

IBIS 2016

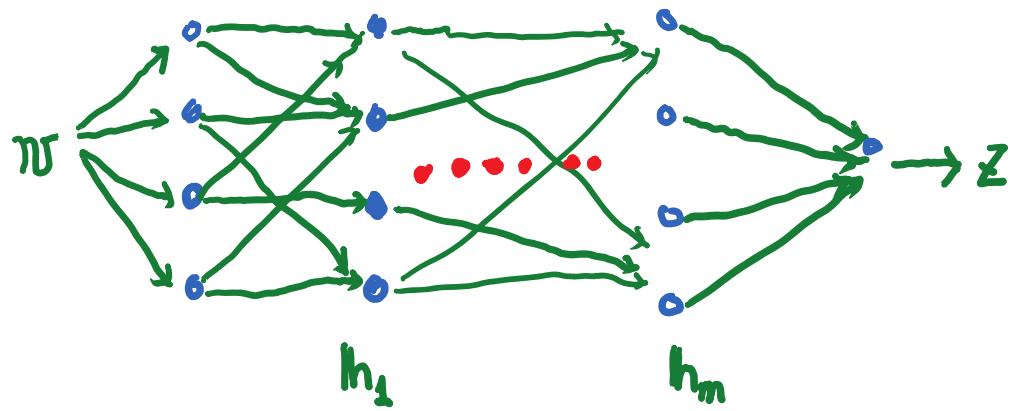
深層学習の基礎：自己組織化と教師付学習

---- 深層学習の理解に向けて：私の試み

甘利俊一 理化学研究所脳科学総合研究センター

深層学習 自己組織化と確率降下教師付学習

RBM: 制約ボルツマン機械、オートエンコーダ、再帰結合回路



tricks !!
ideas !
ドロップアウト
雑音、コンボリューション
bi-directional

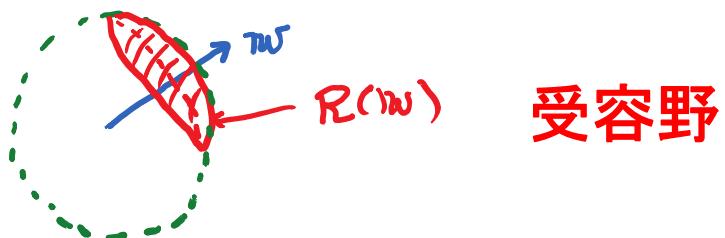
ニューロンのヘブ自己組織化

$$h = f(w \cdot v - \tau) : p(v)$$

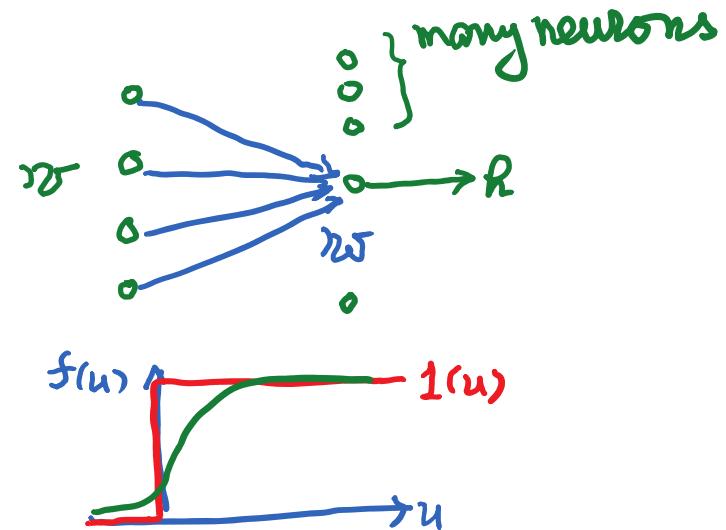
receptive field

$$R(w) = \{v \mid w \cdot v - \tau_0 > 0\}$$

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



受容野



自己組織化學習

$\omega : \text{dynamics of } R(\omega)$

$\eta\omega : P(\eta\omega)$

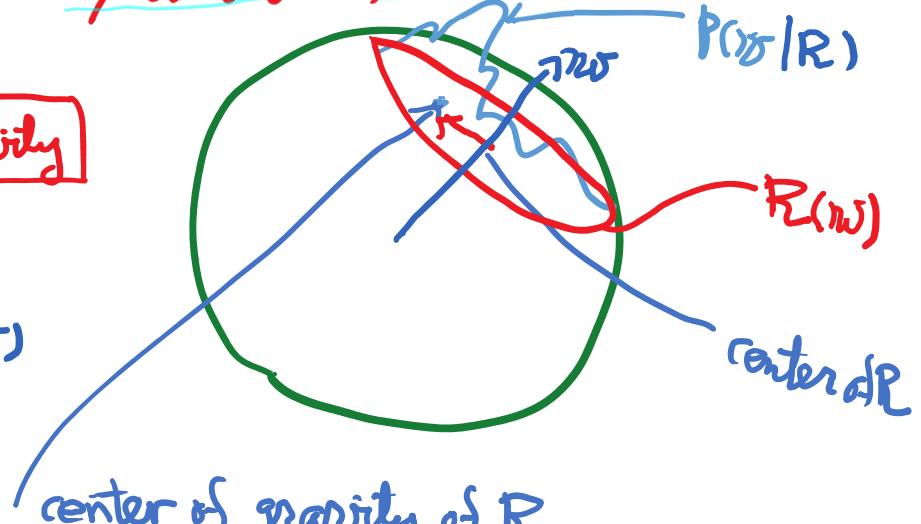
$\dot{\eta}\omega = \langle \mathbf{f}(-\eta\omega + C\eta\omega) \rangle_{P(\eta\omega)}$

$\frac{1}{P_R} \dot{\eta}\omega = -\eta\omega + C \langle \mathbf{f}\eta\omega \rangle_E$

$\text{Prob}\{\eta\omega \in R\} = \langle 1(\eta\omega \cdot \mathbf{v} - v_0) \rangle_{P(\eta\omega)}$

