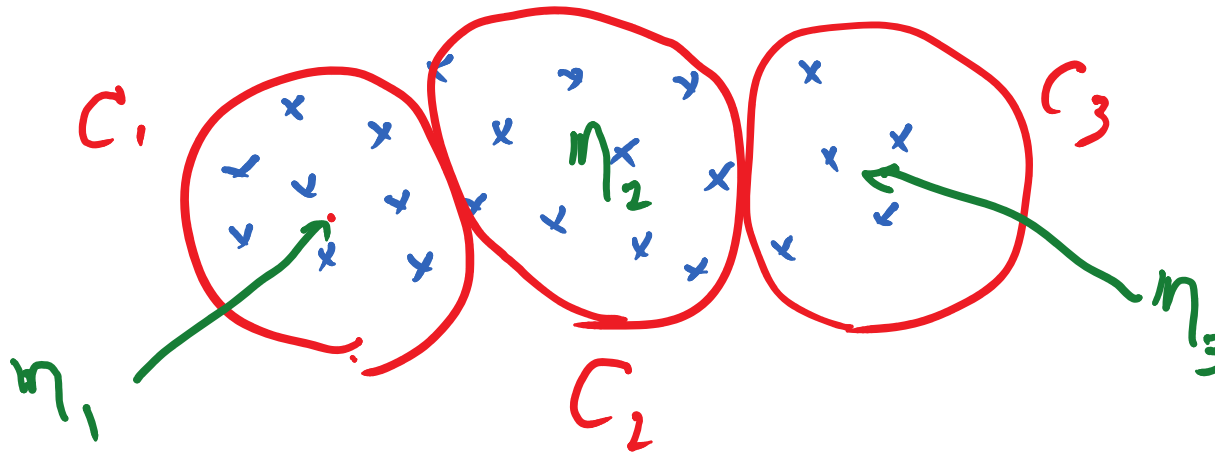
$$\eta^* = \frac{1}{m} \sum \eta_i$$

k-means: clustering algorithm



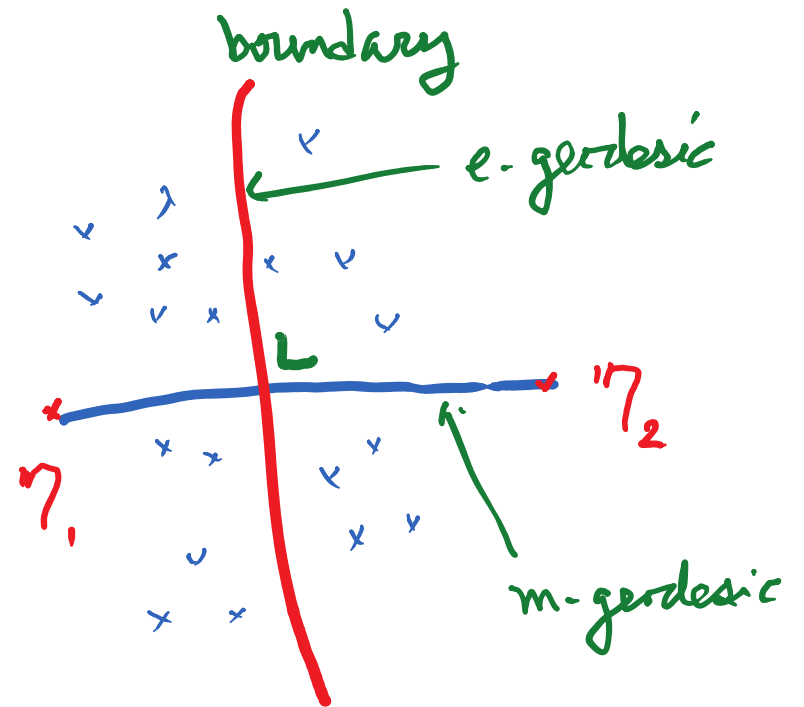
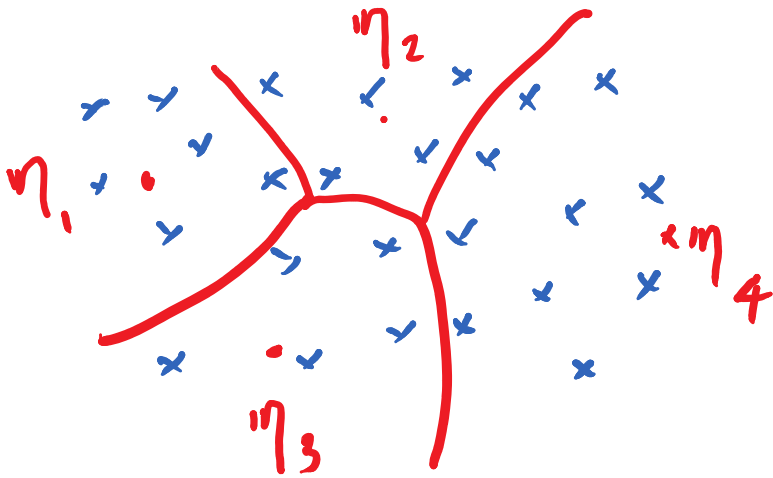
# k-means: clustering algorithm

1. Choose  $k$  cluster centers
2. Classify patterns of  $D$  into clusters
3. Calculate new cluster centers



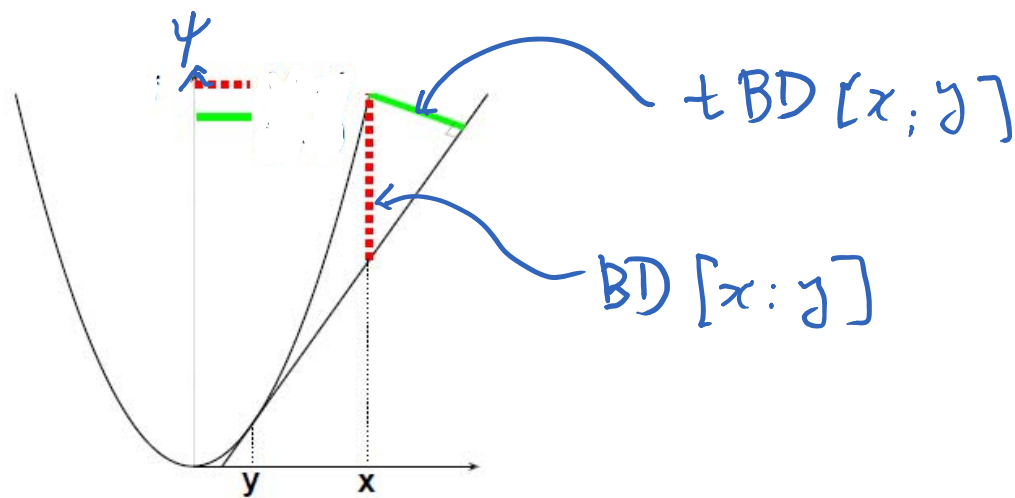
# Voronoi Diagram: Boundaries of Clusters

$$D[\eta : \eta_1^*] = D[\eta : \eta_2^*]$$



## Total Bregman divergence (Vemuri)

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



**$t$ -center of  $C$**   $\eta^*$

$$\eta^* = \frac{\sum w_i \eta_i}{\sum w_i}$$

$$w_i = \frac{1}{\sqrt{1 + \|\nabla \psi(\eta_i)\|^2}}$$



## Conformal change of divergence

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

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

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

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

# ***t*-center is robust**

$$E^* = \{\eta_1, \dots, \eta_n; \mathbf{y}\}$$

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

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

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

$$\text{minimize } \frac{1}{N+1} \left( \sum \frac{D[\eta : \eta_i]}{w_i} + \frac{D[\eta : \mathbf{y}]}{w_{N+1}} \right)$$

Robust:  $\mathbf{z}$  is bounded

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

$$w(\mathbf{y}) > 1$$

$$\mathbf{z}(\eta^*, \mathbf{y}) = G^{-1} \frac{\nabla \psi(\theta_y) - \nabla \psi(\eta^*)}{w(\mathbf{y})}$$

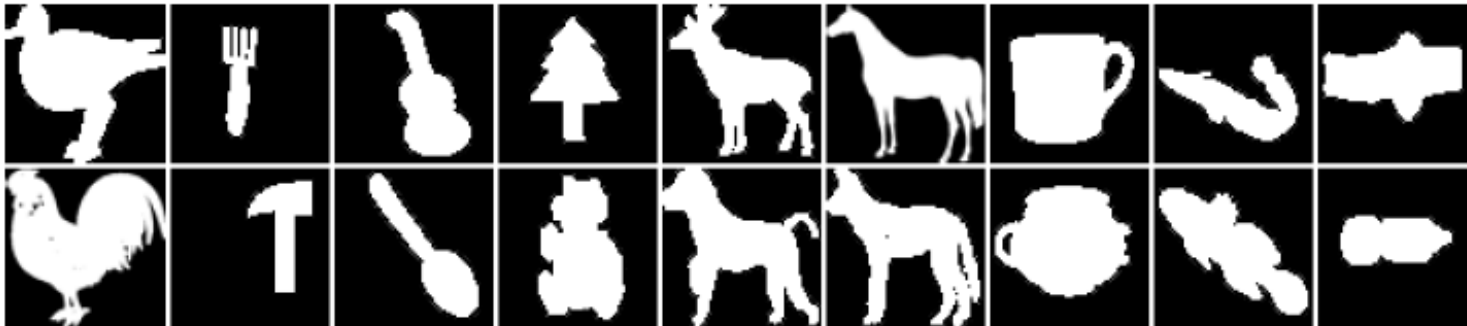
$$\mathbf{z}(\eta^*, \mathbf{y}) = G^{-1} \frac{\mathbf{y}}{\sqrt{1 + \|\mathbf{y}\|^2}}$$

$$\mathbf{z}(\eta^*, \mathbf{y}) = \mathbf{y}$$

Euclidean case  $f = \frac{1}{2}|\mathbf{x}|^2$

# MPEG7 database

- Great intraclass variability, and small interclass dissimilarity.



# Shape representation

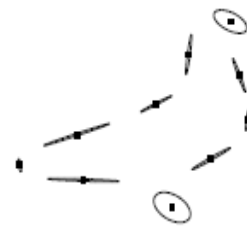
A shape is represented using a mixture of Gaussians from the aligned boundary points.



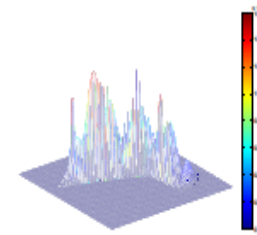
(a)



(b)



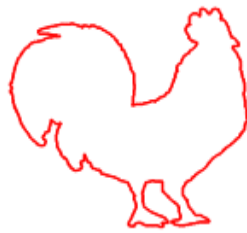
(c)



(d)



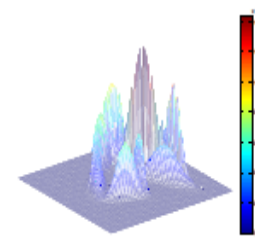
(e)



(f)



(g)



(h)

# First clustering then retrieval

## Retrieval in the whole MPEG-7 database

First clustering and then retrieval.

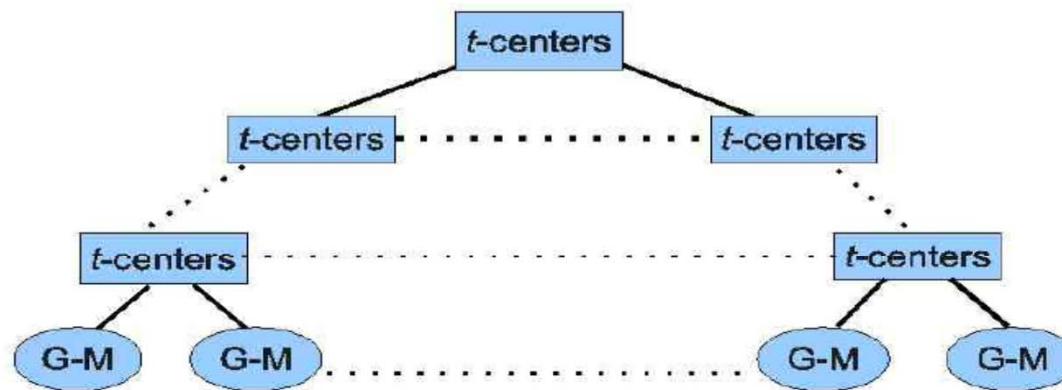


Figure: *k*-Tree diagram. G-M: mixture of Gaussians. Every key is a mixture of Gaussians. Each key in the inner nodes is the *t*-center of all keys in its children nodes. The key of a leaf is a mixture of Gaussians corresponding to an individual shape.

# 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

# Information Geometry of Stochastic Reasing

## Belief Propagation in Graphical Model

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

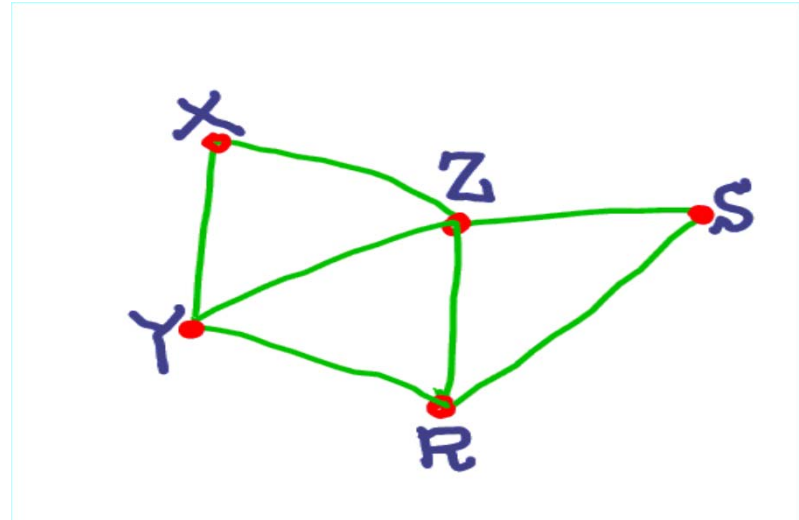


# Stochastic Reasoning: Graphical Model

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

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

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



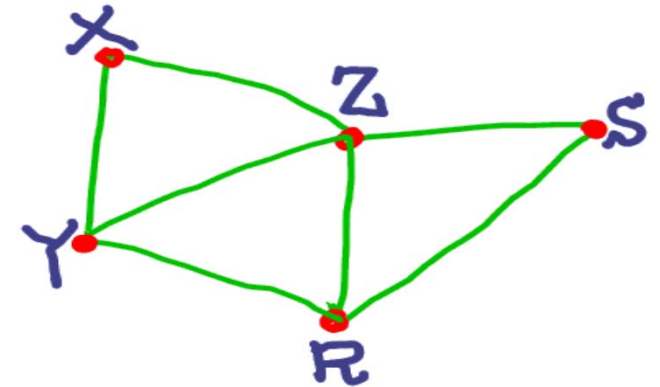
# Stochastic Reasoning

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

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

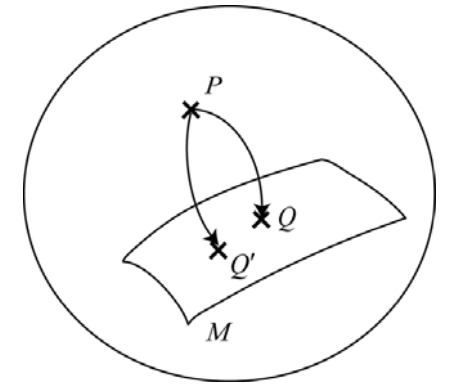
$\mathbf{X} = \text{argmax } q(x_1, x_2, x_3, \dots)$     maximum likelihood

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



# Mean Value

Marginalization:  
projection to independent distributions

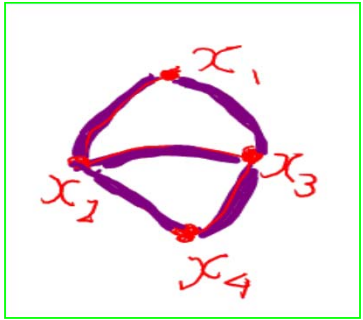


$$\Pi_0 q(\mathbf{x}) = q_1(x_1)q_2(x_2)\dots q_n(x_n) = q_0(\mathbf{x})$$

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

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

# cliques



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

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

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

**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

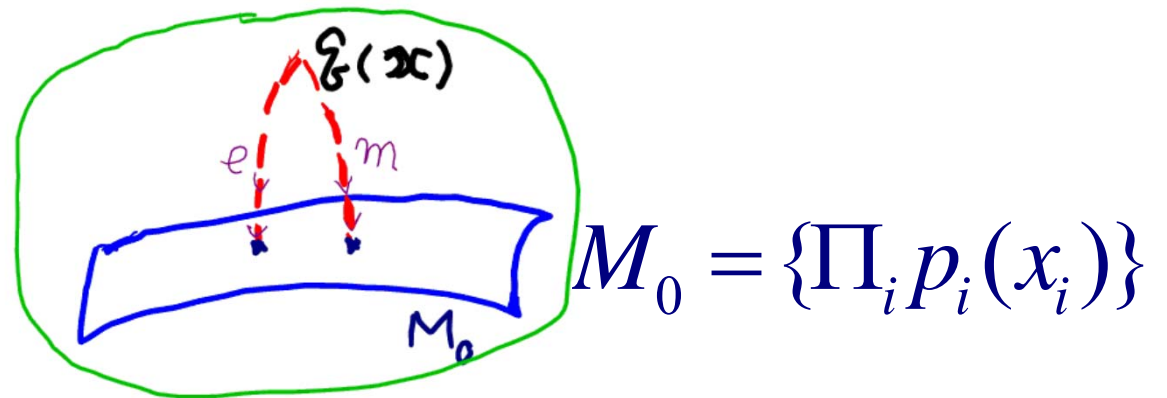
- m-projection
- e-projection

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

$$\Pi_0^m q = \operatorname{argmin}_{q \in M_0} D[q:p]$$

$$\Pi_0^e q = \operatorname{argmin}_{q \in M_0} D[p:q]$$

$$p(x) \in M_0$$



## m-projection keeps the expectation of x

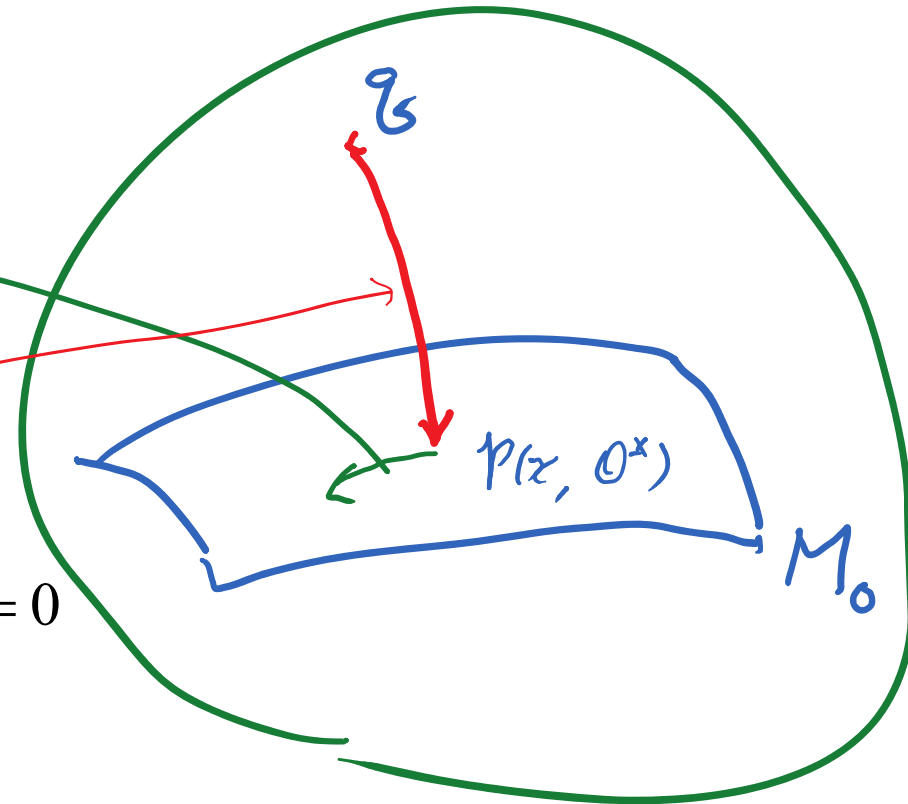
$$p(x, \theta^*) = \prod_{q \in M_0}$$

$$\frac{\partial}{\partial \theta} p(x, \theta^*) = x - \eta^*$$

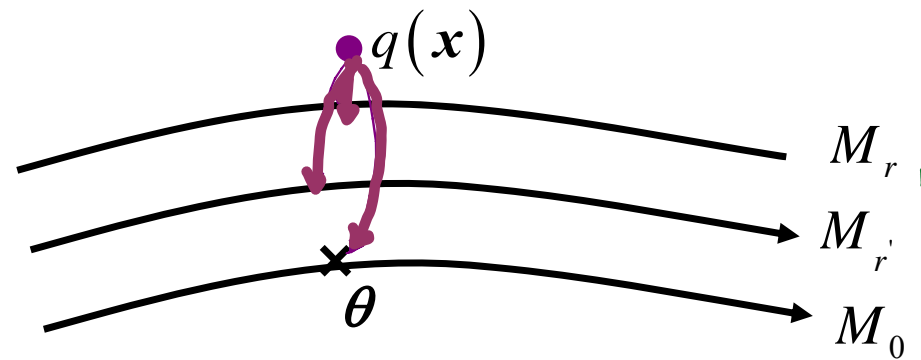
$$t(x) = \frac{q(x) - p^*(x)}{p^*(x)}$$

$$\langle t(x), x - \eta^* \rangle_{p^*} = \sum_x (x - \eta^*) q(x) - p^*(x) = 0$$

$$E_q[x] = E_{p^*}[x]$$



# Information Geometry



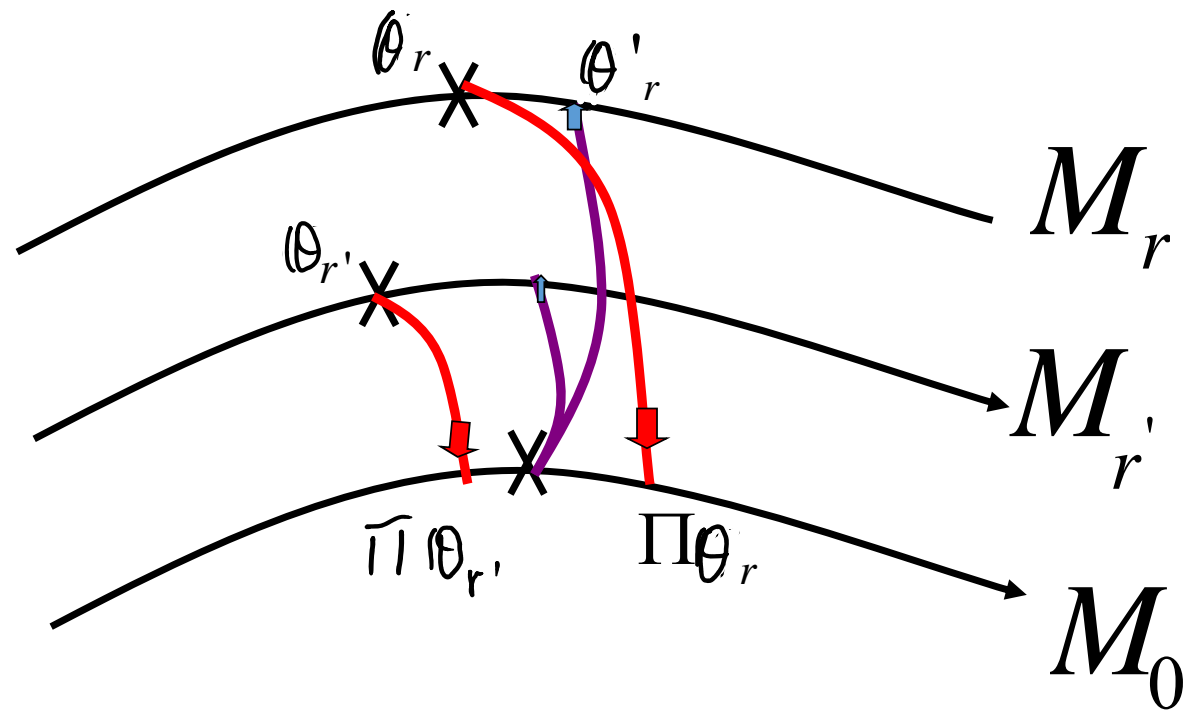
$$q(x) = \exp \left\{ \sum c_r(x) - \phi \right\}$$

$$M_0 = \left\{ p_0(\mathbf{x}, \boldsymbol{\theta}_0) \right\} = \exp \left\{ \boldsymbol{\theta}_0 \cdot \mathbf{x} - \psi_0 \right\}$$

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



# Belief Prop Algorithm



# Belief Propagation

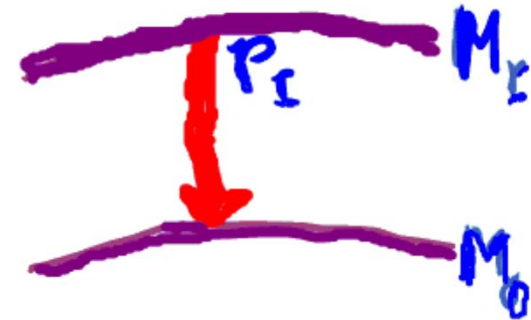
$$p(x, \theta_r) = \exp\{c_r(x) + \theta_r \cdot x - \psi_r\}$$

$$\Pi_0 p_r(x, \theta_r^t) \Rightarrow p_r(x, \tilde{\theta}_0^t)$$

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

$$\theta_r^{t+1} = \sum_{r' \neq r} \xi_{r'}^{t+1}$$

$$\theta_0^{t+1} = \sum \xi_r^{t+1}$$



# Equilibrium of BP

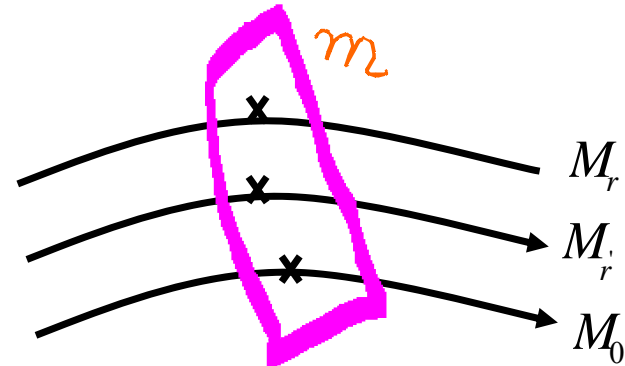
$$\left( \theta^*, \xi_r^* \right)$$

1)  $m$ -condition

$$\theta^* = \Pi_0 p_r \left( \mathbf{x}, \theta_r^* \right)$$

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

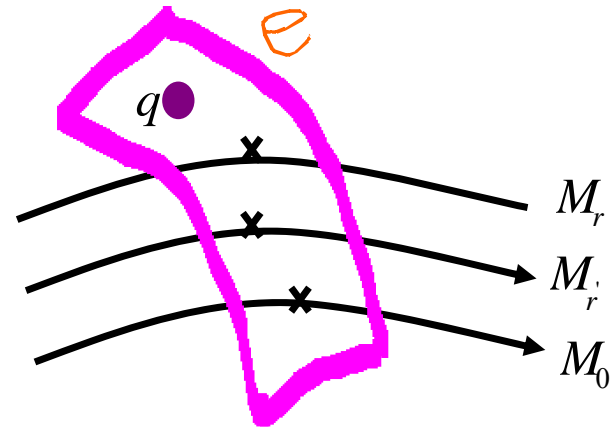
$\xi_r(\theta^*)$



2)  $e$ -condition

$$\theta^* = \sum \xi_r^*$$

$q(\mathbf{x}) \in e$ -flat submanifold



## Belief Propagation e-condition OK

$$(\boldsymbol{\theta}; \xi_1, \xi_2, \dots, \xi_L), \quad \boldsymbol{\theta}' = \sum \xi'_r$$
$$(\theta_1, \theta_2, \dots, \theta_L) \rightarrow (\theta'_1, \theta'_2, \dots, \theta'_L)$$

## CCCP m-condition OK

$$\boldsymbol{\theta}_0 \rightarrow \boldsymbol{\theta}_r : \prod p_0(x, \theta_r) = p_r(x, \theta_r)$$

$$\boldsymbol{\theta}'_0 = \sum \xi_r$$

## Convex-Concave Computational Procedure (CCCP) : A. Yuille

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

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

**Elimination of double loops**

# Geometry of Support Vector Machine

## modification of kernel

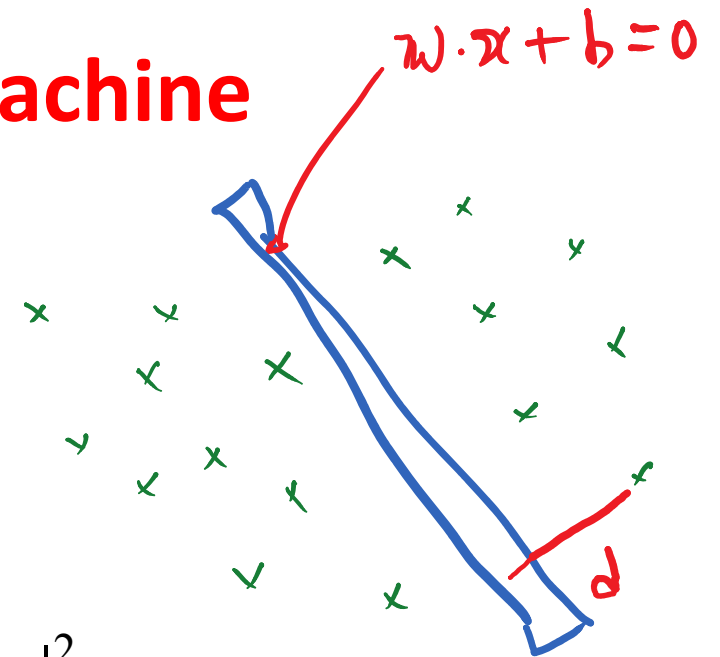
Simple perceptron: decision surface

$$f(x) = w \cdot x + b$$

$$d = \frac{|w \cdot x + b|}{|w|}$$

minimize  $|w|^2$

constraint  $y_i (w \cdot x_i + b) \geq 1$



$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum \alpha_i (w \cdot x_i + b)$$

$$\sum \alpha_i y_i = 0: \quad \text{support vector } x_i : \alpha_i \neq 0$$

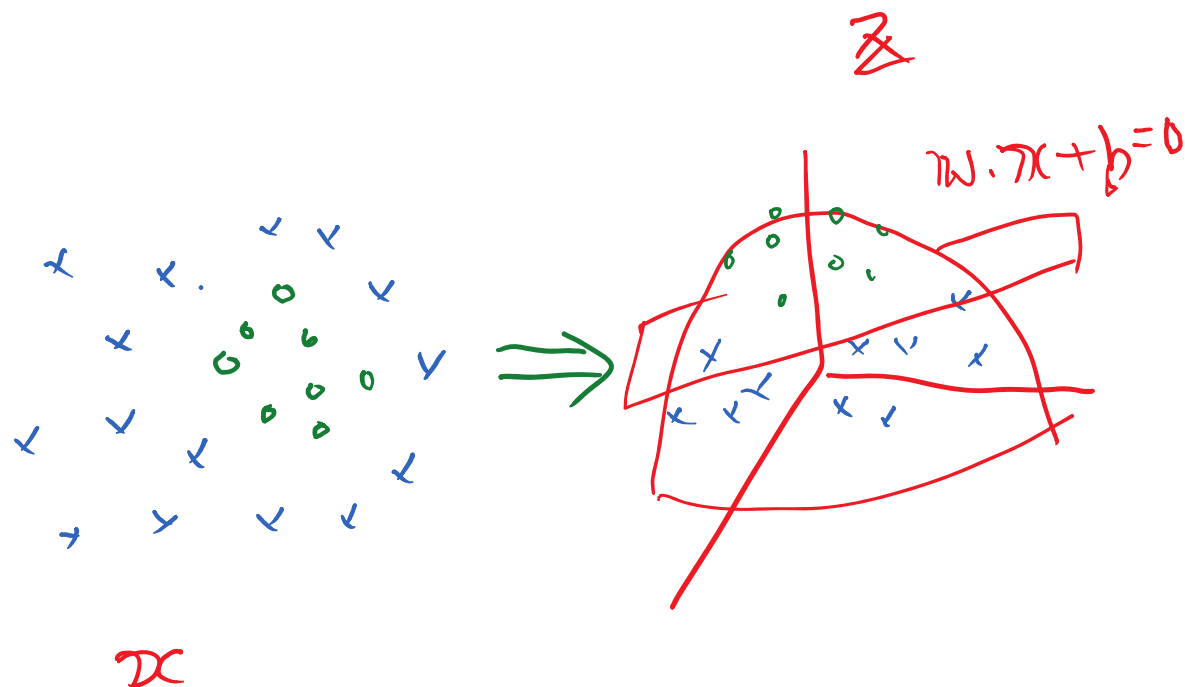
$$w = \sum \alpha_i y_i x_i \quad f(x, w) = \sum \alpha_i y_i x_i + b$$

# Embedding in higher-dimensional space

$$x \Rightarrow z = \varphi(x)$$

$$f(x) = w \cdot \varphi(x) + b$$

**z: infinite-dimensional**





# Kernel trick

$$K(x, x') = \varphi(x) \cdot \varphi(x')$$

$$\int K(x, x') \tilde{\varphi}_i(x') dx' = \lambda_i \tilde{\varphi}_i(x)$$

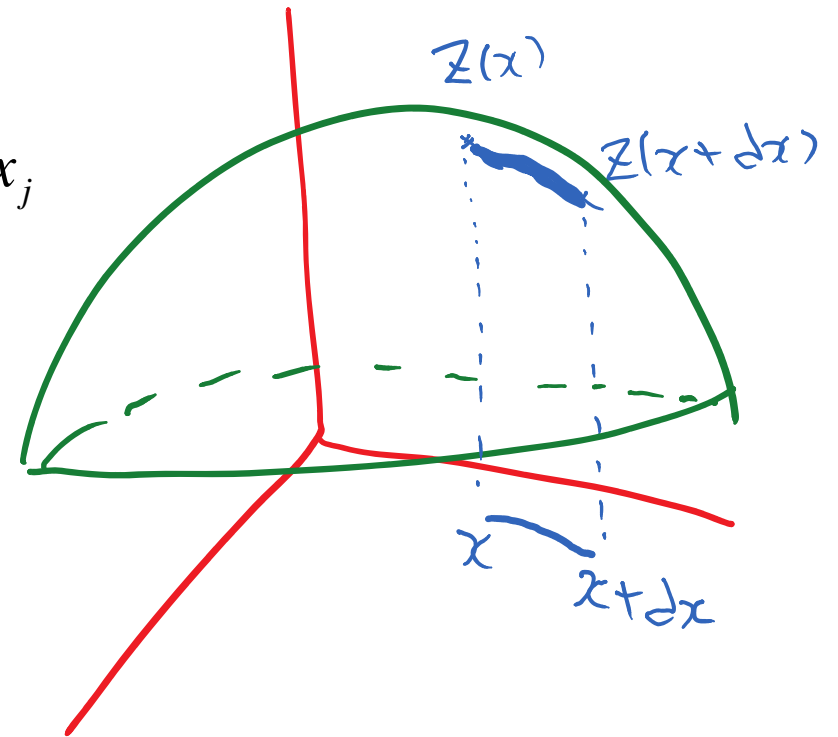
$$\varphi(x) = \frac{1}{\sqrt{\lambda_i}} \tilde{\varphi}_i(x)$$

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

$$ds^2 = |z(x+dx) - z(x)|^2 = \sum \frac{\partial}{\partial x_i} z(x) \frac{\partial}{\partial x_j} z(x) dx_i dx_j$$

$$g_{ij}(x) = \sum \frac{\partial}{\partial x_i} \varphi(x) \frac{\partial}{\partial x_j} \varphi(x)$$

$$g_{ij}(x) = \frac{\partial^2}{\partial x_i \partial x_j} K(x, x') \Big|_{x'=x}$$



# Conformal Transformation of Kernel

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

$$\tilde{K}(x, x') = \sigma(x)\sigma(x')K(x, x')$$

$$\sigma(x) = \sum \exp\{-\kappa_i |x - x_i^*|^2\}$$

$$g_{ij}(x) = \{\sigma(x)\}^2 g_{ij}(x) + \partial_i \sigma(x) \partial_j \sigma(x)$$

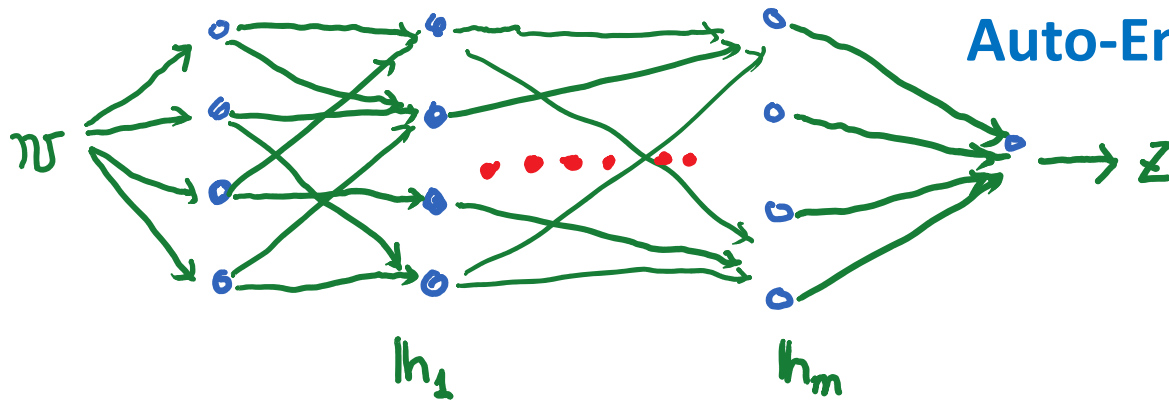
# Basic Principles of Unsupervised and Supervised Learning Toward Deep Learning

Shun-ichi Amari (RIKEN Brain Science Institute)  
collaborators: R. Karakida, M. Okada (U. Tokyo)

# Deep Learning

Self-Organization + Supervised Learning

RBM: Restricted Boltzmann Machine  
Auto-Encoder, Recurrent Net



tricks!!  
ideas!

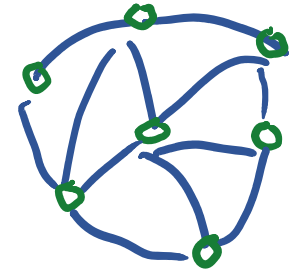
Dropout  
Contrastive divergence

bi-directional

# Boltzmann Machine

Markov chain  $P(\mathbf{z}_t | \mathbf{z}_{t-1})$

$$P(\mathbf{z} = 1) = f(u) \quad , \quad u_i = \sum w_{ij} z_j - b_i$$



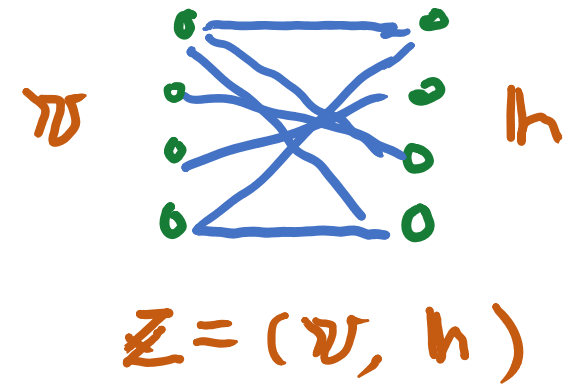
$$w_{ij} = w_{ji}$$

stable distr.  $P(\mathbf{z}) = \exp \left\{ \sum s_i z_i + \frac{1}{2} \sum w_{ij} z_i z_j - \psi \right\}$

# RBM: Restricted Boltzmann Machine

$$P(v, h) = \exp \left\{ b \cdot v + c \cdot h + h^T W v - \psi - \frac{1}{2} \|v\|^2 - \frac{1}{2} \|h\|^2 \right\}$$