$\boxed{\text{center} \Rightarrow \text{c. of gravity}}$

$E[\mathbf{f}\eta\omega | \eta\omega \in R]$



$P(\eta\omega | R)$

$R(\eta\omega)$

$\text{center of } R$

受容野の重心=中心

学習の収束: 平衡点

中心 = 重心 大きさ τ

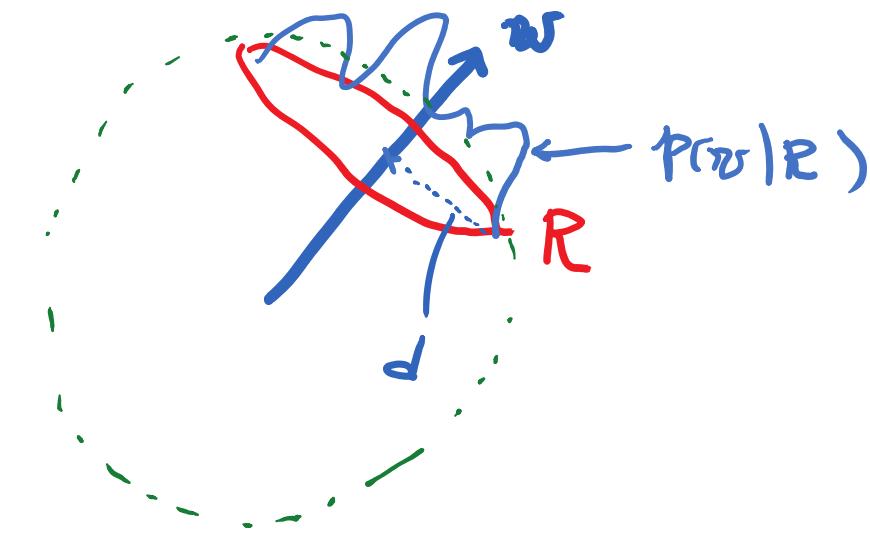
$$\tau_w = \langle \langle h \omega \rangle \rangle_R$$

↑
center \propto C. of gravity

$$\tau_w \cdot \tau_d = \tau$$

d : radius of $R(\omega)$

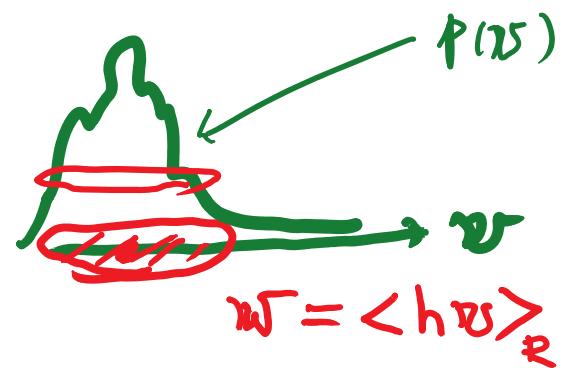
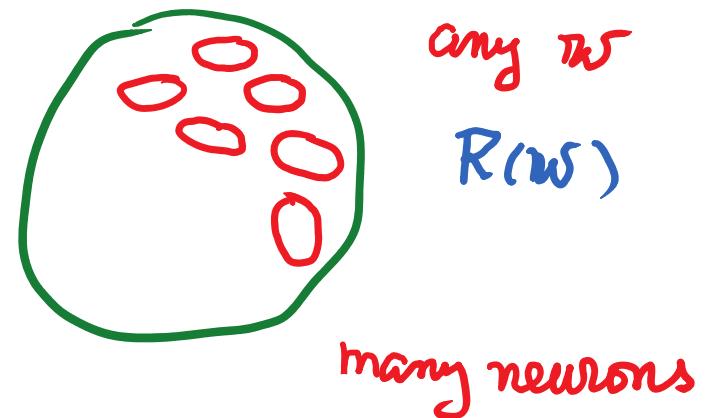
τ : small $\Leftrightarrow d$: large



例: 単純な場合

1. $P(\nu) = c$: uniform

2. $P(\nu)$: single cluster



パターンのクラスター

1. two clusters

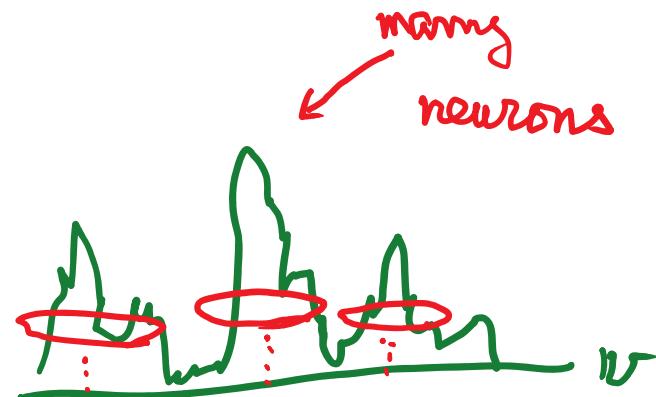
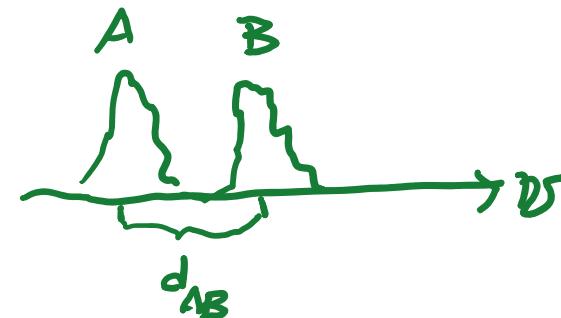
$$d_{AB} > d_0 \quad \mathcal{W} = c \mathcal{V}_A + c \mathcal{V}_B$$