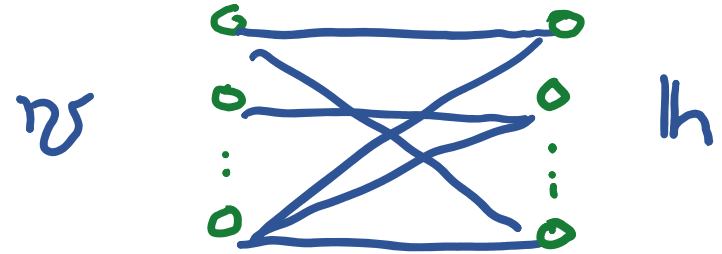
$\nearrow$   $\sum h_i w_{ij} v_j$        $\uparrow$   $\sum v_i^2$



energy machine :  $b, c$  : neglect

# RBM

$$z(v, h) = \exp \{ h^T W v - \psi \}$$



marginal distribution:  $z_v(v) = \sum_h z(v, h) \approx P(v)$

$$z_h(h) = \sum_v z(v, h)$$

conditional distribution

$$z(h|v) = \prod_i f \left( \sum_j W_{ij} v_j - \bar{t}_i \right),$$

$$z(v|h) = \prod_j f \left( \sum_i h_i W_{ij} - \bar{t}_j \right),$$

conditionally  
independent



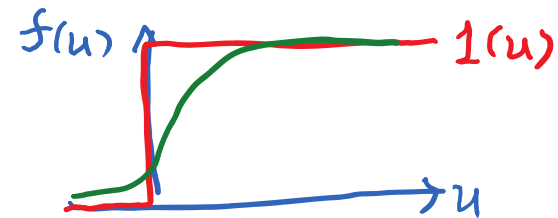
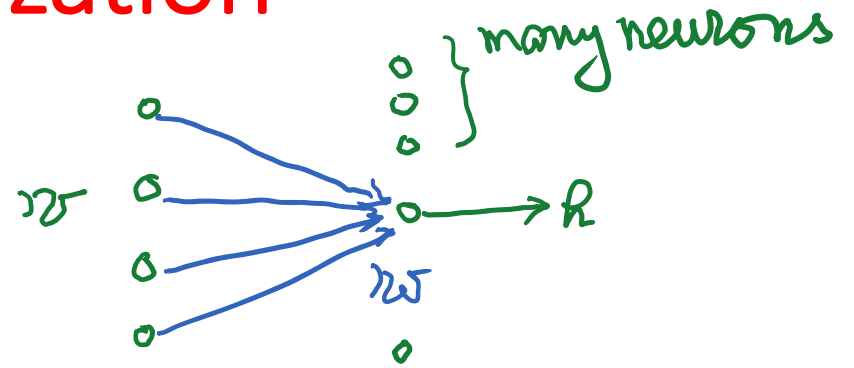
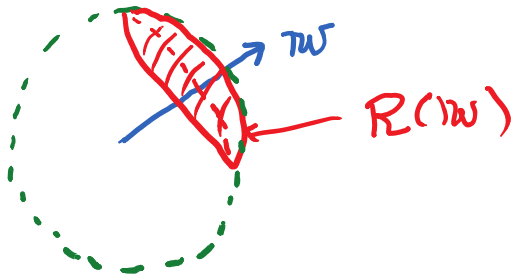
# Simple Hebbian Self-Organization

$$h = f(\mathbf{w} \cdot \mathbf{x} - \tau)$$

receptive field

$$R(\mathbf{w}) = \{ \mathbf{x} \mid \mathbf{w} \cdot \mathbf{x} - \tau > 0 \}$$

$$|\mathbf{w}|^2 = \text{const}$$



# self-organization of $\mathcal{W}$ : dynamics of $R(\mathcal{W})$

$\mathcal{W} : P(\mathcal{W})$

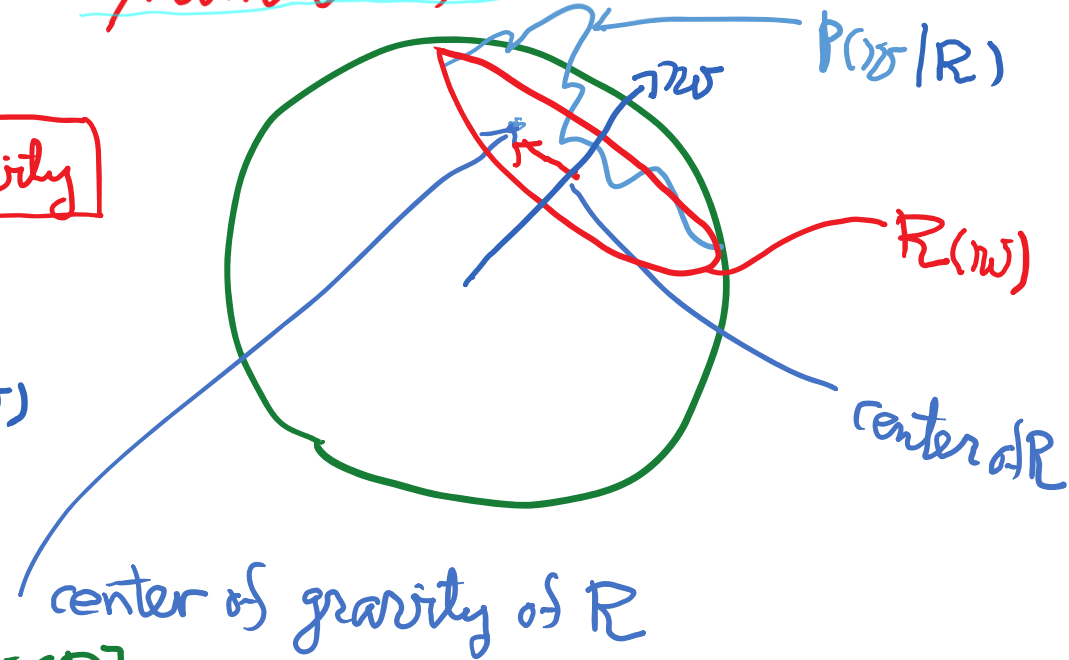
center  $\Rightarrow$  c. of gravity

$$\dot{\mathcal{W}} = \langle h(-\mathcal{W} + c\mathcal{W}) \rangle_{P(\mathcal{W})}$$

$$\frac{1}{P_R} \dot{\mathcal{W}} = -\mathcal{W} + c \langle h\mathcal{W} \rangle_R$$

$$P_{\text{rob}}\{\mathcal{W} \in R\} = \langle 1(\mathcal{W} \cdot \mathcal{W} - v_0) \rangle_{P(\mathcal{W})}$$

$$E[h\mathcal{W} | \mathcal{W} \in R]$$



# Self-Organization

$$\underset{W}{\text{minimize}} \text{KL} [ P(\nu) : \mathcal{Q}_\nu(\nu; W) ]$$

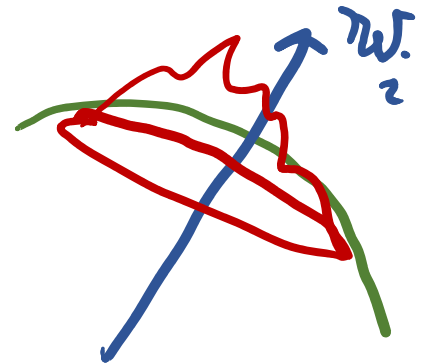
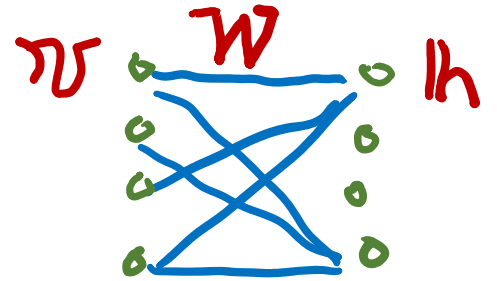
$$= \int P(\nu) \log \frac{P(\nu)}{\mathcal{Q}_\nu(\nu; W)} d\nu$$

$$\dot{w}_i = \varepsilon \{ -\langle R_i \nu \rangle_{\mathcal{Q}} + \langle R_i \nu \rangle_P \}$$

center  $w_i$  ?

center of gravity of  $R(w_i)$

$$\langle h \nu \rangle_{\mathcal{Q}} = \frac{\partial}{\partial w} \psi(w)$$

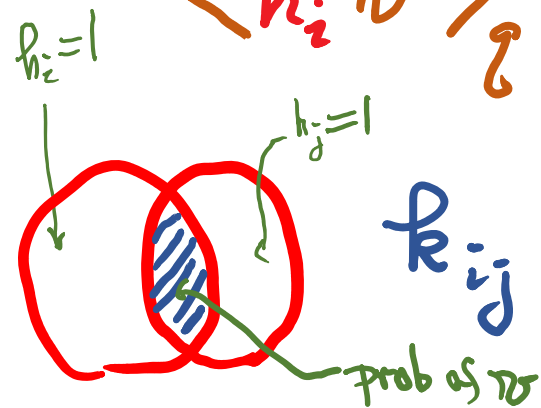


# Interaction of Hidden Neurons

$$\langle h_i v \rangle_{\mathcal{G}} = \frac{\partial \Psi(W)}{\partial W} \quad ; \quad \mathcal{G}(h, v) = \exp \left\{ h^T W v - \underbrace{\Psi(W)}_{-\frac{1}{2} |v|^2} \right\}$$

$$\begin{aligned} \Psi(W) &= \log \sum_{\mathbf{h}} \int \exp \left\{ -\frac{1}{2} |v|^2 + h^T W v \right\} dv \\ &= \log \sum_{\mathbf{h}} \exp \left\{ \frac{1}{2} |h^T W|^2 \right\} + C \end{aligned}$$

$$\langle h_i v \rangle_{\mathcal{G}} = \sum R_{ij} W_{ij} \quad ; \quad \text{interaction}$$



$$R_{ij} = E_{\mathcal{G}} [h_i h_j]$$

← joint firing prob.

$v$ : analog, Gaussian

# Bayesian Duality in Exponential Family

Data  $x$       Parameter (higher-order concepts)  $\theta$

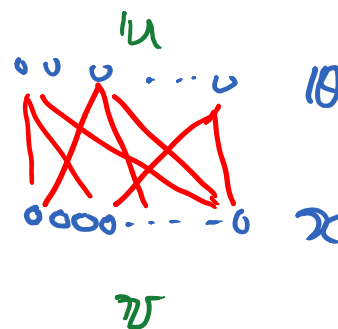
$$P(x|\theta) = \exp\{\theta \cdot x - \bar{R}(x) - \Psi(\theta)\}$$

$$P(x, \theta) = \exp\{\theta \cdot x - \bar{R}(x) - \bar{\Psi}(\theta)\} \quad \Psi - \log \pi(\theta)$$

$$P(\theta|x) = \exp\{\theta \cdot x - \bar{R}(x) - \bar{\Psi}(\theta)\}$$

Curved exponential family

$$\mathcal{X}(\mathcal{U}), \quad \Theta(\mathcal{U}) : \exp\{\theta(\mu) \cdot x(\mu) - \bar{R} - \bar{\Psi}\}$$



**RBM** : curved exp. family

$$\theta = h, \quad x = Wv$$

$$x = v \quad \theta = hW$$

$$P(x, h) = \exp \left\{ \underbrace{\theta \cdot x}_{h^T W v} - \bar{R} - \bar{\Psi} \right\}$$

# Gaussian Boltzmann Machine

$$g(\boldsymbol{w}, \boldsymbol{h}) = \exp\left\{-\frac{1}{2}|\boldsymbol{w}|^2 - \frac{1}{2}|\boldsymbol{h}|^2 + \boldsymbol{h}^T \boldsymbol{W} \boldsymbol{w} - \psi(\boldsymbol{w})\right\}$$

$$\langle \boldsymbol{h} \boldsymbol{w} \rangle_{\boldsymbol{h}} = \boldsymbol{W} \boldsymbol{C} \quad \boldsymbol{C} = \int \boldsymbol{w} \boldsymbol{w}^T p(\boldsymbol{w}) d\boldsymbol{w}$$

*center of gravity* *covariance matrix*

$$\langle \boldsymbol{h} \boldsymbol{w} \rangle_{\boldsymbol{w}} = (\boldsymbol{I} - \boldsymbol{W} \boldsymbol{W}^T)^{-1} \boldsymbol{W}$$

*interactions of  $\boldsymbol{h}$  neurons*

# Equilibrium Solution

$$WC = (I - WW^T)^{-1} W$$

$$W = \begin{pmatrix} w_1 \\ \vdots \\ w_k \end{pmatrix}$$

solution :  $w_i \cdot w_j = 0$  ( $i \neq j$ ),  $|w_i|^2 = m_i$

$$w_i C = \underbrace{(1 - m_i)^{-1}}_{\lambda_i} w_i$$

$$w_i C = \lambda_i w_i$$

PCA!

$$m_i = 1 - \frac{1}{\lambda_i}$$



# Equilibrium Solution

$$WC = (I - WW^T)^{-1}W$$

**General Solution**  $W = U \text{diag} \left( \sqrt{1 - \frac{1}{\lambda_1}}, \dots, \sqrt{1 - \frac{1}{\lambda_m}}, 0, \dots, 0 \right) V$

- orthogonal matrix :  $U, V$
- $C$  diagonalized by  $V$   $C = V^T \text{diag}(\lambda_1, \dots, \lambda_n) V$

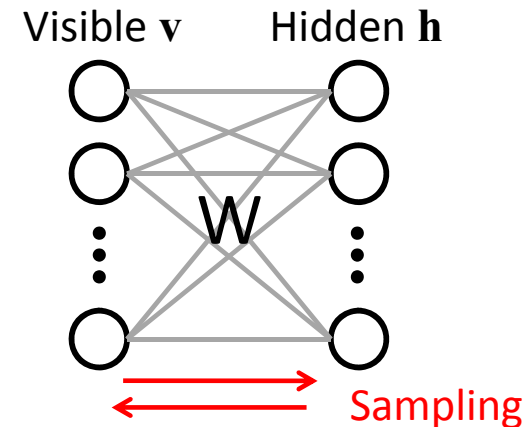
You can choose  $m(\leq k)$  eigen values form  $\lambda_1, \dots, \lambda_k > 1$

**Stable Solution** the case of  $m = k$

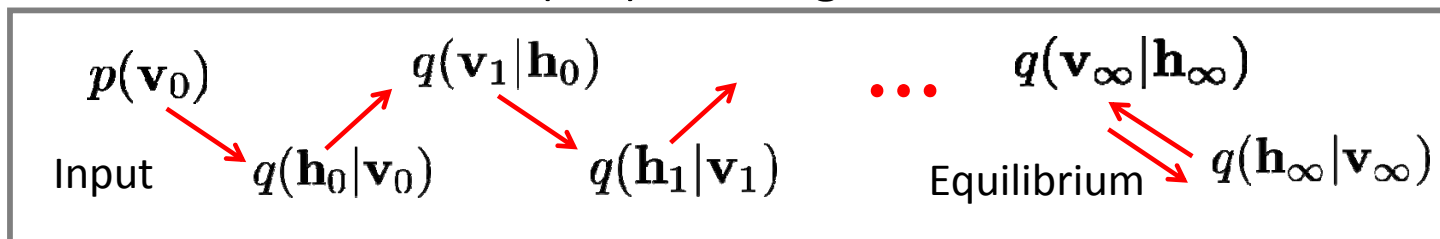
# Contrastive Divergence

- RBM
  - 2-layered probabilistic neural network
  - No connections within layers

$$q(\mathbf{h}, \mathbf{v}; W) = \exp(-\mathbf{h}^T W \mathbf{v}) / \sum_{\mathbf{h}, \mathbf{v}} \exp(-\mathbf{h}^T W \mathbf{v})$$



- How to train RBM  
Maximum Likelihood (ML) learning is hard



Many iterations of Gibbs Sampling demand **too much computational time**

# Contrastive Divergence Solution

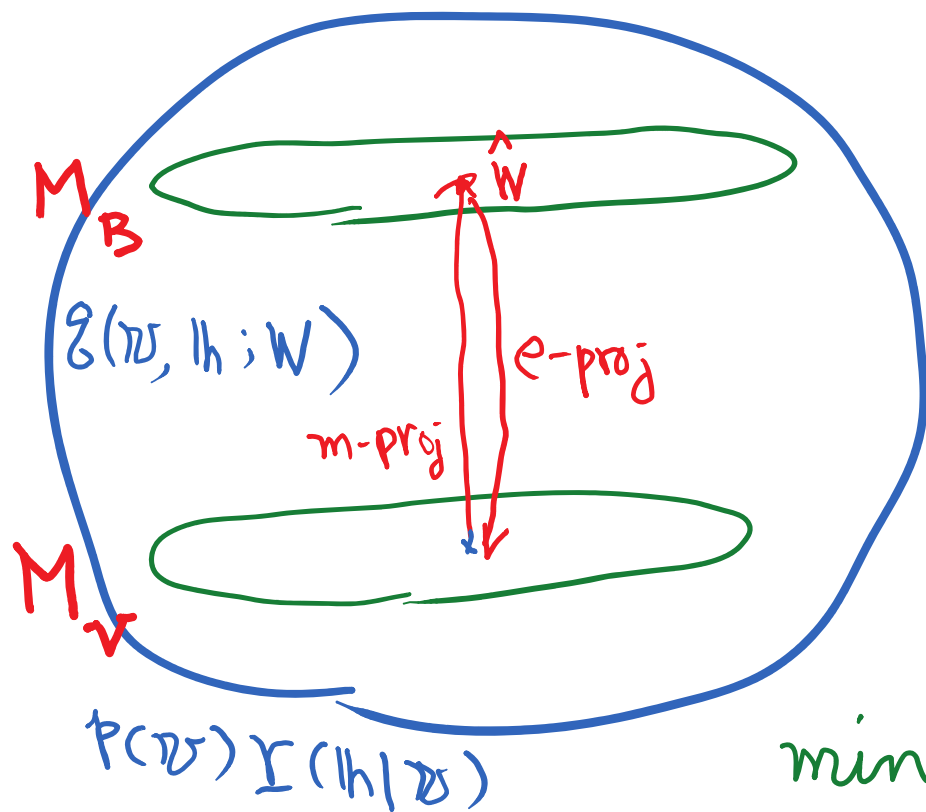
$$\langle h v \rangle_{\tilde{q}} = W C W^T W + W$$

$$W C = W C W^T W + W$$

$$\pi w_i C = \frac{1}{1 - m_i} \pi w_i$$

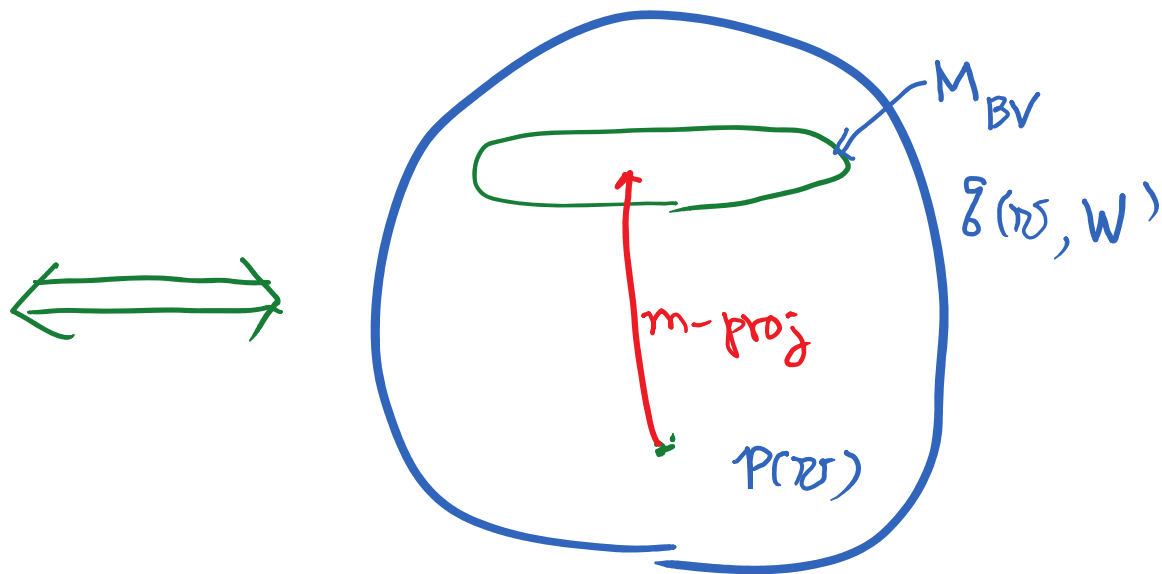
PCA!

# Two Manifolds

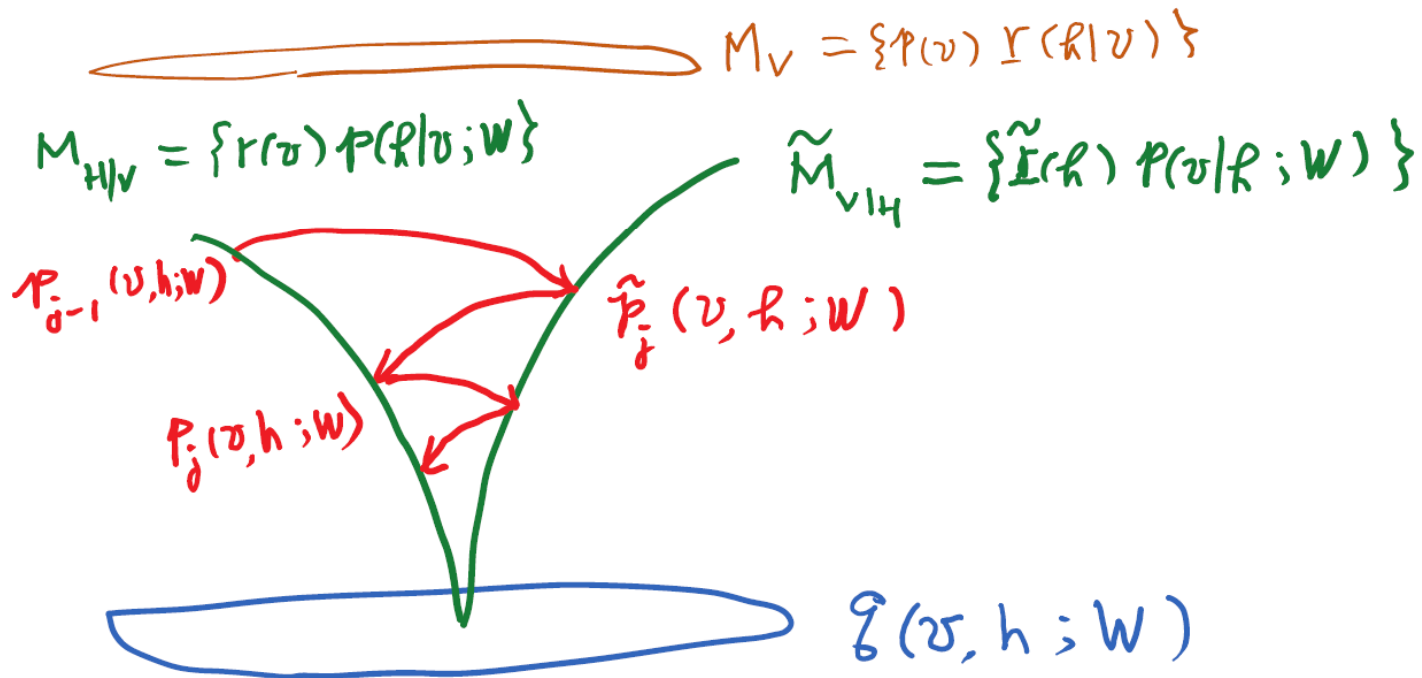


$$\min_{I, W} KL [ P(v) \sqcup (h|v) : g(v, h; W) ]$$

$$\min_{-W} KL [ P(v) : g_v(v; W) ]$$



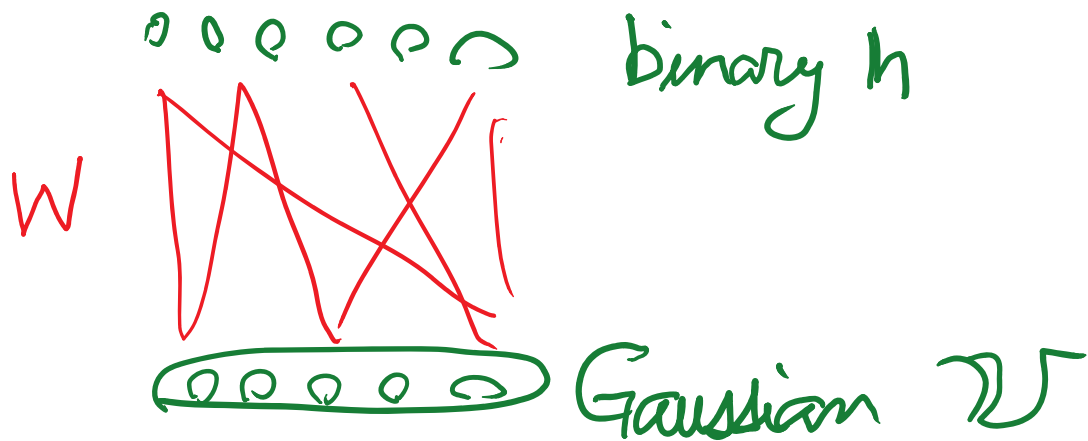
## Geometry of CDn (contrastive divergence)