$$d_{AB} < d_0 \quad \mathcal{W} = c_1 \mathcal{V}_A + c_2 \mathcal{V}_B$$

2. many clusters : multi-stable

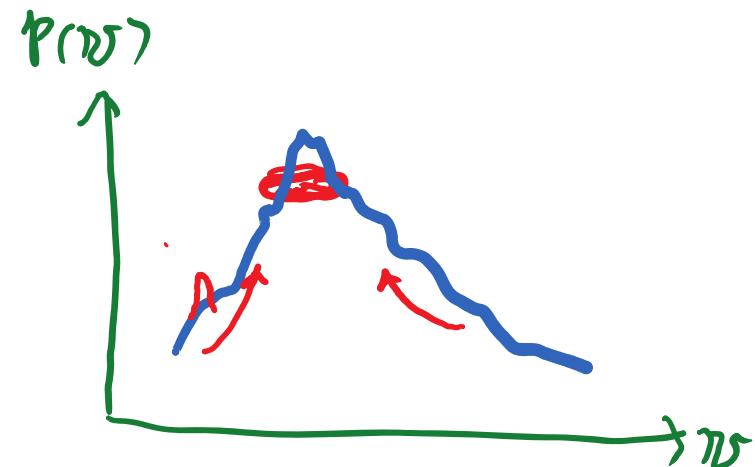
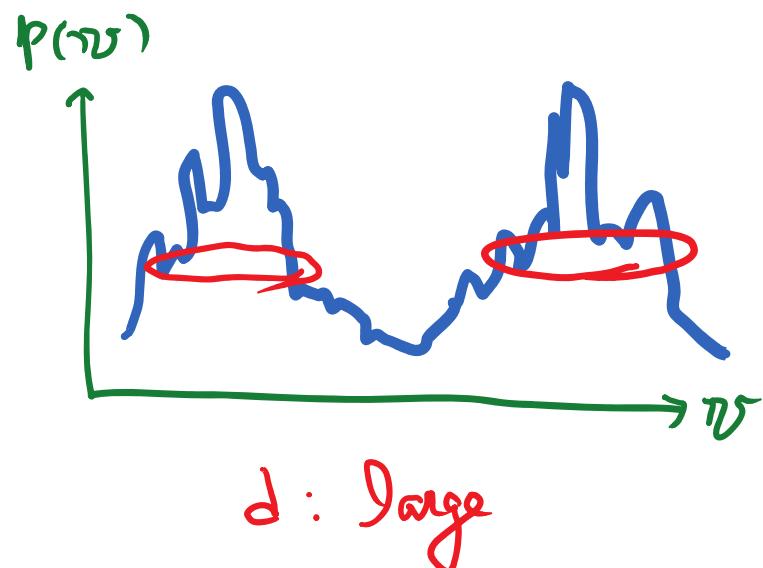
adequate size

center = c. gravity



學習力学

$$\tilde{w} = \langle h w \rangle_k - w \approx \nabla p(w) : d \text{ small}$$



Lyapunov Function

$$F(w, v) = c(v \cdot w - \tau)^+ - \frac{1}{2} \|w\|^2$$

$$\dot{w} = h \nabla_w F$$

$$u^+ = \begin{cases} u, & u > 0, \\ 0, & u \leq 0 \end{cases}$$

$$\dot{F} = \nabla F \cdot \dot{w} = h |\nabla F|^2 > 0$$

$$\langle \dot{F} \rangle_{P(v)} = \int |\nabla F|^2 P(v|R) dv > 0$$

convergence to equilibria

Further Problems

Dimension reduction; PCA, ICA

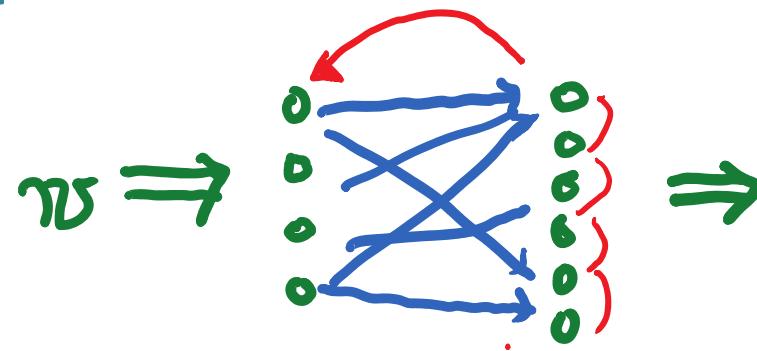
Distributed τ small clusters; large clusters

Mutual interactions among h-neurons neural field

Localized receptive fields

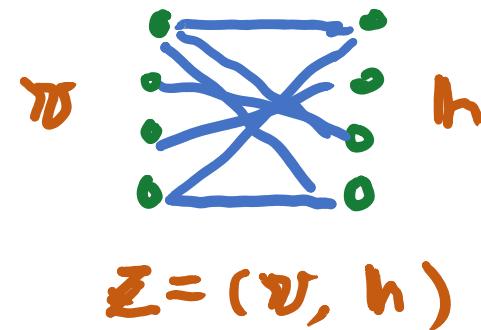
invariance: convolution

feedback



RBM: 制約ボルツマン機械

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

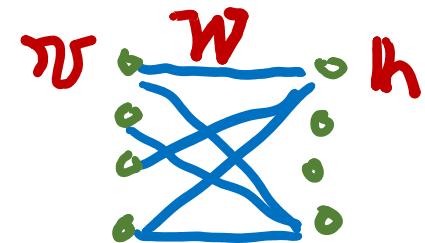


energy machine :

自己組織化

$$: p(v)$$

$\underset{w}{\text{minimize}}$ $\text{KL}[p(v) : \mathcal{E}_v(v; w)]$

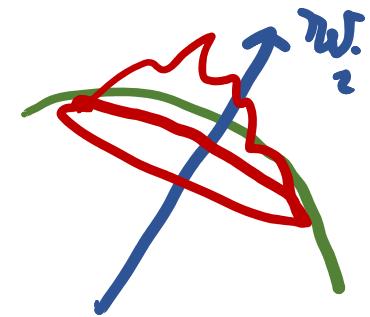


$$= \int p(v) \log \frac{p(v)}{\mathcal{E}_v(v; w)} dv$$

$$\dot{w}_i = \epsilon \left\{ -\langle h_i v \rangle_z + \langle h_i v \rangle_p \right\}$$

center w_i ?

center of gravity of $P(w_i)$



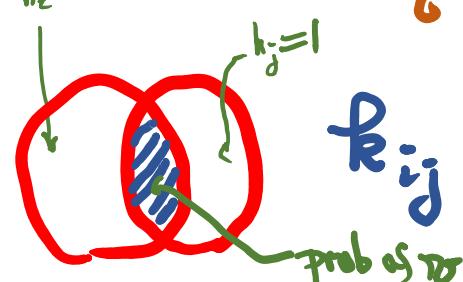
$$\langle h v \rangle_z = \frac{\partial}{\partial w} \Psi(w)$$

ニューロン間の相互作用

$$\langle h_i v \rangle_g = \frac{\partial \Psi(W)}{\partial W} : g(h, v) = \exp\{h^T W v - \Psi(W)\}$$

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

$$\langle h_i v \rangle_i = \sum R_{ij} w_j : \text{interaction}$$



$$R_{ij} = E_i[h_i h_j] \xrightarrow{\text{joint firing prob}} v: \text{analog. Gaussian}$$

受容野は周りのニューロンの受容野に影響される

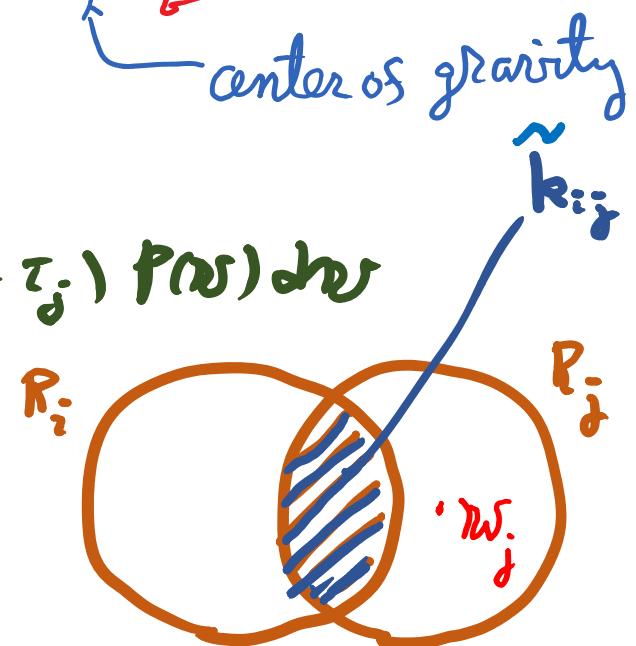
$$w_j = \sum (\hat{k}^{-1})_{j,i} \langle h_i v \rangle_g$$

v : Gaussian

$$\hat{k}_{ij} = \int f(\mathbf{w}_i \cdot \mathbf{v} - \tau_i) f(\mathbf{w}_j \cdot \mathbf{v} - \tau_j) P(\mathbf{v}) d\mathbf{v}$$

$$\langle h_i v \rangle_{g,j} = \sum \tilde{k}_{ij} w_j$$

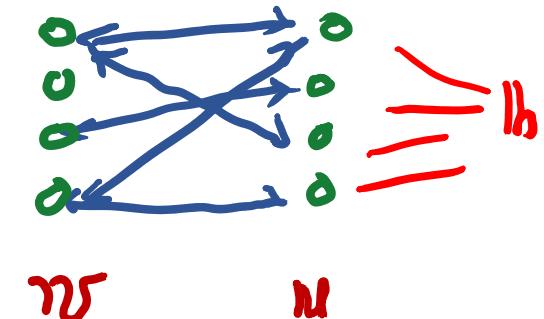
↑
center gravity
↑
interaction



Recurrent Net (Auto-Encoder)

$$\begin{cases} \dot{u} = -u + f(Wv) \\ \dot{v} = -v + f(S^T h) + a \end{cases} \quad \text{multi-stable}$$

pr(a)



$$\dot{w} = -w + c \langle h v^\top \rangle_{\text{pr}(a)}$$

$$\dot{s} = -s + c \langle v h^\top \rangle_{\text{pr}(a)}$$

$$s_{:,j} = w_{:,j} \quad (\text{equilibrium})$$

ガウス型ボルツマン機械

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

$$\langle h_i w_j \rangle_c = \boldsymbol{w}^T \boldsymbol{C} \boldsymbol{h}_j$$

center of gravity

$$\boldsymbol{C} = \int \boldsymbol{w} \boldsymbol{w}^T p(\boldsymbol{w}) d\boldsymbol{w}$$

covariance matrix

$$\langle h_i w_j \rangle_i = (\mathbf{I} - \boldsymbol{W} \boldsymbol{W}^T)^{-1} \boldsymbol{W}$$

interactions of h neurons

平衡状態の解

(唐木田;東大)

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

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 $C = V^T \text{diag}(\lambda_1, \dots, \lambda_n) V$

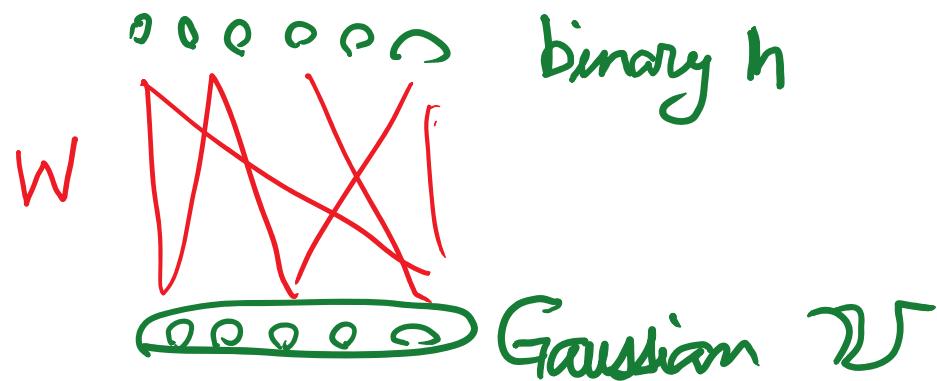
You can choose $m (\leq k)$ eigen values from $\lambda_1, \dots, \lambda_k > 1$