# Bernoulli-Gaussian RBM

ICA

R. Karakida

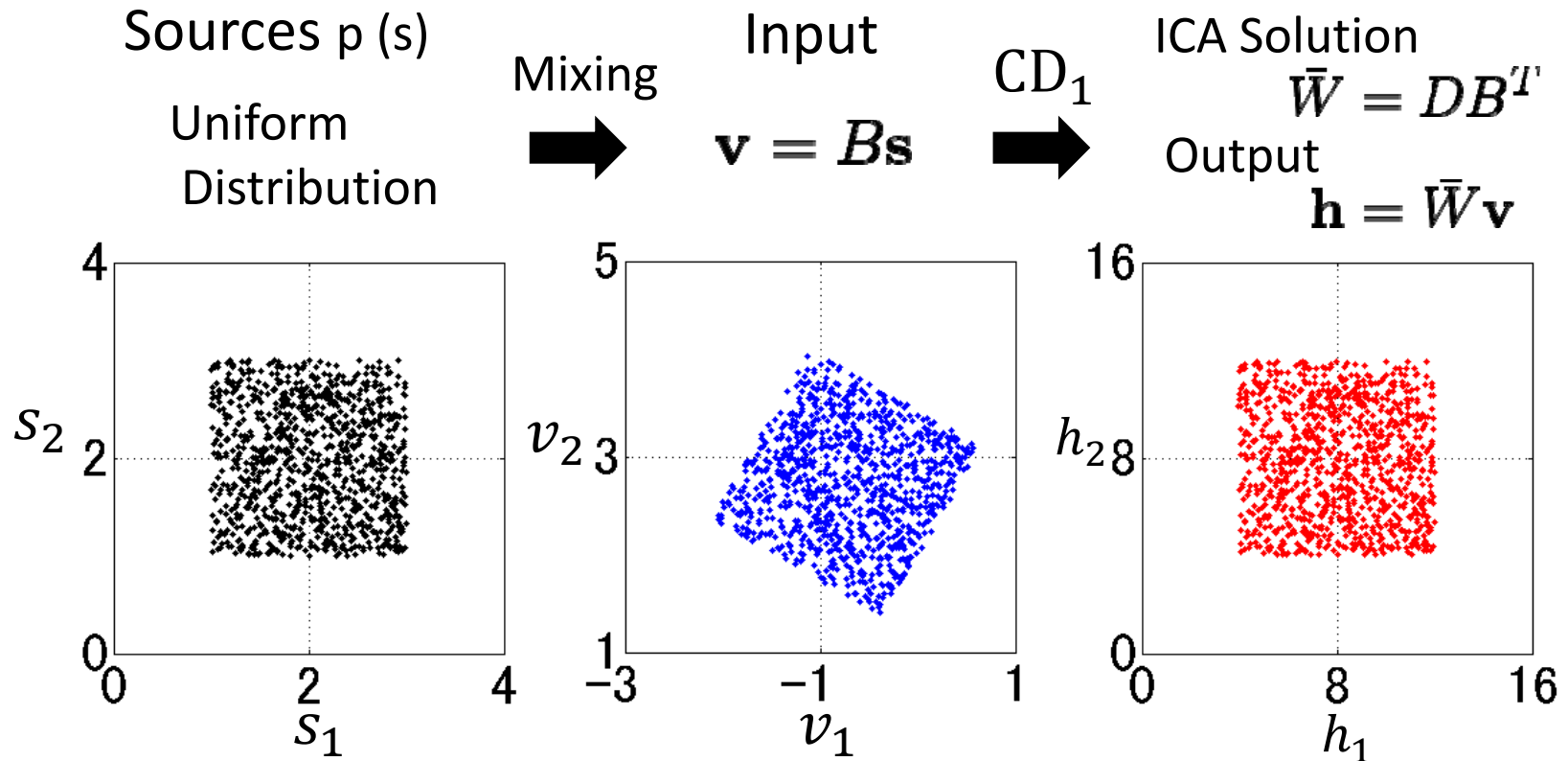


$$P(v) : \quad v = \underline{0D5}$$

$$W = A^{-1}$$

# Simulation

The number of Neurons:  $N = M = 2$ ,  $\sigma = 1/2$



Independent sources are extracted in G-B RBM

# Information Geometry of Neuromanifolds

Natural Gradient and Singularities

Multilayer Perceptron

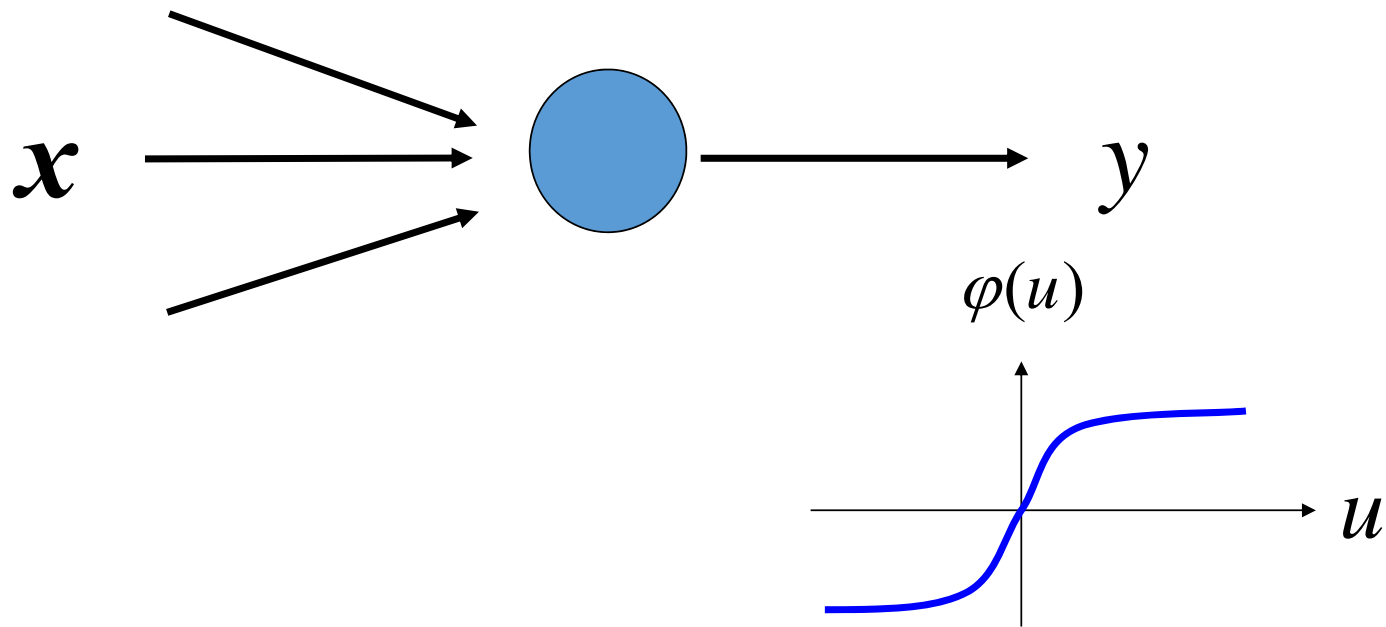
Shun-ichi Amari

RIKEN Brain Science Institute



# Mathematical Neurons

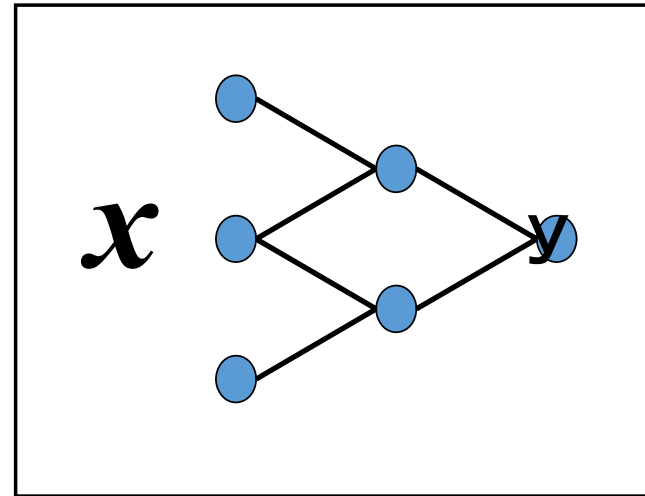
$$y = \varphi\left(\sum w_i x_i - h\right) = \varphi(\mathbf{w} \cdot \mathbf{x})$$



## Multilayer Perceptrons

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

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

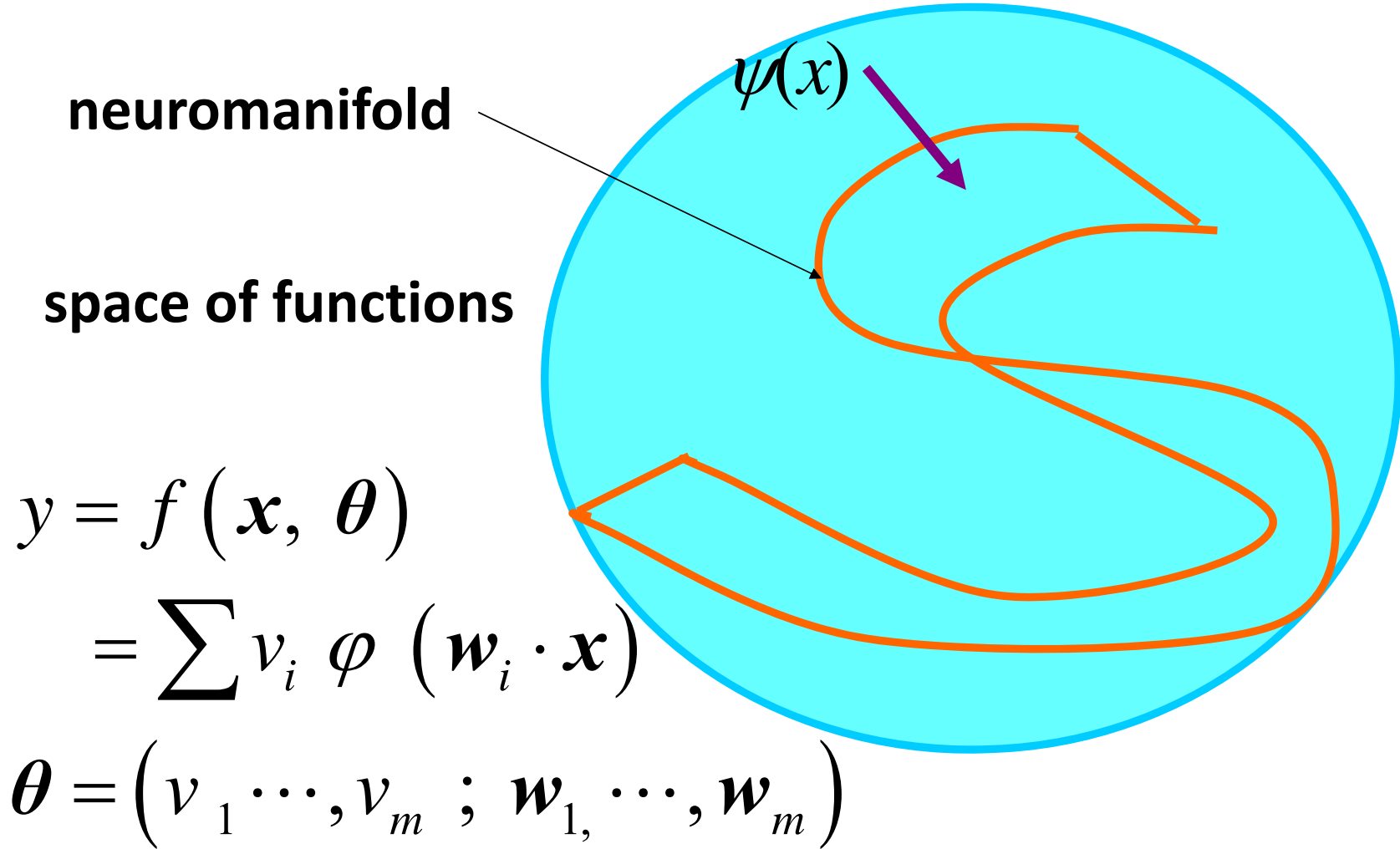


$$p(y|\mathbf{x};\boldsymbol{\theta}) = c \exp\left\{-\frac{1}{2}(y - f(\mathbf{x},\boldsymbol{\theta}))^2\right\}$$

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

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

# Multilayer Perceptron: Neuromanifold



# Universal Function Approximator

$$\psi(\mathbf{x}) \approx \sum_{i=1}^N v_i a_i(\mathbf{x})$$

$$\psi(\mathbf{x}) \approx \sum_{i=1}^m v_i \varphi_i(\mathbf{x})$$

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

# Learning from examples

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

examples  $\cdots D = \{(\mathbf{x}_1, y_1), \cdots, (\mathbf{x}_n, y_n)\}$

learning ; estimation

# Backpropagation ---gradient learning

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

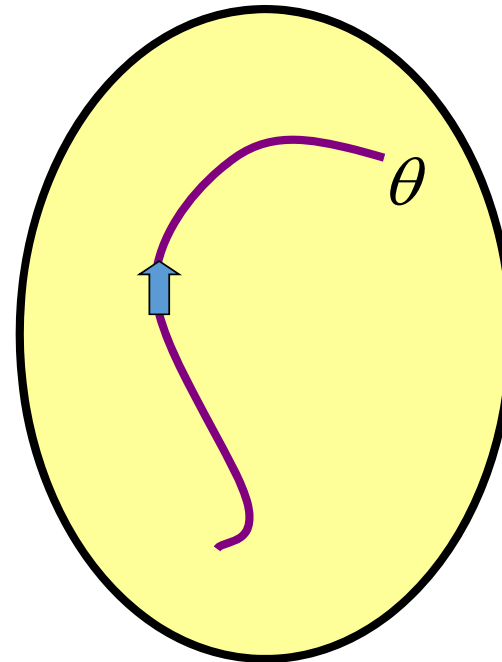
$$y = f(x, \theta) + n$$

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

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

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

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

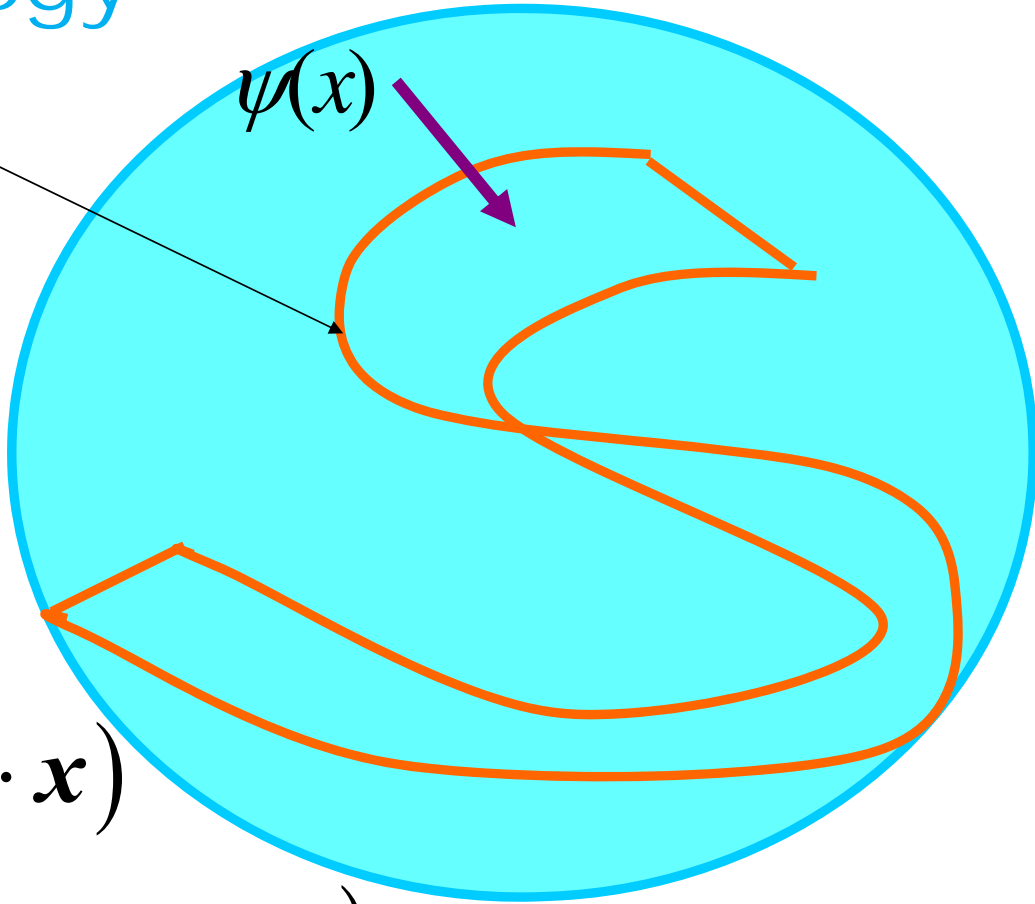


# Multilayer Perceptron: Neuromanifold

## Metric and Topology

neuromanifold

space of functions

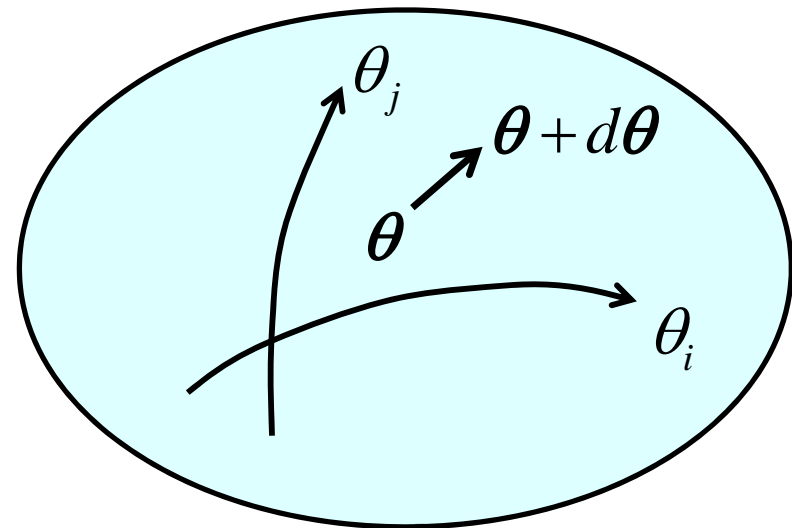


$$y = f(\mathbf{x}, \boldsymbol{\theta})$$
$$= \sum v_i \varphi(\mathbf{w}_i \cdot \mathbf{x})$$
$$\boldsymbol{\theta} = (v_1 \cdots, v_m ; \mathbf{w}_1, \cdots, \mathbf{w}_m)$$

# Metric: Riemannian manifold

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

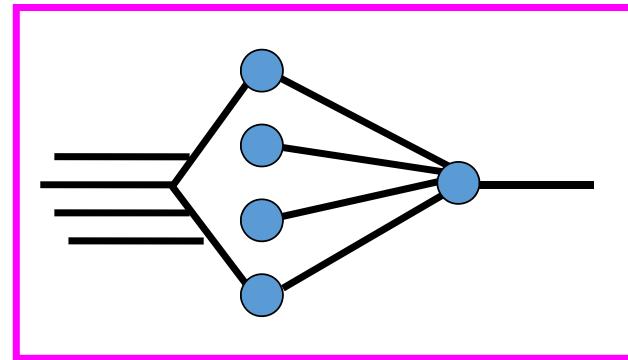
$$\begin{aligned} ds^2 &= |d\boldsymbol{\theta}|^2 \\ &= \sum g_{ij}(\boldsymbol{\theta}) d\theta_i d\theta_j \\ &= d\boldsymbol{\theta}^T G(\boldsymbol{\theta}) d\boldsymbol{\theta} \end{aligned}$$