Stable Solution the case of $m = k$

ガウス一ベルヌイ機械

ICA: 独立成分分析

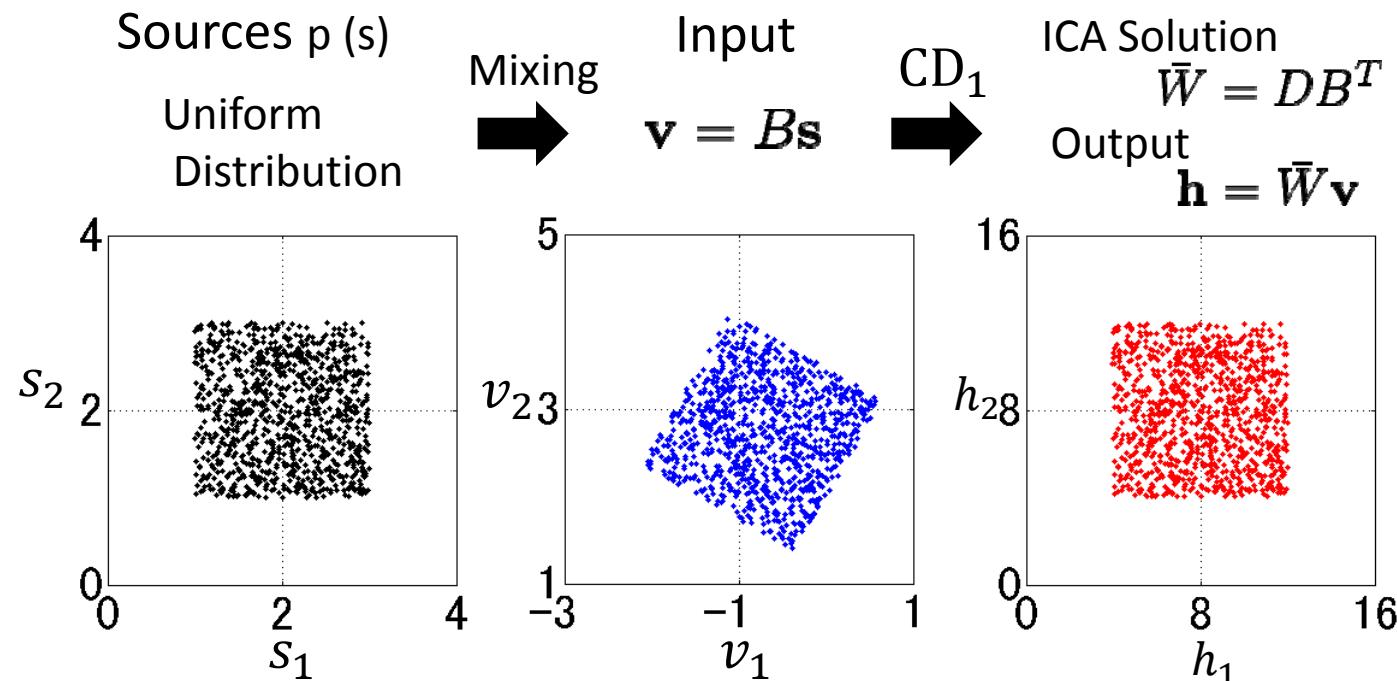
唐木田



$$P(w) : \mathcal{W} = \underbrace{\text{ODS}}_{B^{-1}} \quad \text{独立}$$

シミュレーション例

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



Independent sources are extracted in G-B RBM

何が問題か : 画像情報の構造、生成モデル(確率分布)

Uniform : no structure

Aggregate of clusters : Hebb self-organization

PCA : Gaussian RBM submanifolds : $p(\mathbf{v})$

ICA : Bernoulli-Gaussain sparse

スペース生成のクラスター(低次元)

階層構造と深層学習

invariancy

logical structure

hierarchies of hierarchy

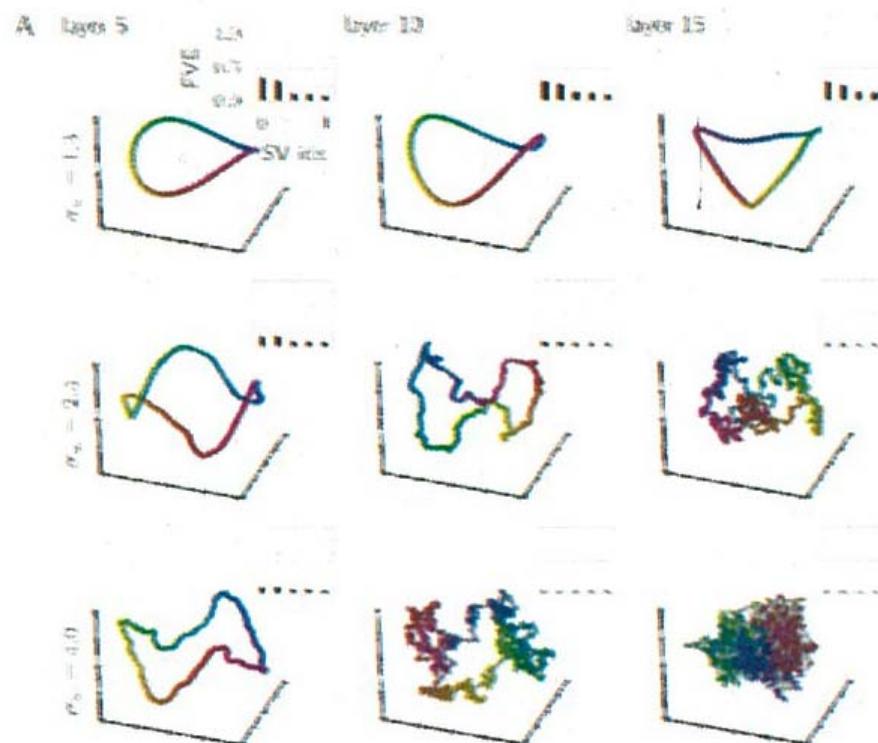
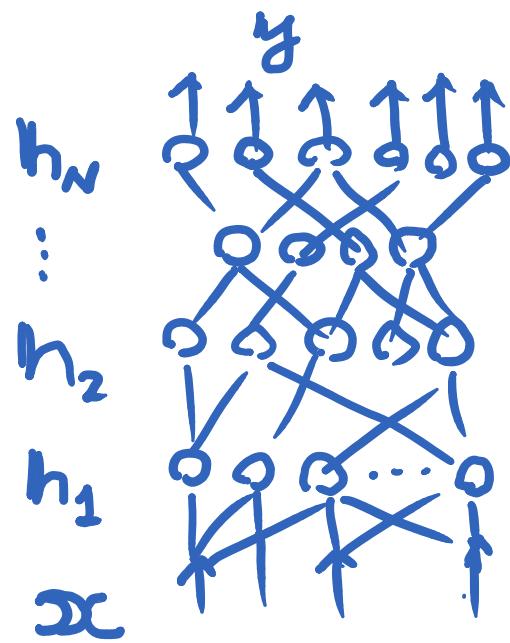
大規模系の特徴