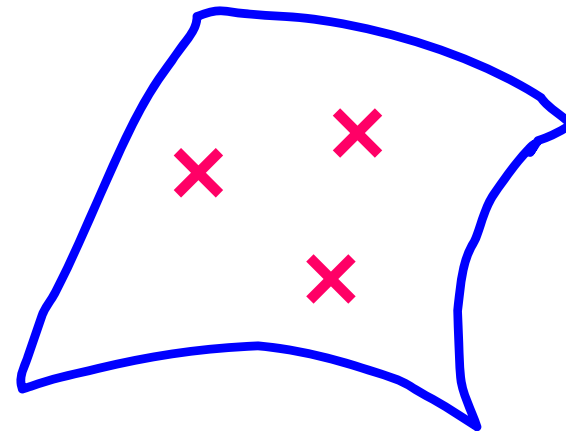
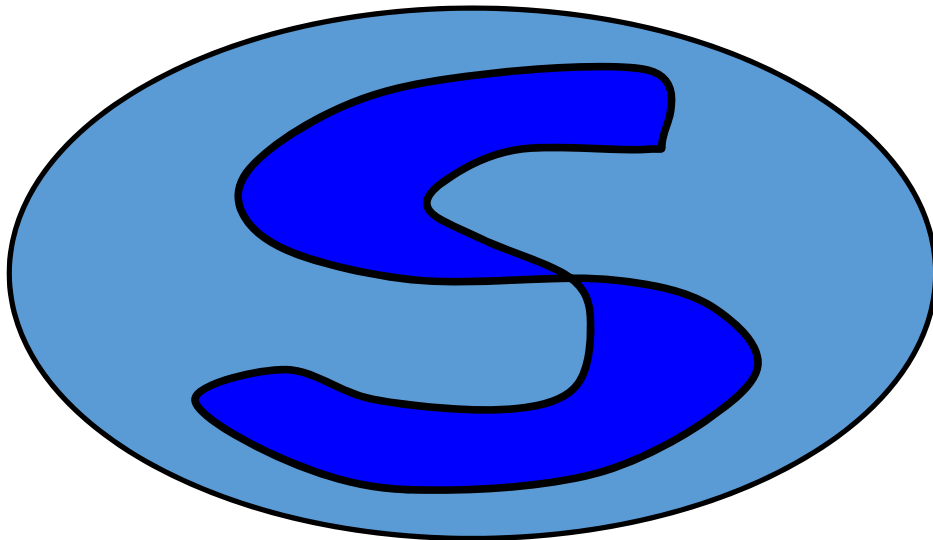


# Topology: Neuromanifold

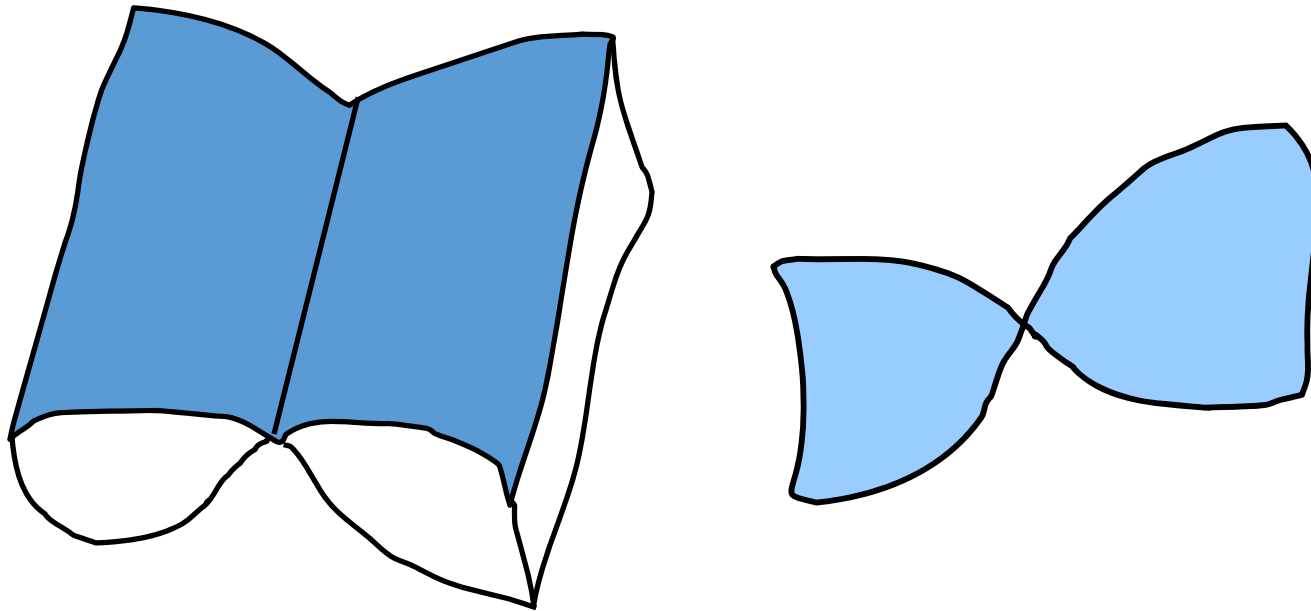
- Metrical structure
- Topological structure



$\theta$



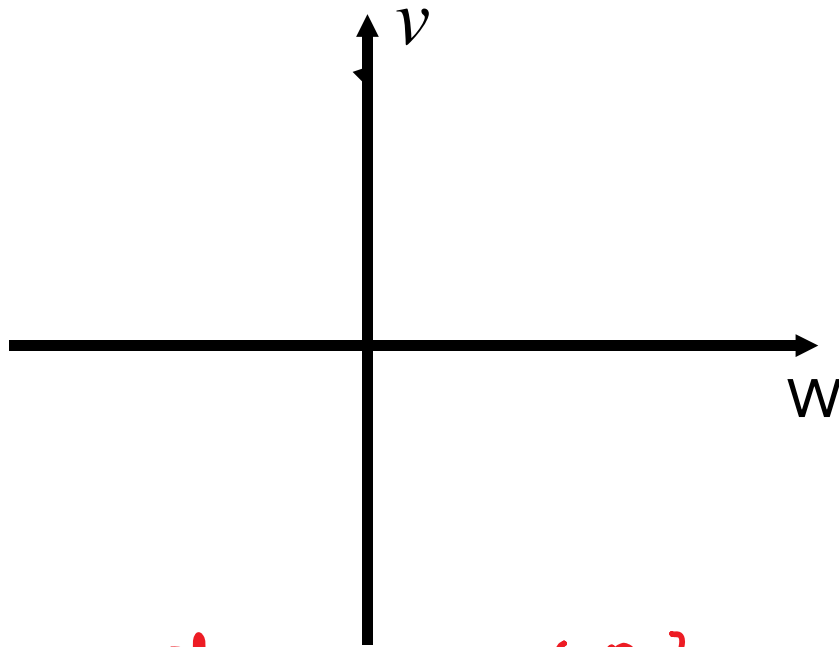
# singularities



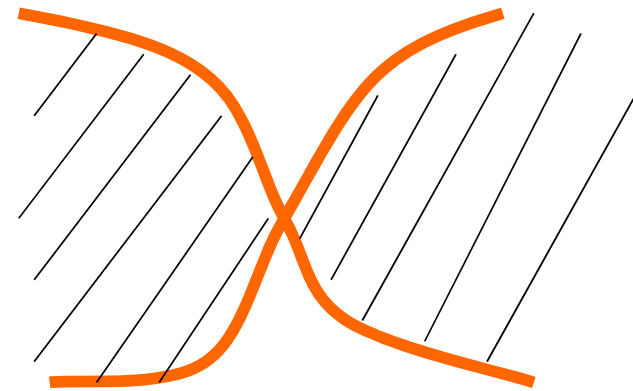
# Geometry of singular model

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

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



parameter space  $\{\theta\}$



behavioral space  $\{f(x, \theta)\}$

## Parameter Space

$S$

$$S = \{\boldsymbol{\theta}\}$$

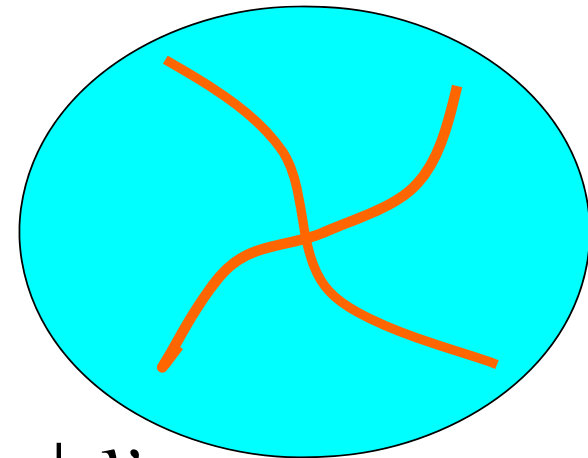
$$y = \sum v_i \varphi(\mathbf{w}_i \cdot \mathbf{x}) + n$$

## Equivalence

1)  $v_i \mathbf{w}_i = \mathbf{0}$

2)  $\mathbf{w}_i = \mathbf{w}_j \Rightarrow v_i + v_j$

$$M = S / \approx$$

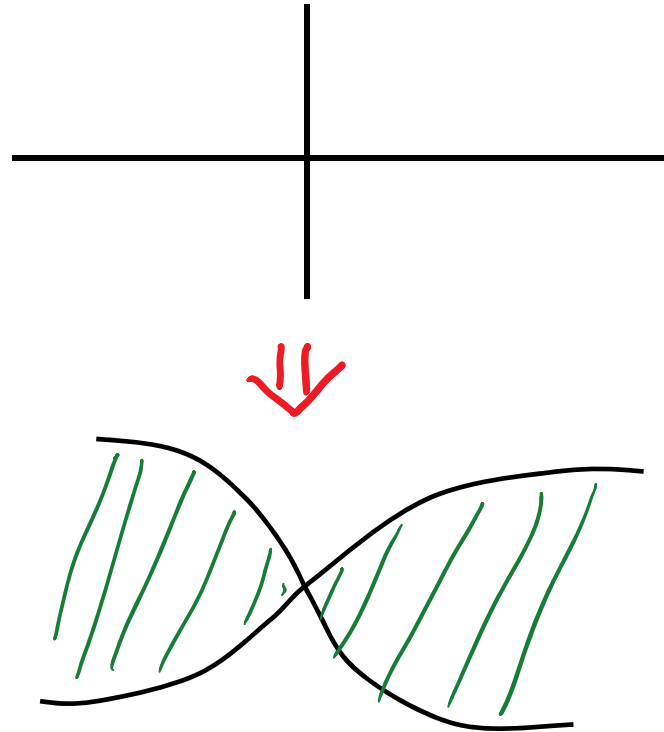


## Topological Singularities

$S$  : parameter space

$$M = S / \approx$$

$M$  : behavioral space

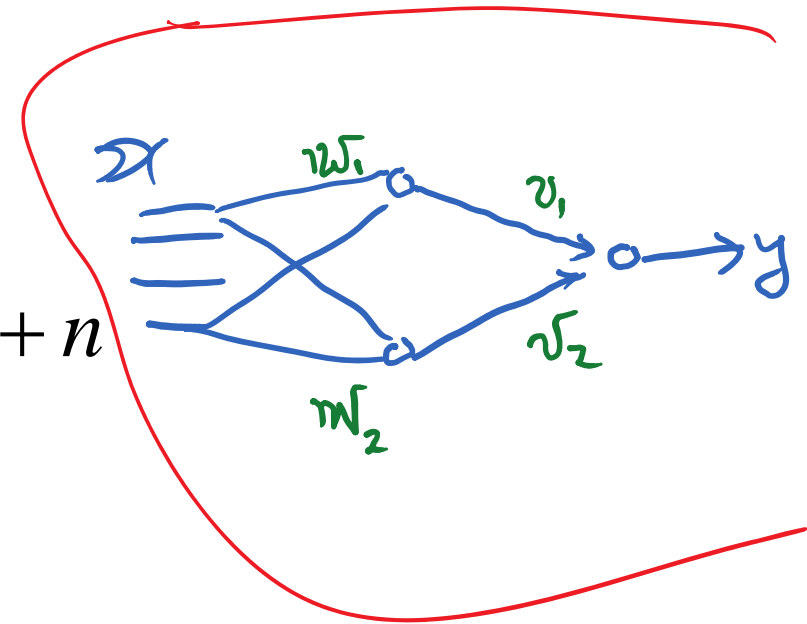


## 2 hidden-units

$$y = v_1 \varphi(\mathbf{w}_1 \cdot \mathbf{x}) + v_2 \varphi(\mathbf{w}_2 \cdot \mathbf{x}) + n$$

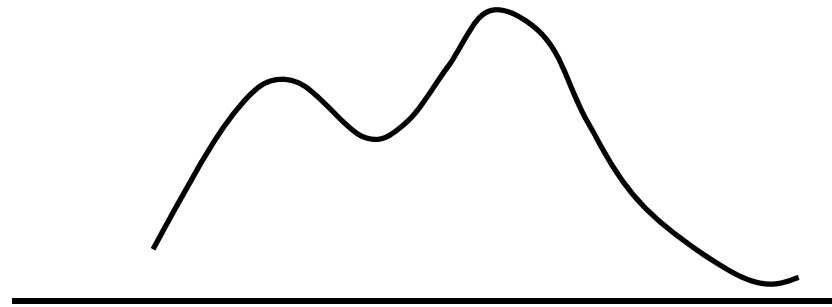
$$S : v_1 v_2 \|\mathbf{w}_1 - \mathbf{w}_2\| \|\mathbf{w}_1 + \mathbf{w}_2\| = 0$$

$$(1 - v) \varphi(x - w_1) + v \varphi(x - w_2)$$



## Gaussian mixtures

$$p(x) = \sum v_i \exp \left\{ -\frac{1}{2} (x - w_i)^2 \right\}$$



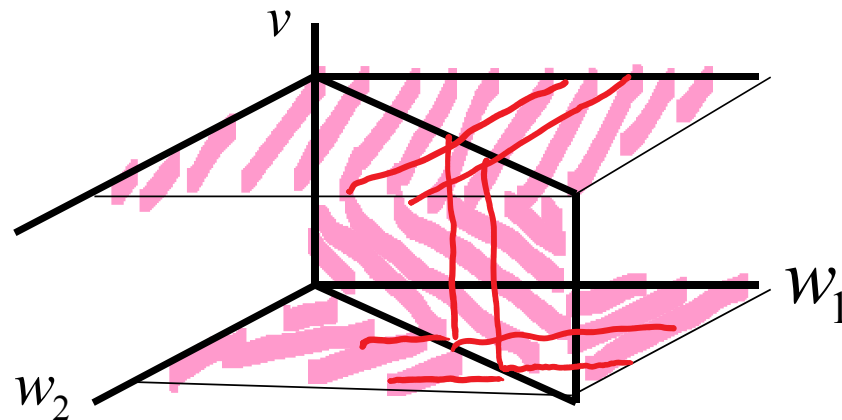
## Gaussian mixture

$$p(x; v, w_1, w_2) = (1-v)\varphi(x-w_1) + v\varphi(x-w_2)$$

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}x^2\right\}$$

singular:

$$w_1 = w_2, \quad v(1-v) = 0$$





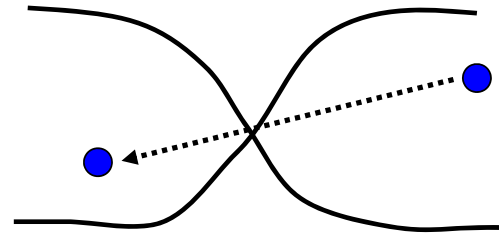
# Learning, Estimation, and Model Selection

$$E_{\text{gen}} = D \left[ p_0(y|\mathbf{x}) : p(y|\mathbf{x}; \hat{\boldsymbol{\theta}}) \right]$$

$$E_{\text{train}} = D \left[ p_{\text{emp}}(y|\mathbf{x}; \hat{\boldsymbol{\theta}}) \right]$$

$$E_{\text{gen}} = \frac{d}{2n} \quad d : \text{dimension}$$

$$E_{\text{gen}} = E_{\text{train}} + \frac{d}{n}$$



$$\frac{\log n}{n}, \frac{\log \log n}{n}, \dots$$

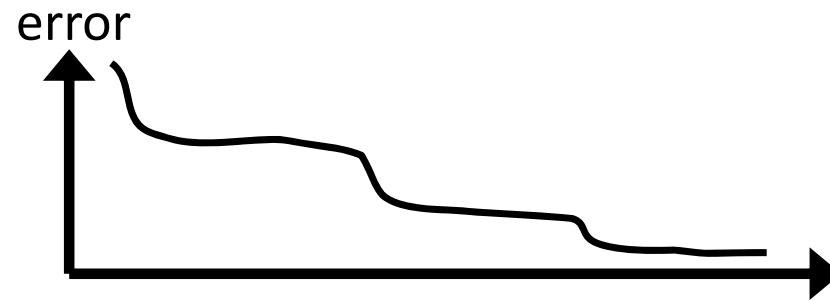
# Problem of Backprop

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

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

# Flaws of MLP

**slow convergence : plateau**

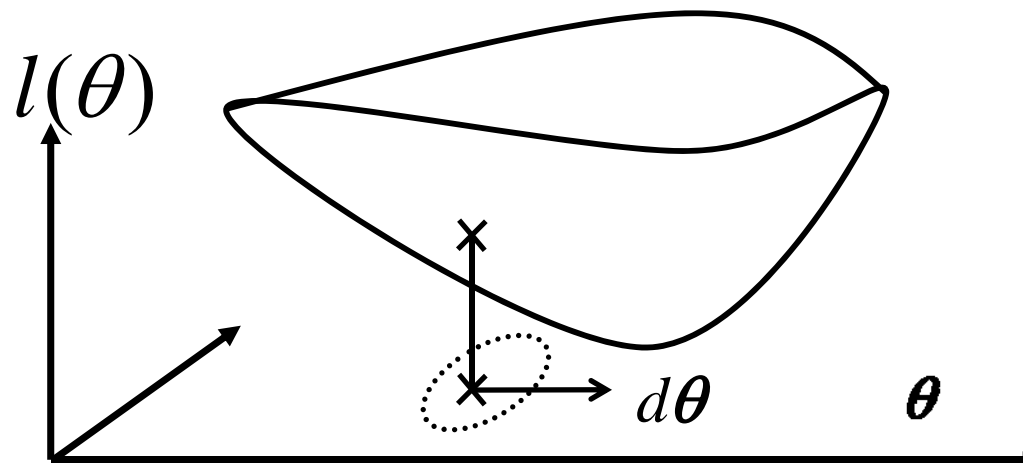


**local minima**



**→ Boosting, Bagging, SVM**

# Steepest Direction --- Natural Gradient



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

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

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

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

# Natural Gradient

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

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

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

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

# Information Geometry of MLP

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

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

Adaptive natural gradient learning

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

## Regular statistical model

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

$G$ : Fisher information

$$E[\Delta\theta\Delta\theta^T] = \frac{1}{n}G^{-1}$$

$$E\left[KL\left[p(x, \theta_0) : p(x, \hat{\theta})\right]\right] \approx \frac{1}{2n}G \cdot E[\Delta\theta\Delta\theta]$$
$$\approx \frac{d}{2n}$$

AIC, MDL

## Landscape of error at singularity

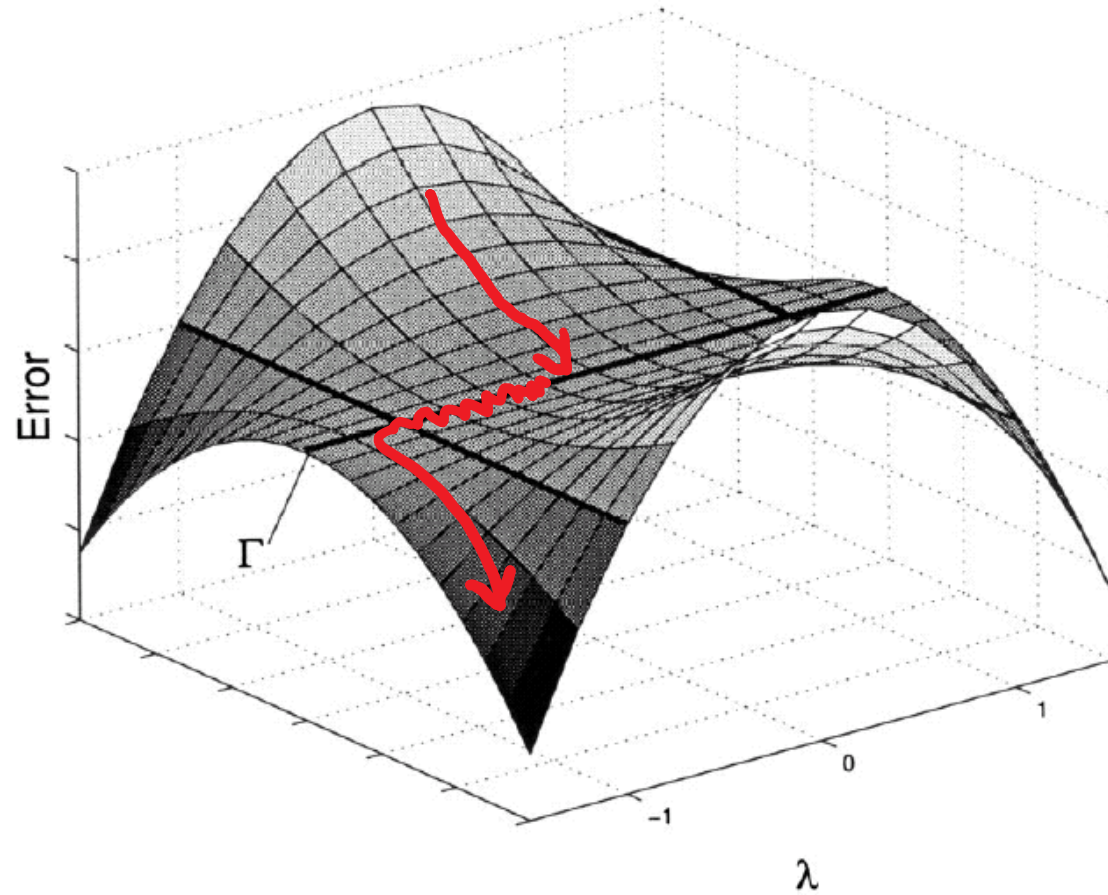
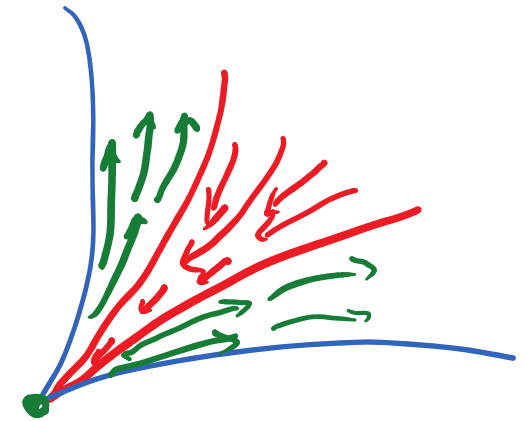


Fig. 5. Critical set with local minima and plateaus.