ランダム行列 A の固有値の分布
ほとんどが鞍点（極小解なし!!）

大規模回路

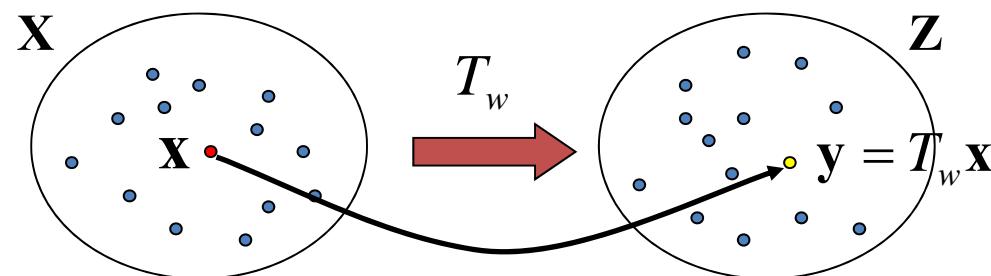
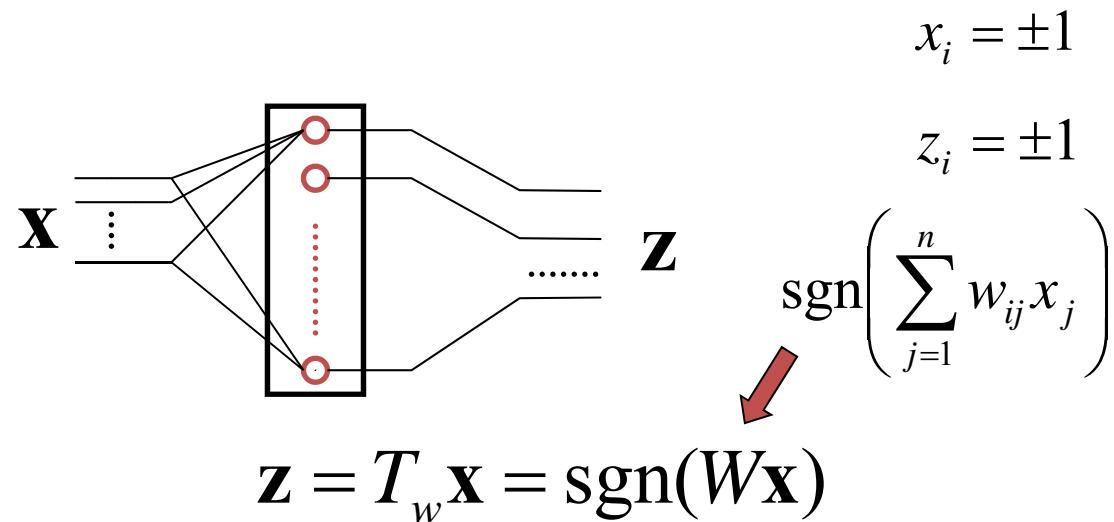
極小解は最小解の付近に集まる