## Milner attractor





# Dynamics of Learning

$$\frac{d\theta}{dt} = -\eta \nabla l, \quad \frac{d\theta}{dt} = -\eta G^{-1} \nabla l$$

$$\frac{du}{dt} = f(u, z), \quad \frac{dz}{dt} = k(u, z)$$

$$\frac{du}{dz} = \frac{f(u, z)}{k(u, z)}, \quad u^2 = z^2 - \frac{1}{2} \log |z| + c$$

# Coordinate Transformation

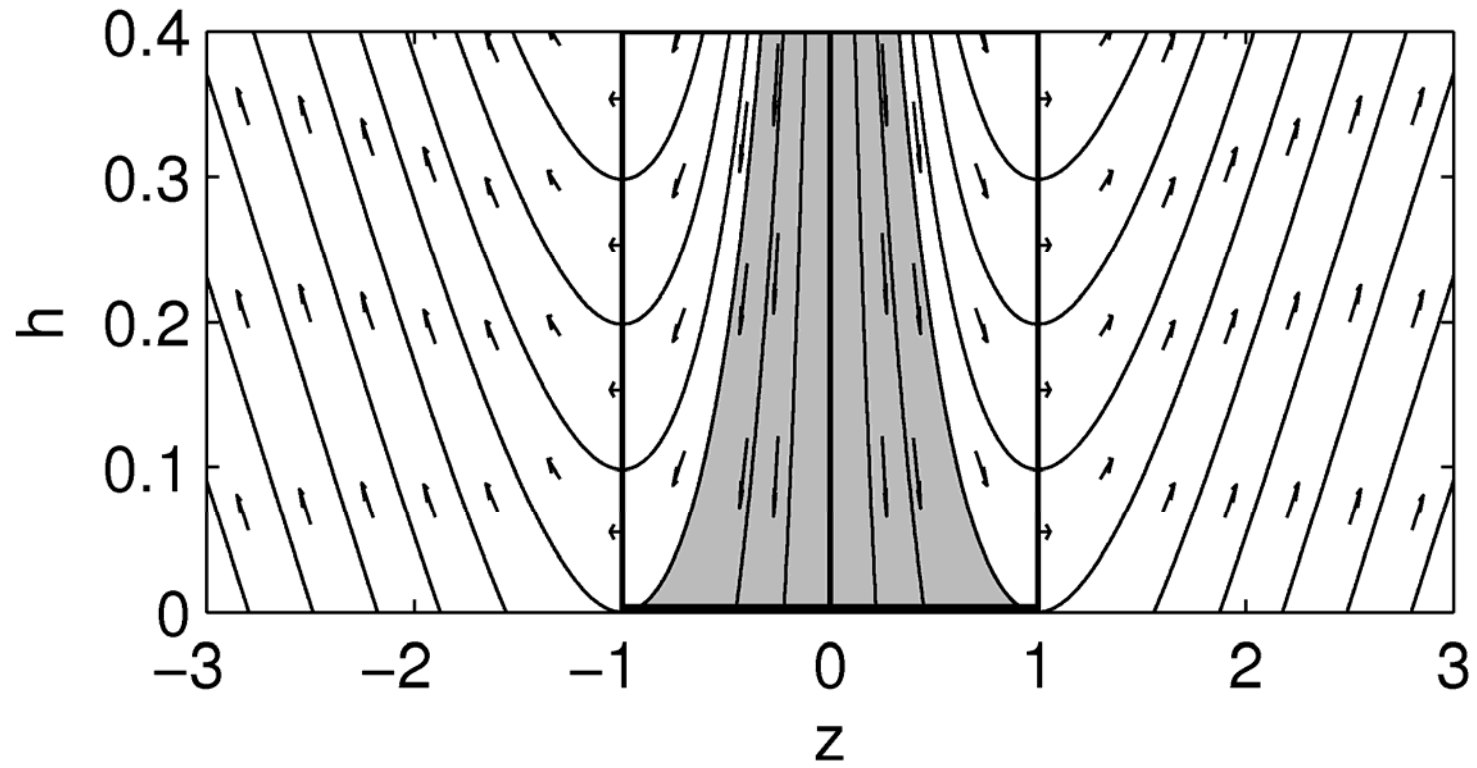
$$\left\{ \begin{array}{l} \mathbf{u} = \mathbf{w}_2 - \mathbf{w}_1 \\ \mathbf{w} = \frac{v_1 \mathbf{w}_1 + v_2 \mathbf{w}_2}{v} \\ v = v_1 + v_2 \\ z = \frac{v_2 - v_1}{v} \end{array} \right. \quad \begin{array}{l} : \mathbf{u} = 0 \\ \mathbf{w} = \mathbf{w}^* \\ v = v^* \\ z = \pm 1 \end{array} \quad \begin{array}{l} \mathcal{R}_1 \\ \\ \\ \mathcal{R}_2 \end{array}$$

# Dynamics of Learning: Natural gradient works!!

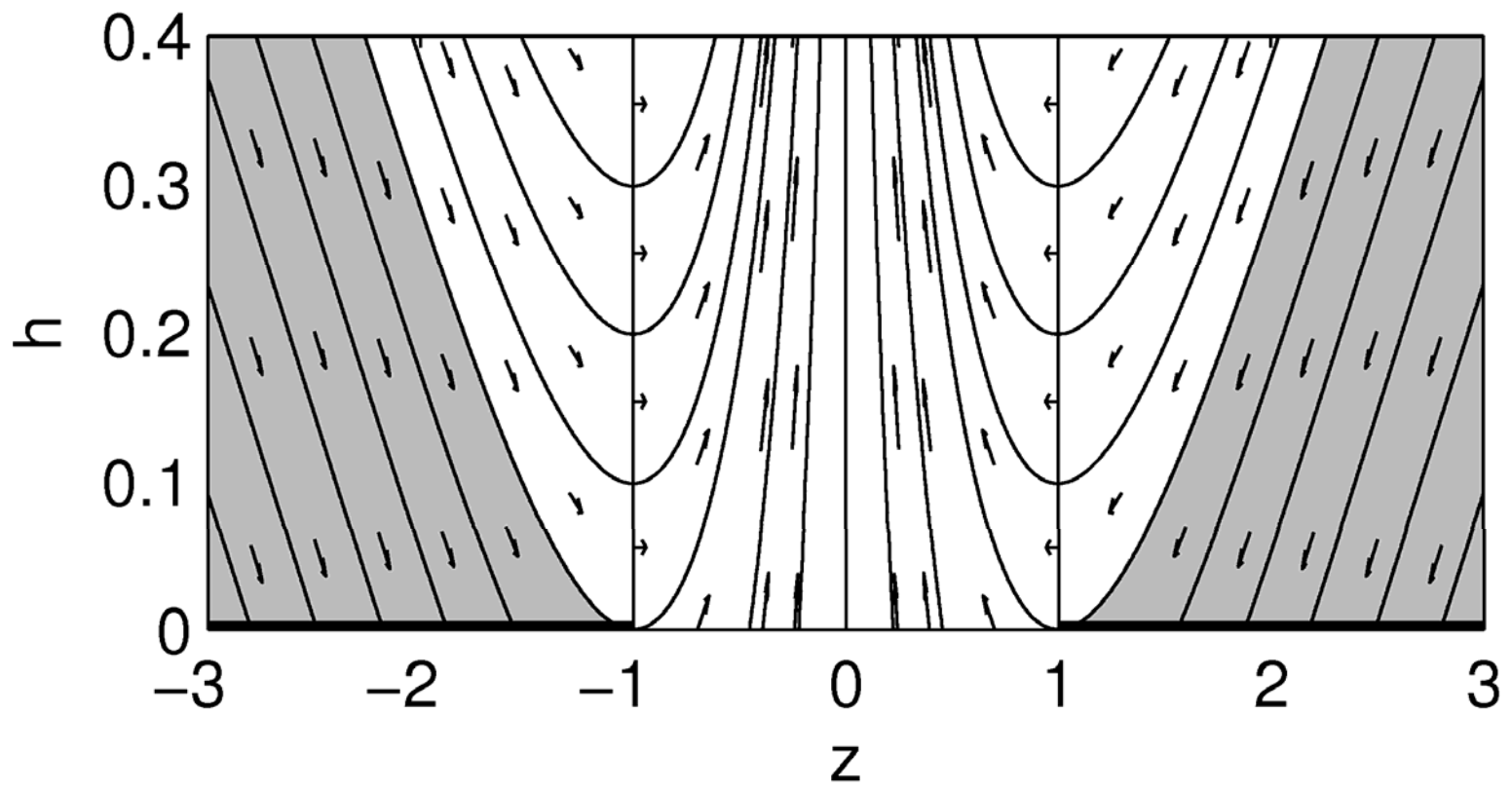
$$\frac{d\theta}{dt} = -\eta \nabla l, \quad \frac{d\theta}{dt} = -\eta G^{-1} \nabla l$$

$$\frac{du}{dt} = f(u, z), \quad \frac{dz}{dt} = k(u, z)$$

$$\frac{du}{dz} = \frac{f(u, z)}{k(u, z)}, \quad u^2 = z^2 - \frac{1}{2} \log |z| + c$$

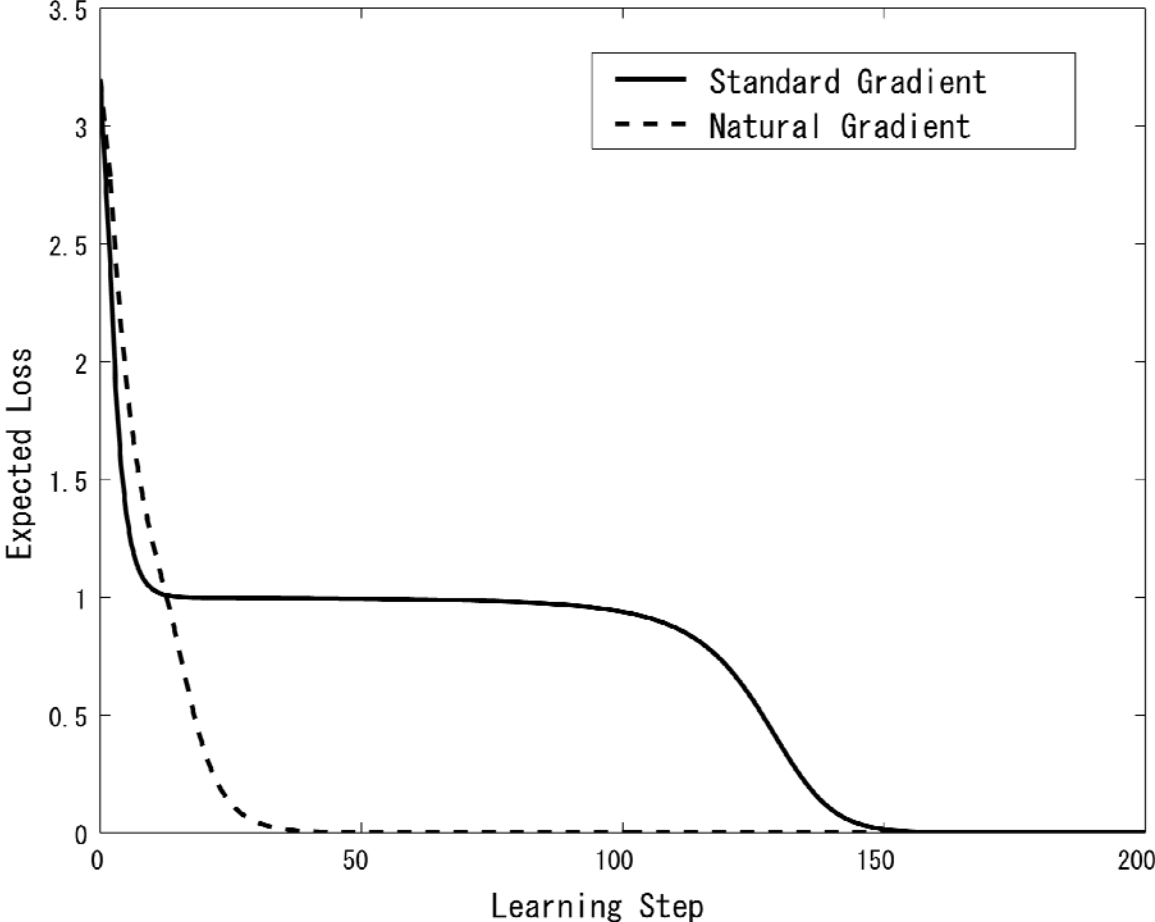


Dynamic vector fields: General case ( $|z| < 1$  part stable)



Dynamic vector fields: General case ( $|z| > 1$  part stable )

# Adaptive Natural Gradient works well



# Signal Processing

## ICA : Independent Component Analysis

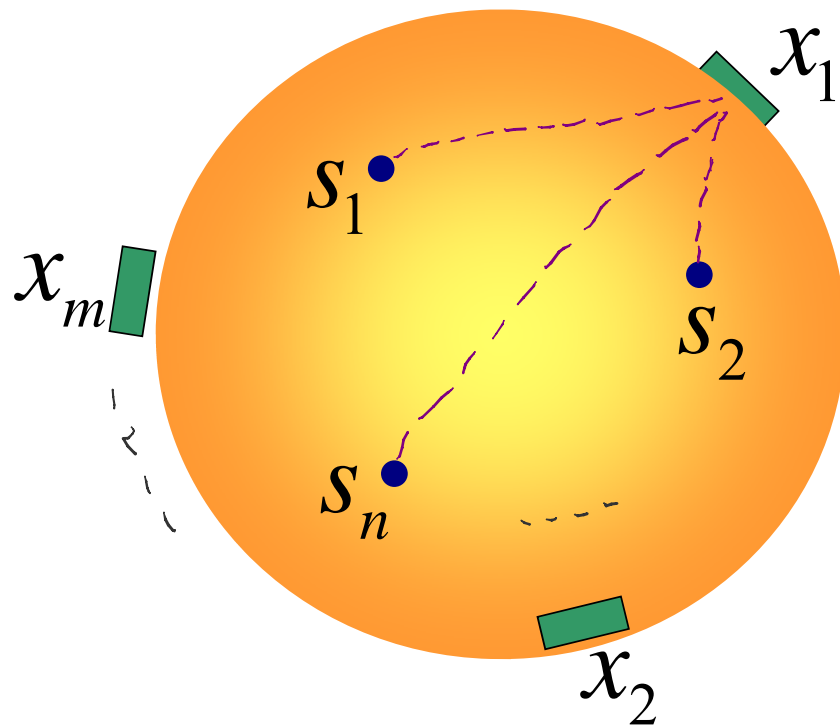
$$\mathbf{x}_t = A\mathbf{s}_t \quad \mathbf{x}_t \rightarrow \mathbf{s}_t$$

$$\begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix} = (A) \begin{pmatrix} s_1 \\ \vdots \\ s_n \end{pmatrix}$$

**sparse component analysis**

**positive matrix factorization**

# mixture and unmixture of independent signals



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

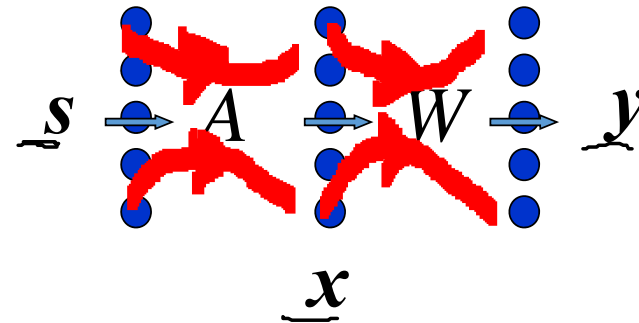
$$\mathbf{x} = \mathbf{A}\mathbf{s}$$



# Independent Component Analysis

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

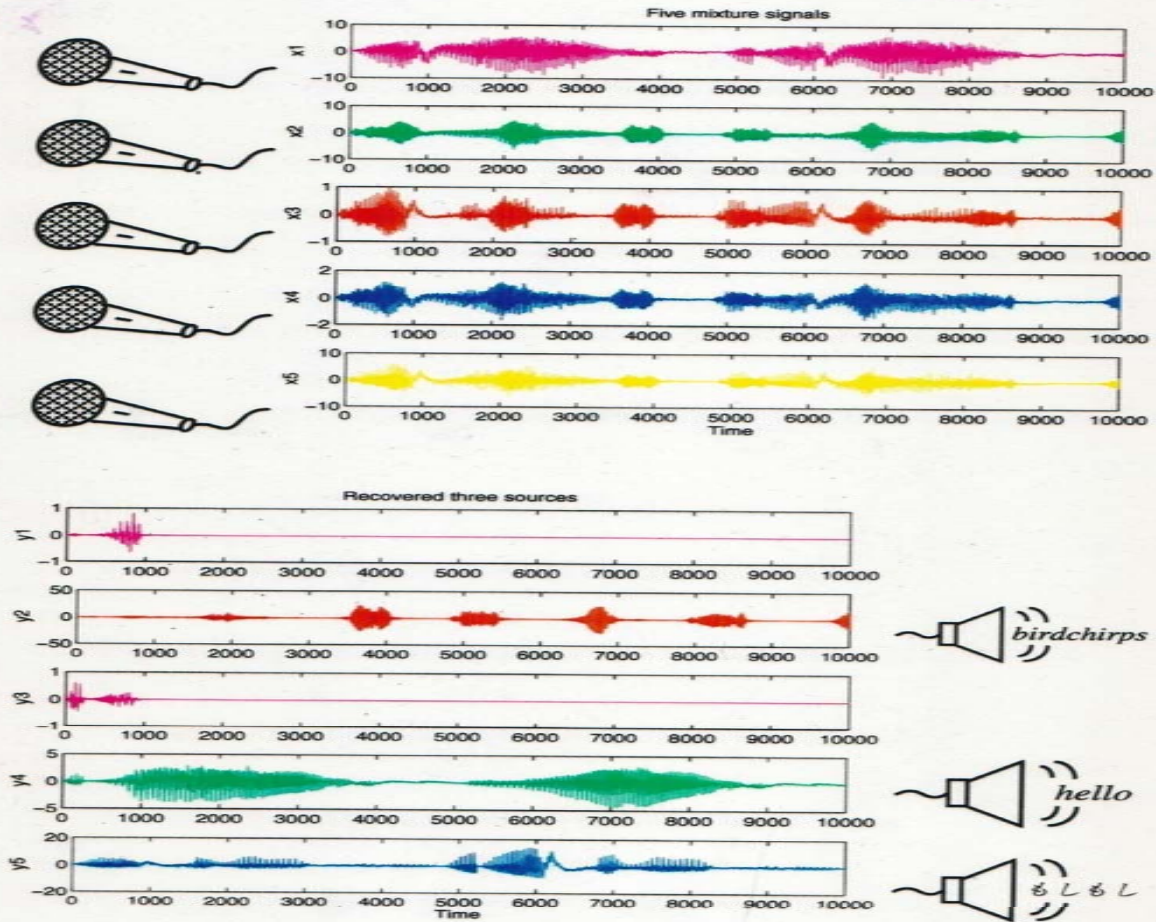
$$\mathbf{y} = W\mathbf{x} \quad W = A^{-1}$$



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

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

## Cocktail party experiment



- 5 microphones (sensors) and only 3 speakers

**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

# Semiparametric Statistical Model

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

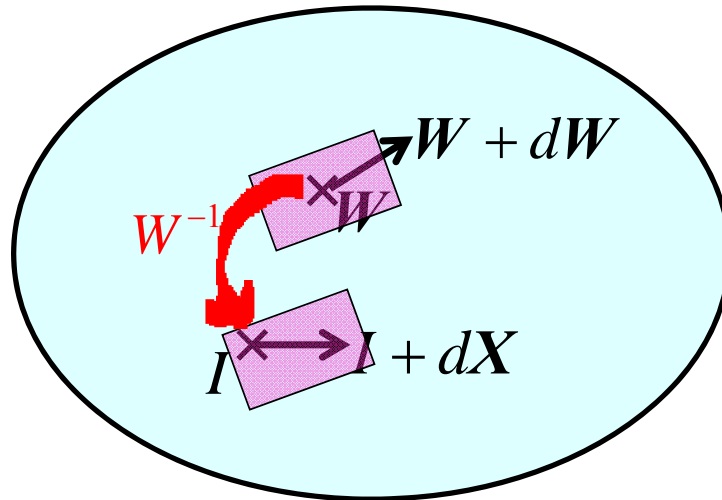
$$\mathbf{W} = \mathbf{A}^{-1}, \quad r(s): \text{ unknown} \quad r = \prod r_i$$

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

# Natural Gradient

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

# Space of Matrices : Lie group



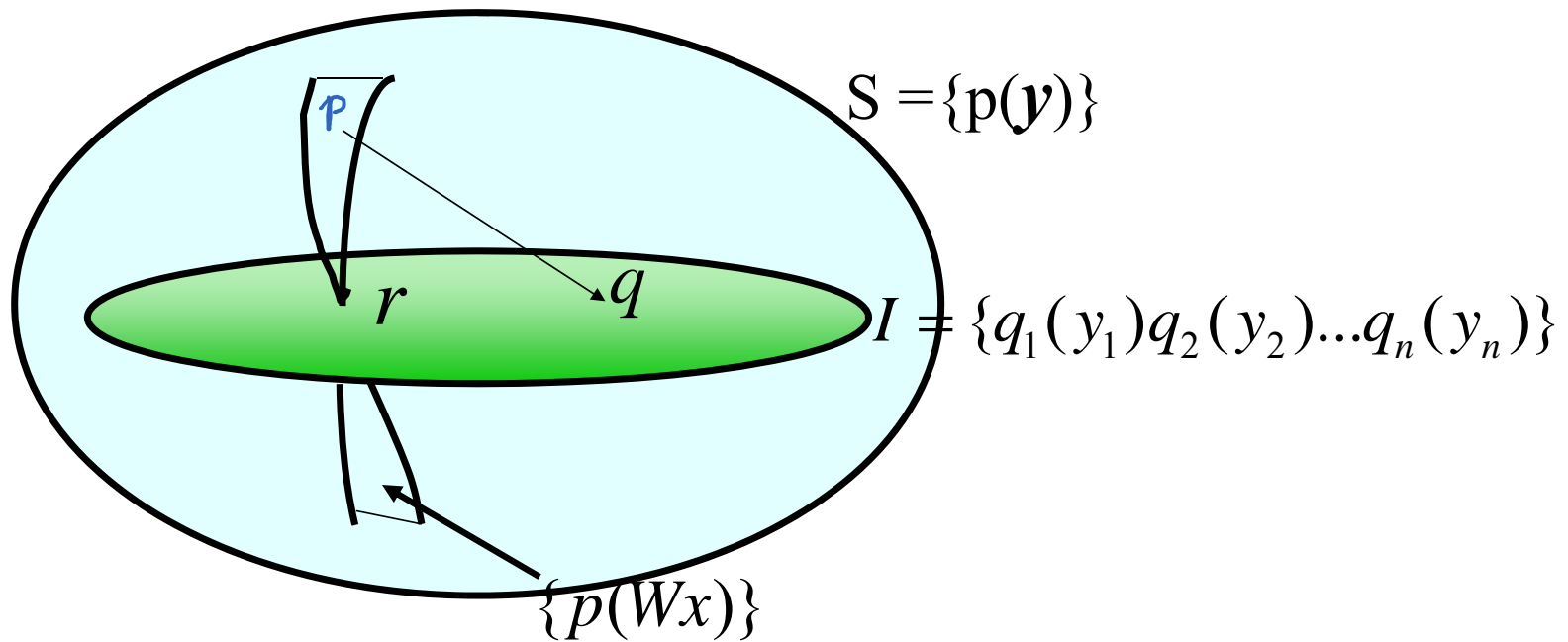
$$dX = dW W^{-1}$$

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

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

$dX$  : non-holonomic basis

# Information Geometry of ICA



**natural gradient**  
**estimating function**  
**stability, efficiency**

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

$r(\mathbf{y})$

## Estimating Functions

$$\Delta W = -\eta F(y, W)$$

$$E_{W, r} [F(y, W')] \begin{cases} = 0, & W' = W \\ \neq 0, & W' \neq W \end{cases}$$

## estimating equation

$$\sum_t F(y_t, W) = 0 \quad y_t = Wx_t$$

## learning

$$\Delta W_t = -\eta F(W_t x_t)$$



## Admissible class

$$F(\mathbf{y}, \mathbf{W}) = \{I - \varphi(\mathbf{y})\mathbf{y}^T\} \mathbf{W}$$

$$\tilde{F} = \mathbf{R}(\mathbf{W}) F(\mathbf{y}, \mathbf{W}) \quad \{\alpha\varphi(y_i)y_j - \varphi(y_j)y_i\} \mathbf{W}$$

$$\sum \mathbf{R}(\mathbf{W}) F(\mathbf{W}\mathbf{x}_t) = 0: \text{ estimating equation}$$

$$\text{on-line learning} : \Delta \mathbf{W}_t = -\eta \mathbf{R}(\mathbf{W}_t) F(\mathbf{W}_t \mathbf{x}_t)$$

$$\text{canonical estimating function: } \frac{\partial \tilde{F}}{\partial \mathbf{W}} = I$$

# Basis Given: overcomplete case

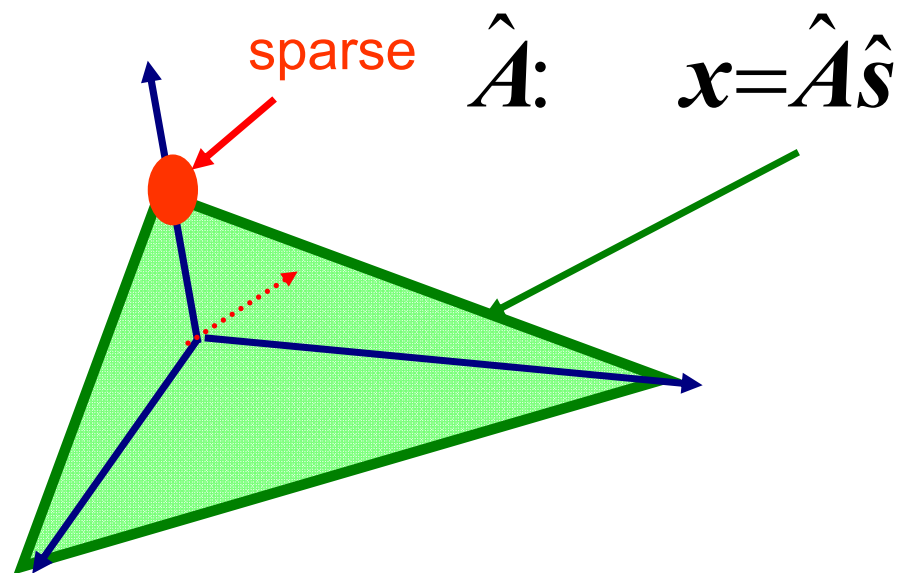
## Sparse Solution

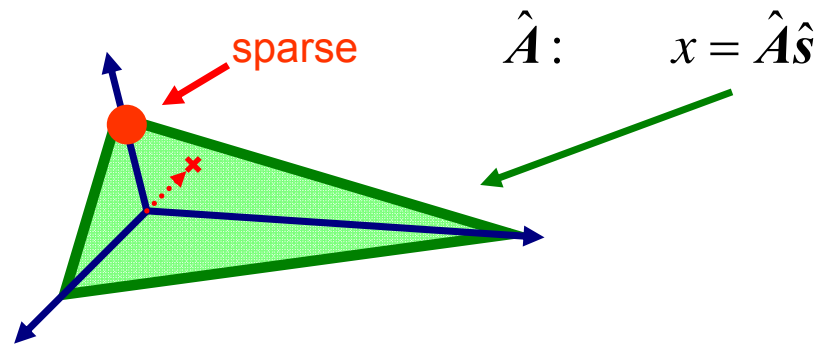
$$\mathbf{x} = A\mathbf{s} = \sum s_i \mathbf{a}_i$$

many solutions

many  $s_i \rightarrow 0$

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



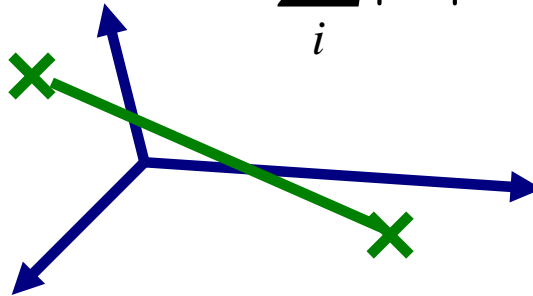


generalized inverse

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

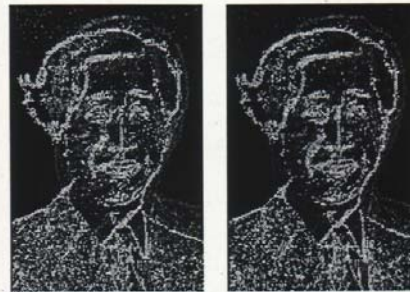
sparse solution

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

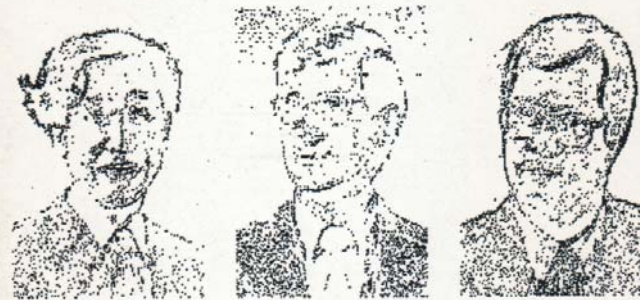




(a) Three binary edge images (reverse images are used in the experiment)



(b) Two edge image mixtures



(c) Reconstructed binary edge images (after reversion)

**Fig. 5:** Example of edge image 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).

# Minkovskian Gradient

*Sparse Signal Analysis*

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

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

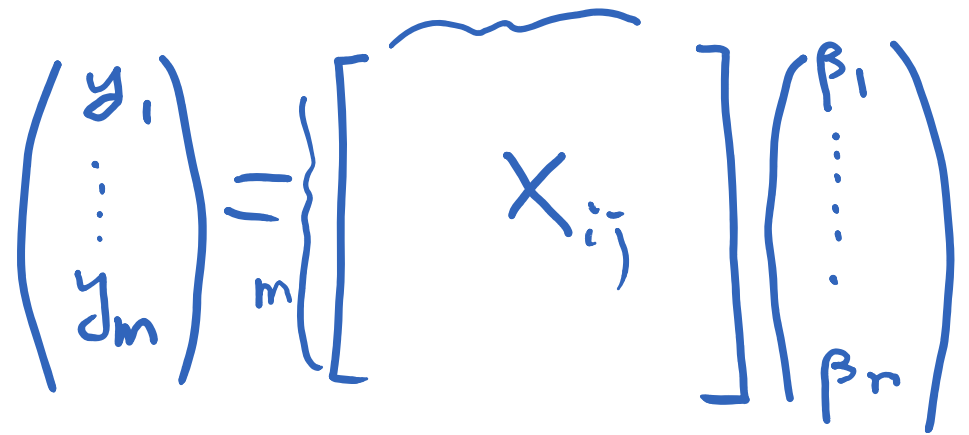
$$F(\boldsymbol{\beta}) = \frac{1}{p} \sum |\beta_i|^p ; \quad p = 2, p = 1, p = 1/2$$

## Optimization under Sparsity Condition:

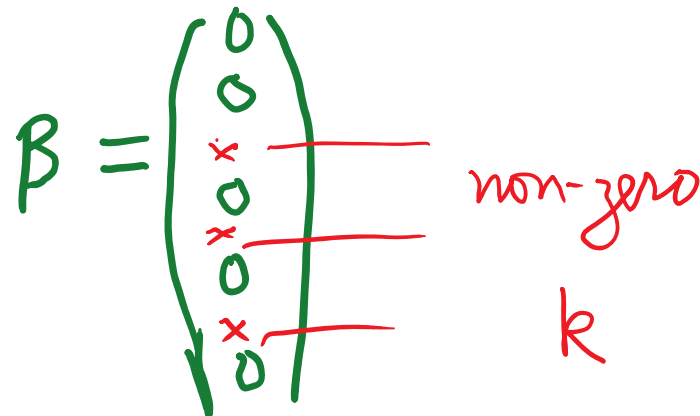
$$y = X\beta \quad (+n)$$

$\beta$ :  $k$ -sparse

$$m > 2k \log n$$



A handwritten matrix equation in blue ink:  $\begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} = \begin{bmatrix} & & \\ & X_{ij} & \\ & & \end{bmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_n \end{pmatrix}$ . The matrix  $X$  is enclosed in large square brackets, and the vector  $y$  is enclosed in large parentheses. A brace under the  $m$  in the vector  $y$  indicates its dimension. A bracket above the  $X$  matrix indicates its dimensions.



A handwritten representation of a sparse vector  $\beta$  in green ink:  $\beta = \begin{pmatrix} 0 \\ 0 \\ 0 \\ x \\ 0 \\ x \\ 0 \\ x \\ 0 \end{pmatrix}$ . Red lines connect the 'x' entries to the text "non-zero" and "k" written in red ink to the right.

# Linear regression

$$y_t = \sum_{i=1}^k x_i^t \beta_i + n_t, \quad t = 1, 2, \dots, N$$

$$y = X\beta + n \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

**Overcomplete:  $N < k$**

**under-determined  $\rightarrow$  infinitely many solutions**

**Sparsity constraint solves the problem**

# Maximum likelihood estimator

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \sum_t \psi(y_t - \boldsymbol{\beta} \cdot \mathbf{x}_t), \quad \psi(u) = \frac{1}{2}u^2$$

$$\sum_t \psi'(y_t - \boldsymbol{\beta} \cdot \mathbf{x}_t) \mathbf{x}_t = 0$$

**Euclidean case (Gaussian noise):**

$$G = \nabla \nabla \psi = X^T X, \quad G \hat{\boldsymbol{\beta}} = X^T \mathbf{y}$$

$$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y} = X^\dagger \mathbf{y} \quad \text{Not sparse}$$



# Sparse Solution

$$\min \psi(\boldsymbol{\beta}) = \min D[\boldsymbol{\beta} : \boldsymbol{\beta}^*] = \psi(\boldsymbol{\beta}) - \psi(\boldsymbol{\beta}^*) - \nabla \psi(\boldsymbol{\beta}^*)(\boldsymbol{\beta} - \boldsymbol{\beta}^*)$$

$$\nabla \psi(\boldsymbol{\beta}^*) = 0$$

penalty  $F_p(\boldsymbol{\beta}) = \sum |\beta_i|^p = c$

$$F_0(\boldsymbol{\beta}) = \#1[\beta_i \neq 0]: \text{ sparsest solution}$$

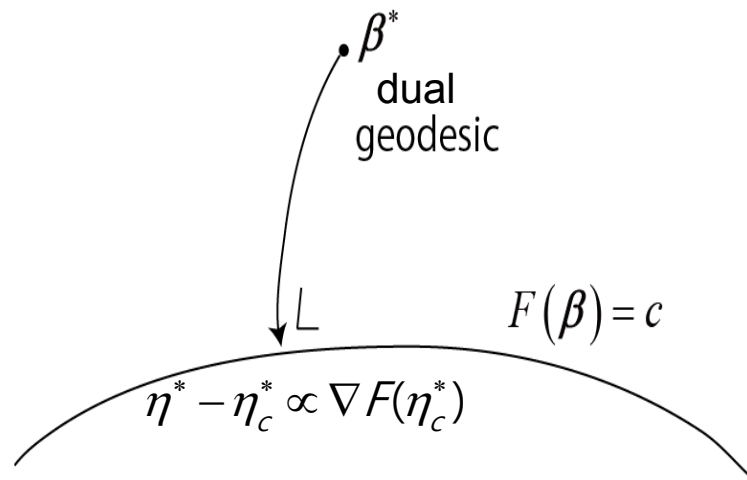
$$F_1(\boldsymbol{\beta}) = \sum |\beta_i|: L_1 \text{ solution}$$

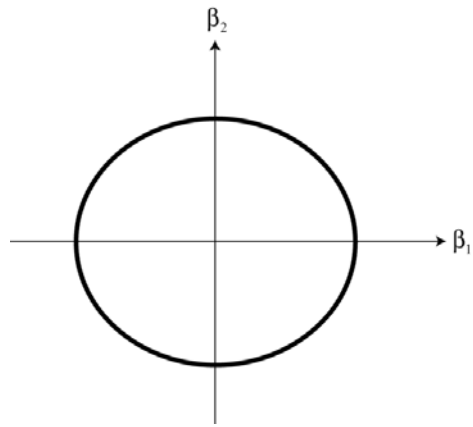
$$F_p(\boldsymbol{\beta}): 0 \leq p \leq 1 \quad \text{Sparse solution: overcomplete case}$$

$$F_2(\boldsymbol{\beta}) = \sum |\beta_i|^2: \text{ generalized inverse solution}$$

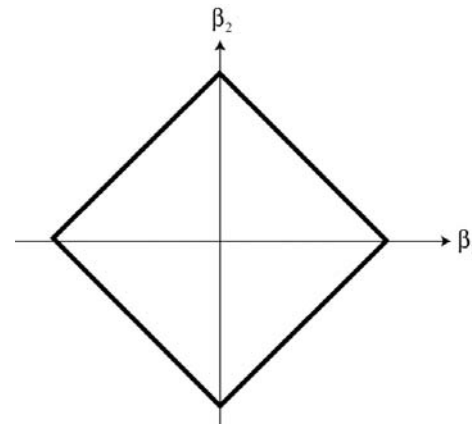
## orthogonal projection, dual projection

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

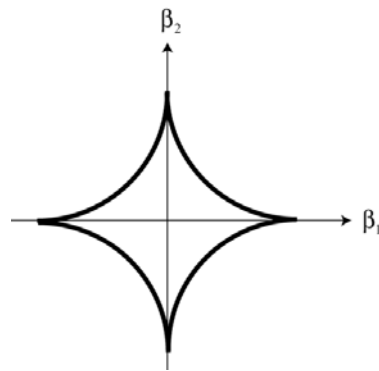




a)  $R_c : n = 2, p > 1$



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



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

**non-convex**

Fig. 1

# L1-constrained optimization

$P_c$  **Problem**       $\min \psi(\beta) \quad \text{under } F(\beta) \leq c$       **LASSO**  
solution  $\beta^*(c) : c = 0 \rightarrow \infty$   
 $\beta_c^* = 0 \rightarrow \beta^*$

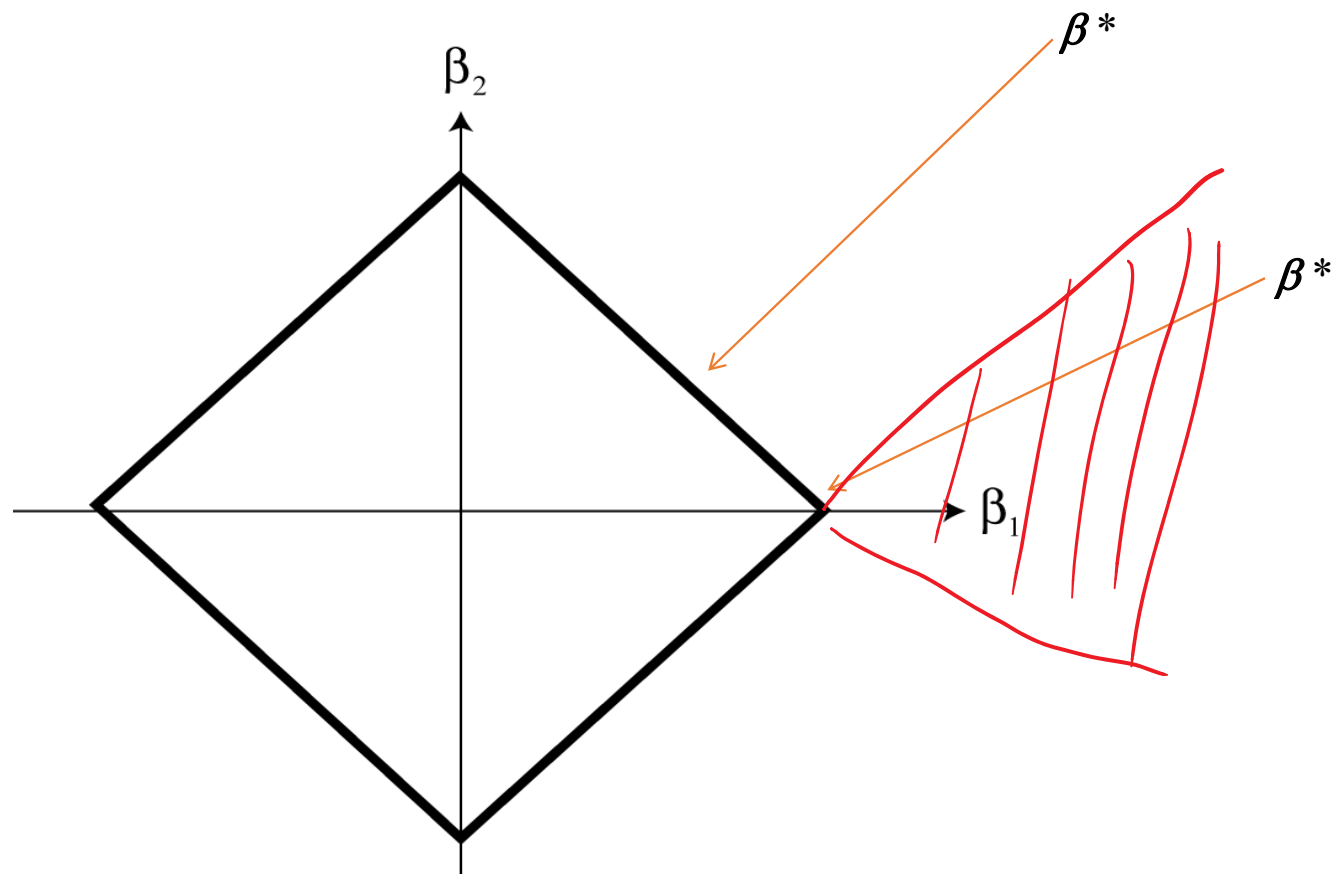
$P_\lambda$  **Problem**       $\min \psi(\beta) + \lambda F(\beta)$       **LARS**  
solution  $\beta^*(\lambda) \quad \lambda = \infty \rightarrow 0$   
 $\beta_\lambda^* = 0 \rightarrow \beta^*$

solutions  $\beta_c^*$  and  $\beta_\lambda^*$  : coincide,  $\lambda = \lambda(c)$ ,  $p \geq 1$