Poole et al (2016)
Deep neural networks

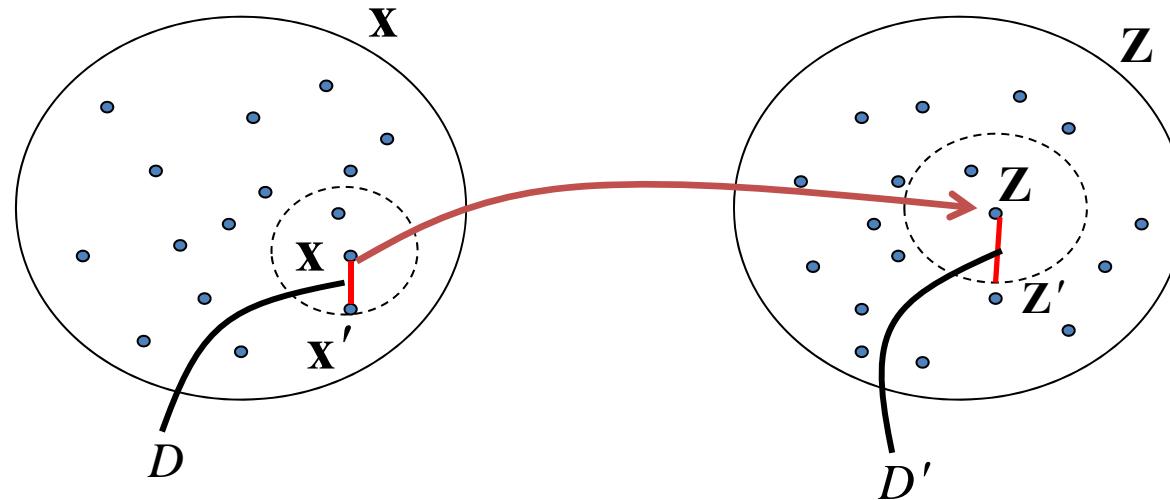


統計神経力学

1-layer network



距離法則 : stability

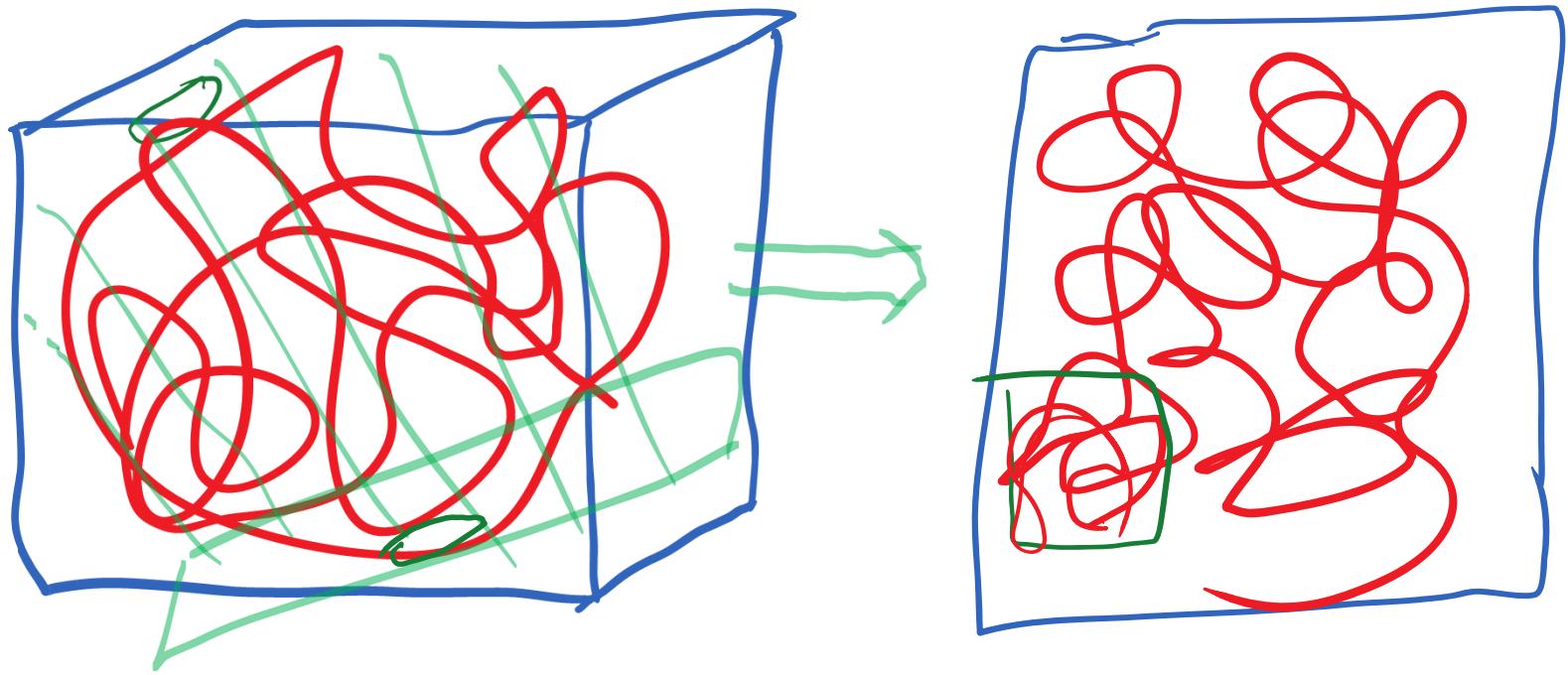


$$D(\mathbf{x}, \mathbf{x}') = \frac{1}{2n} \sum |x_i - x'_i| = D$$

$$\begin{array}{ccc} D(T_w \mathbf{x}, T_w \mathbf{x}') & = & \phi(D) = D' \\ \parallel & \parallel \\ \mathbf{z} & & \mathbf{z}' \end{array}$$

$$\phi(D) = \frac{2}{\pi} \sin^{-1}(c\sqrt{D})$$

次元 小



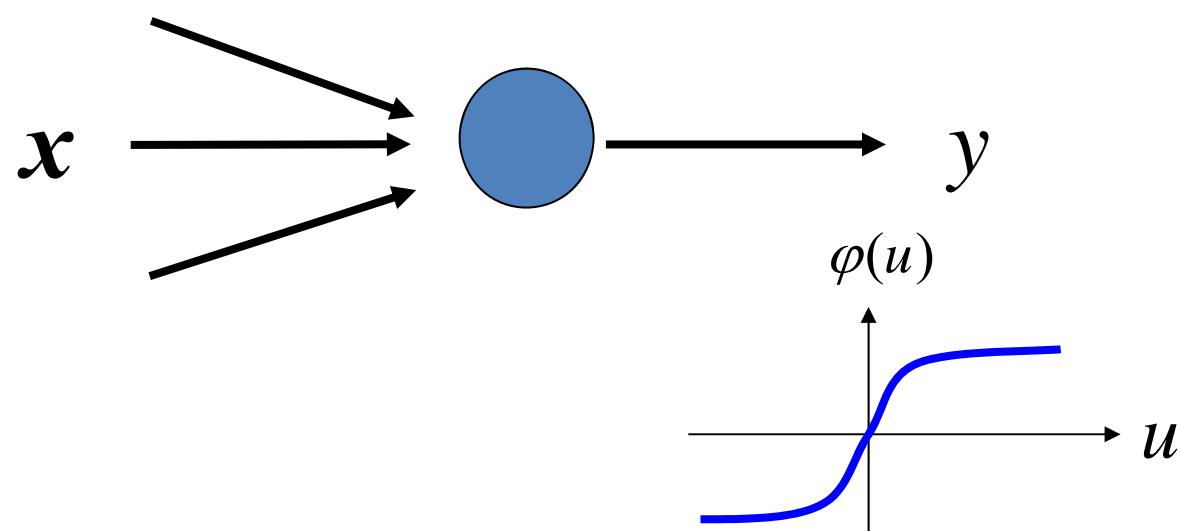
教師付学習：誤差逆伝搬

多層パーセプトロン
特異点!!