$p < 1$  :  $\lambda = \lambda(c)$  multiple, noncontinuous

stability different

# Projection from $\beta^*$ to $F = c$ (information geometry)



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

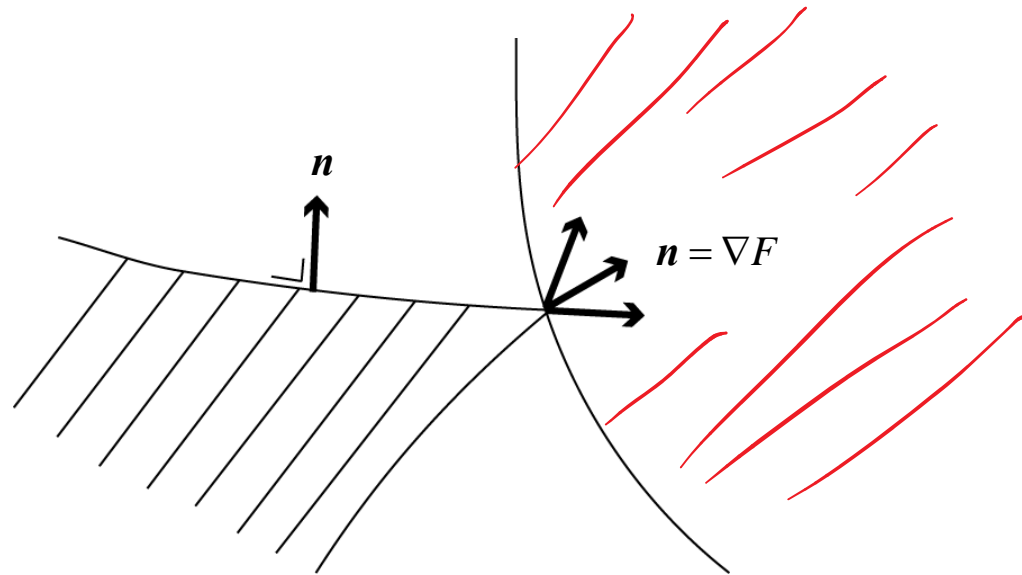


Fig. 5 subgradient

# LASSO path and LARS path (stagewise solution)

$$\min \psi(\boldsymbol{\beta}) : F(\boldsymbol{\beta}) = c$$

$$\min \psi(\boldsymbol{\beta}) + \lambda F(\boldsymbol{\beta})$$

$$\boldsymbol{\beta}^*(c), \boldsymbol{\beta}^*(\lambda)$$

$c \Leftrightarrow \lambda$  correspondence

# Active set and gradient

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

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



## Solution path

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

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

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

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

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

# Solution path in the subspace of the active set

$$\nabla_A \varphi(\boldsymbol{\beta}_\lambda^*) + \lambda \nabla_A F(\boldsymbol{\beta}_\lambda^*) = 0 \quad \nabla_A : \text{active direction}$$

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

turning point  $\mathbf{A} \rightarrow \mathbf{A}'$

# Gradient Descent Method

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

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

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

$$\beta_{t+1} = \beta_t - c \tilde{\nabla} L(\beta_t)$$

# Steepest direction of L

$$G(\mathbf{a}) = \sum g_{ij}(\boldsymbol{\beta}) a_i a_j$$

**Riemannian metric**

$$G(t\mathbf{a}) = |t|^p G(\mathbf{a}) > 0$$

**Minkovskian metric**

$$G(\mathbf{a}) = \sum |a_i|^p. \quad p > 1$$

**Lp-norm**

$$\frac{\delta}{\delta \mathbf{a}} \{ \nabla L \cdot \mathbf{a} - \lambda G(\mathbf{a}) \} = 0$$

$$\delta G(\mathbf{a}) = p(\text{sign } a_i) |a_i|^{p-1}$$

$$f_i = \frac{\partial}{\partial \beta_i} F(\boldsymbol{\beta})$$

$$a_i = c(\text{sgn } f_i) |f_i|^{\frac{1}{p-1}}$$

$$G(\mathbf{a}) = \sum g_{ij} a_i a_j, \quad \tilde{\nabla} F = G^{-1} \nabla f \quad \text{Natural gradient}$$

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

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

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

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

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

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

$$\beta_{t+1} = \beta_t - \eta \tilde{\nabla} F \quad \text{LASSO}$$

**Try for various p, p->1**  
**Try for various noise function**  
**LASSO and flat geometry**

# Extended LARS ( $p = 1$ ) and Minkovskian gradient

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

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

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

$$p = 1^+$$

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

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



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

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

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

$$\beta_{t+1} = \beta_t - \eta \tilde{\nabla} F$$

**LARS**

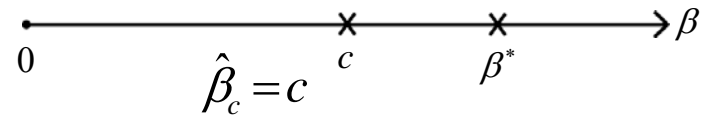
# L1/2 constraint: non-convex optimization

## $\lambda$ -trajectory and-trajectory

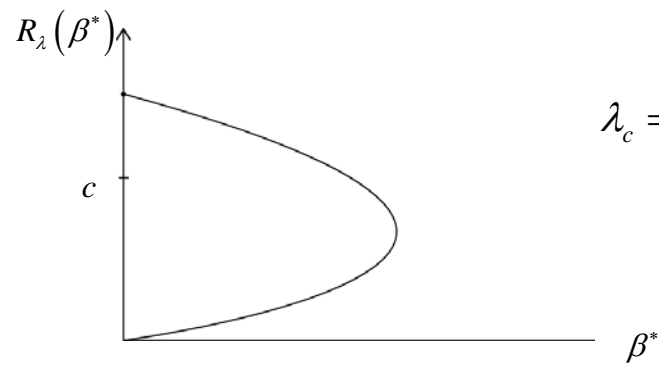
**Ex. 1-dim**  $\varphi(\beta) = \frac{1}{2}(\beta - \beta^*)^2$

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

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

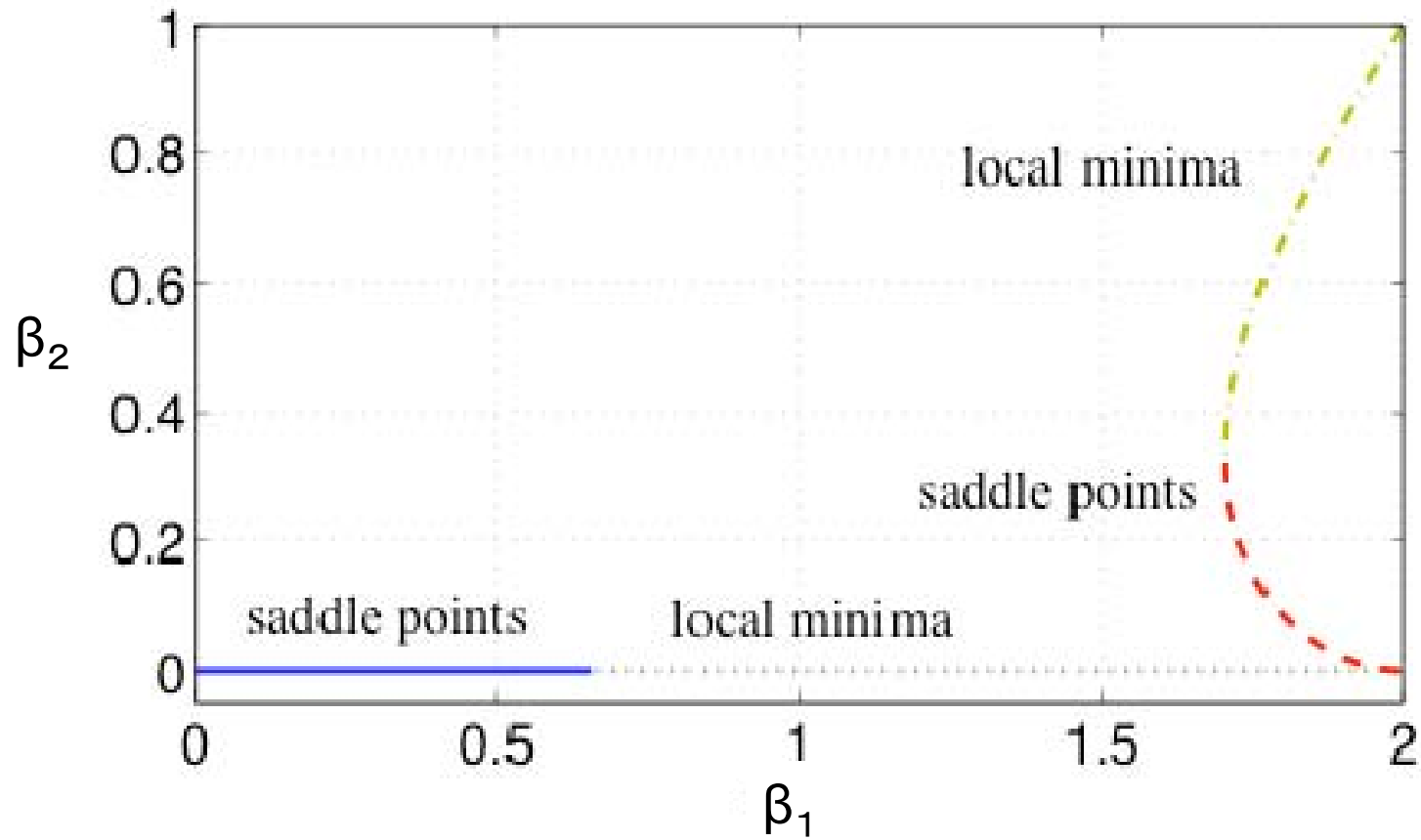


$$P_\lambda : \nabla f_\lambda = 0 \quad \beta - \beta^* + \frac{\lambda}{\sqrt{\beta}} = 0 \quad \hat{\beta} = R_\lambda(\beta^*) \quad : \text{Xu Zongben's operator}$$



$$\lambda_c = \sqrt{c}(\beta^* - c)$$

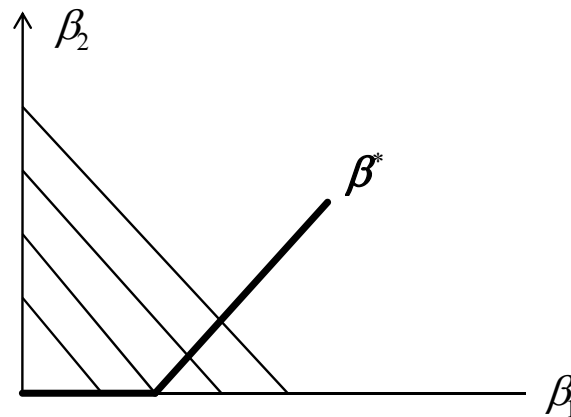
# An Example of the greedy path



## LASSO and LARS :

$p > 1$  :  $\beta_c^*$  is non-sparse

$p = 1$  : sparse

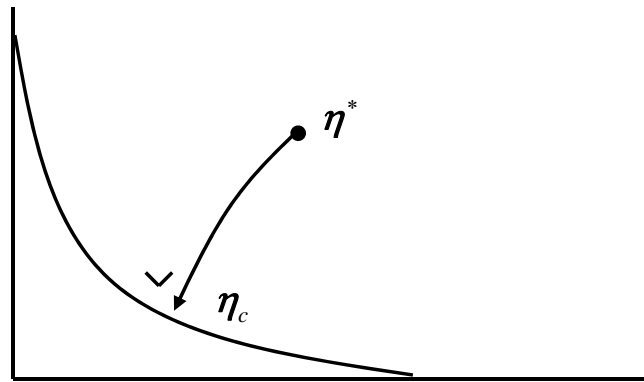


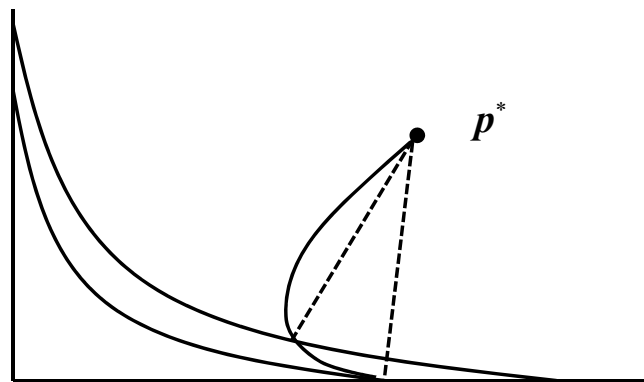
Minkovskian gradient

$$D[\boldsymbol{\beta}_c : \boldsymbol{\beta}^*] - \lambda F(\boldsymbol{\beta})$$

$$\nabla \varphi(\boldsymbol{\beta}_c) - \nabla \varphi(\boldsymbol{\beta}^*) - \lambda_c \nabla F(\boldsymbol{\beta}_c) = 0$$

$$\boldsymbol{\eta}_c - \boldsymbol{\eta}^* = \lambda_c \nabla F(\boldsymbol{\beta}_c)$$





$$\nabla \nabla \varphi(\beta_c) \cdot \dot{\beta}_c - \lambda_c \nabla \nabla F(\beta_c) \cdot \dot{\beta}_c = \dot{\lambda}_c \nabla F(\beta_c)$$

$G$

$H$

$$\dot{\beta}_c = (G - \lambda_c H)^{-1} \nabla F(\beta_c)$$

**Solution Path** :  $\lambda \leftrightarrow c$

**not continuous, not-monotone  
jump**

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

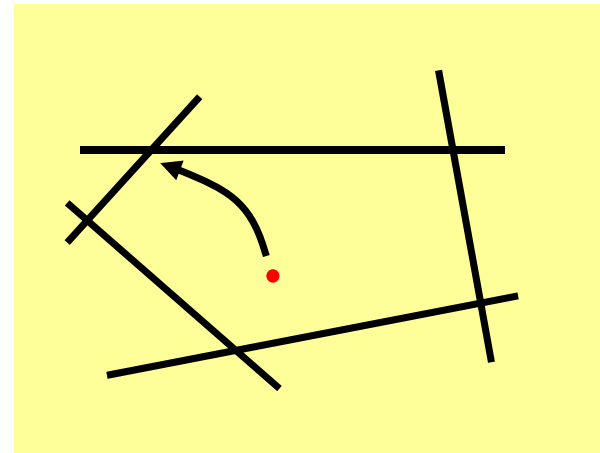


# Linear Programming

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

$$\max \sum c_i x_i$$

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



# Convex Cone Programming

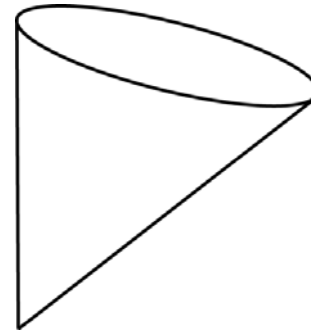
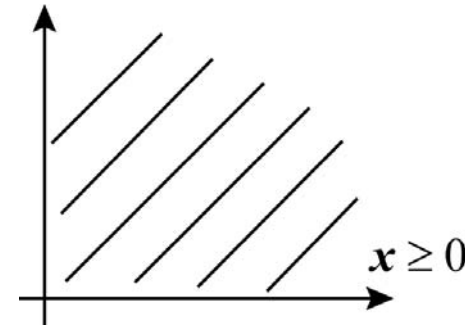
$P$  : positive semi-definite matrix

convex potential function

dual geodesic approach

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

Support vector machine



# Convex Programming

## — Inner Method

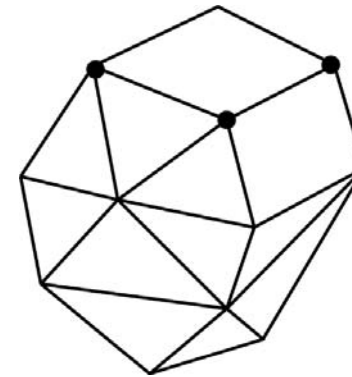
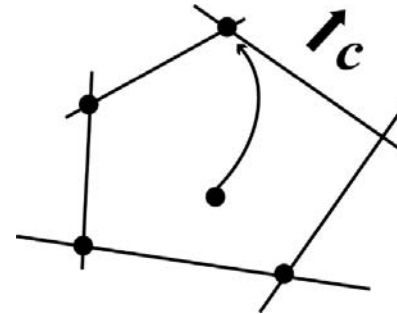
$$LP : Ax \geq b$$

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

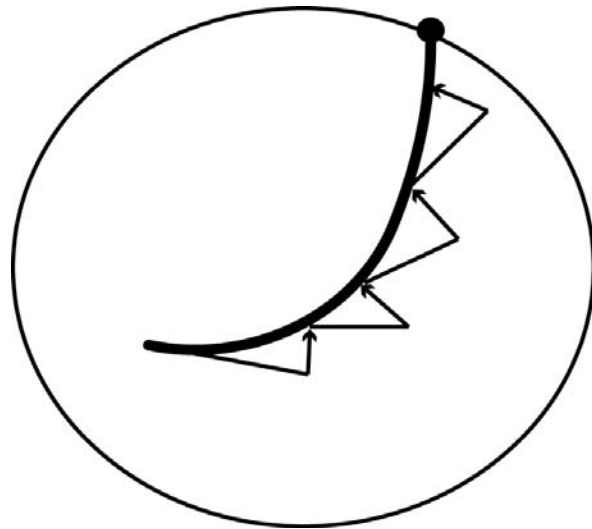
**Barrier function**  $\psi(x) = \sum \log(\sum A_{ij}x_j - b_i)$

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

**Simplex method ; inner method**



# Polynomial-Time Algorithm



curvature : step-size

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

$$\min : tc \cdot \mathbf{x} + \psi(\mathbf{x}) \quad \mathbf{x} = \delta(t)$$

$\nabla^*$  – geodesic

# Integration of evidences:

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

arithmetic mean

geometric mean

harmonic mean

$\alpha$ -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**

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

**$\alpha$ -representation**

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

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

$\alpha$ -mean :  $m_\alpha(p_1(s), p_2(s))$

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

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

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

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

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

# Various Means

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

arithmetic      geometric      harmonic

**Any other mean?**



# $\alpha$ — Family of Distributions

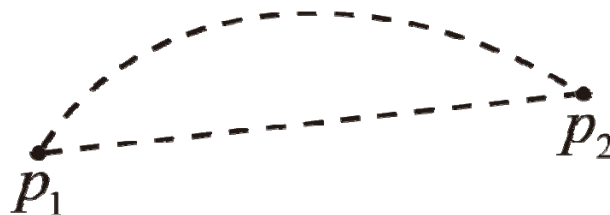
$$\{p_1(s), \dots, p_k(s)\} \quad p(x; \theta) = f_\alpha^{-1} \left\{ \sum \theta_i f_\alpha(p_i(x)) \right\}$$

**mixture family :**

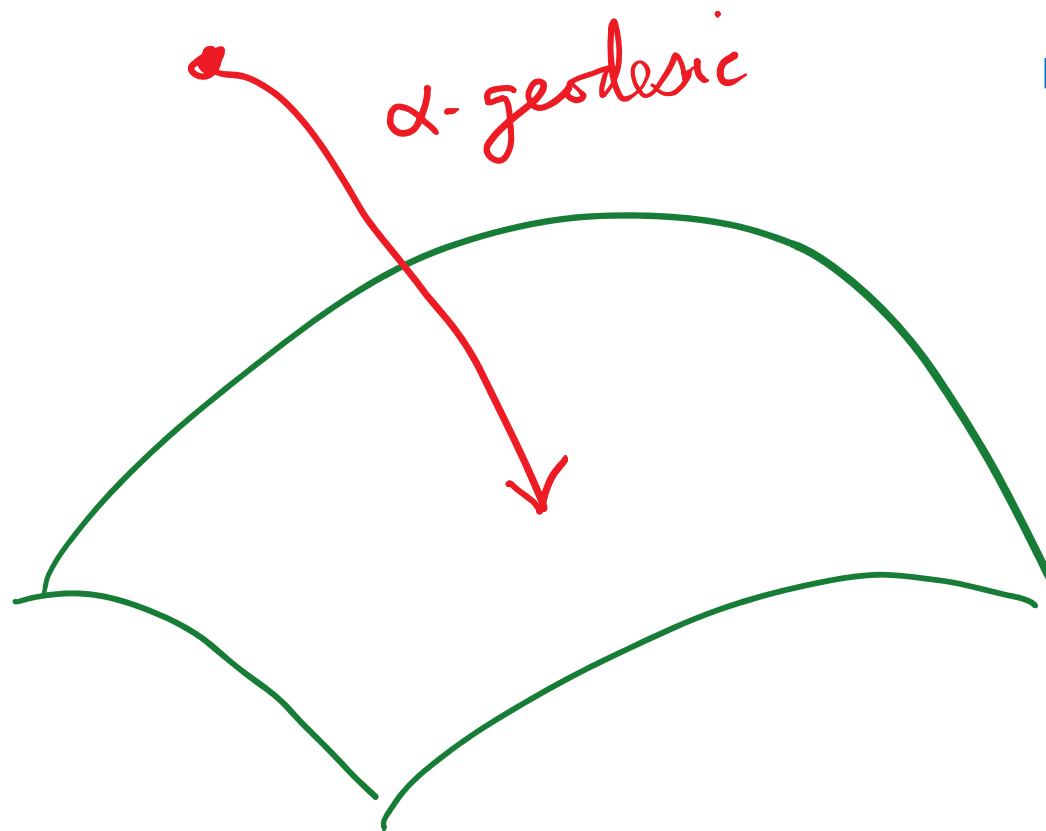
$$p_{mix}(s) = \sum_{i=1}^k t_i p_i(s), \quad \sum t_i = 1$$

**exponential family :**

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



# $\alpha$ -geodesic projection



robust estimator