自然勾配学習：リーマン空間

数理ニューロン

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

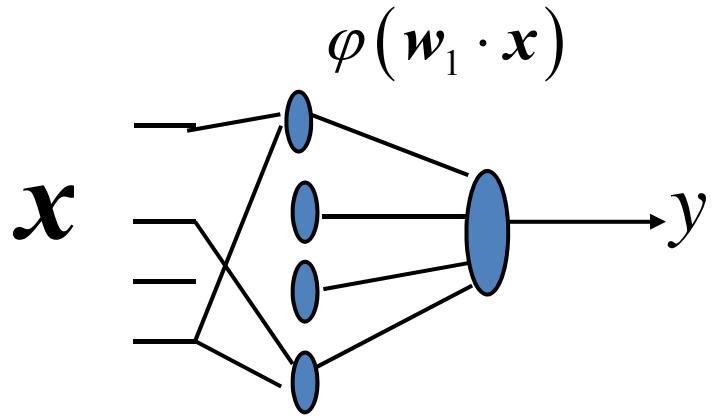


多層パーセプトロン

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

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

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



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

多層パーセプトロン

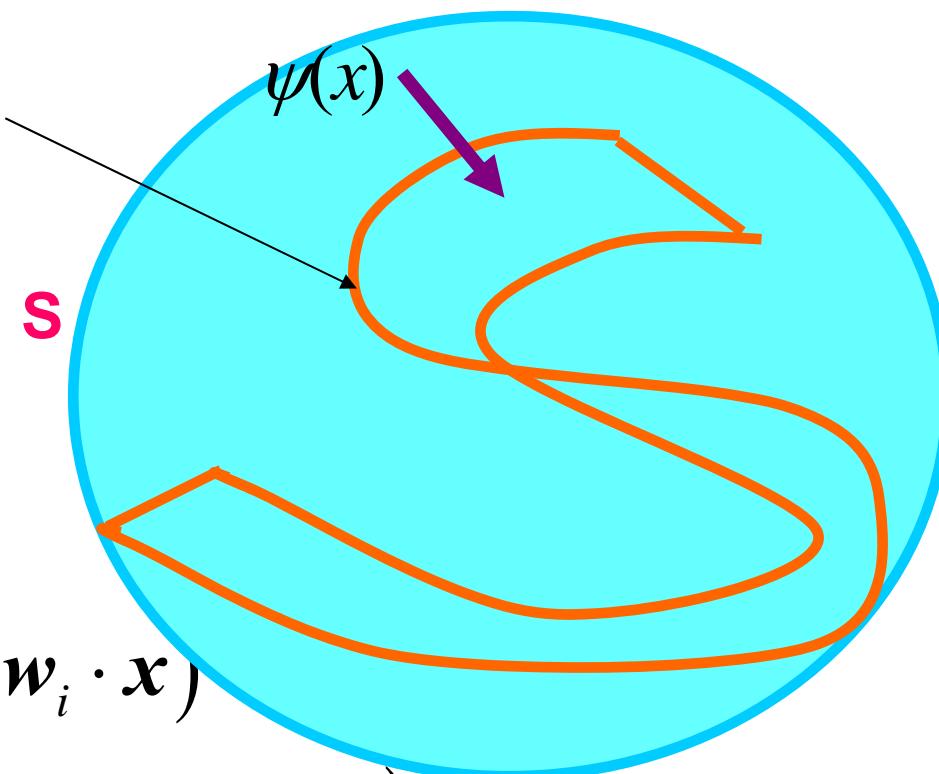
神経多様体

関数の空間 S

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

$$= \sum v_i \varphi(w_i \cdot x)$$

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



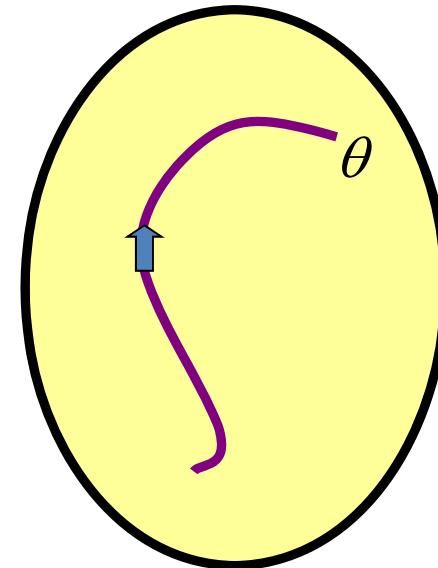
確率降下學習 --- 誤差逆伝搬

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

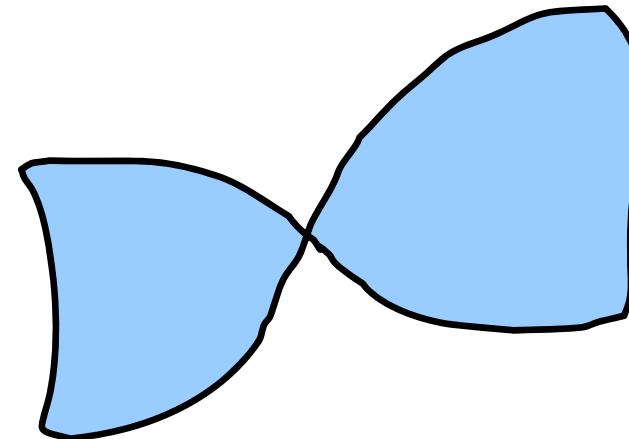
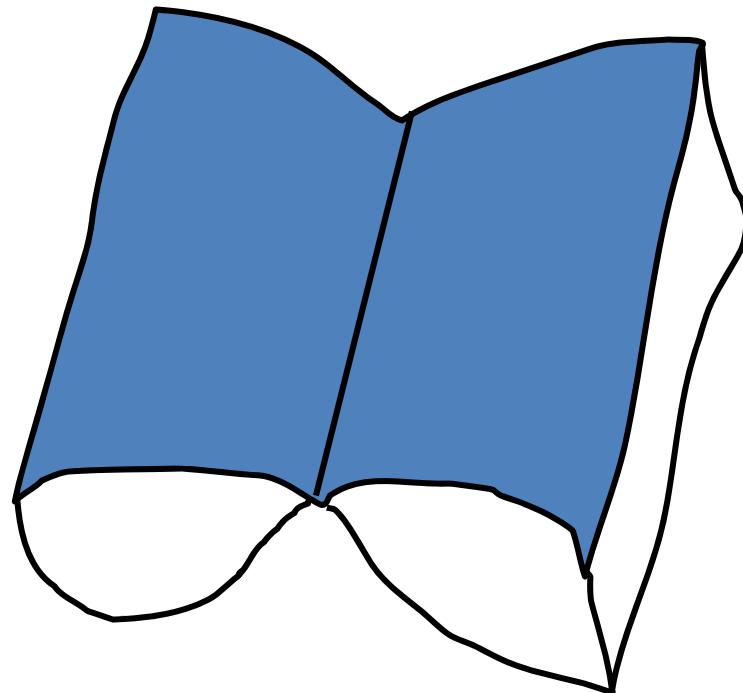
$$\begin{aligned} l(y, \mathbf{x}; \theta) &= \frac{1}{2} |y - f(\mathbf{x}, \theta)|^2 \\ &= -\log p(y, \mathbf{x}; \theta) \end{aligned}$$

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

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

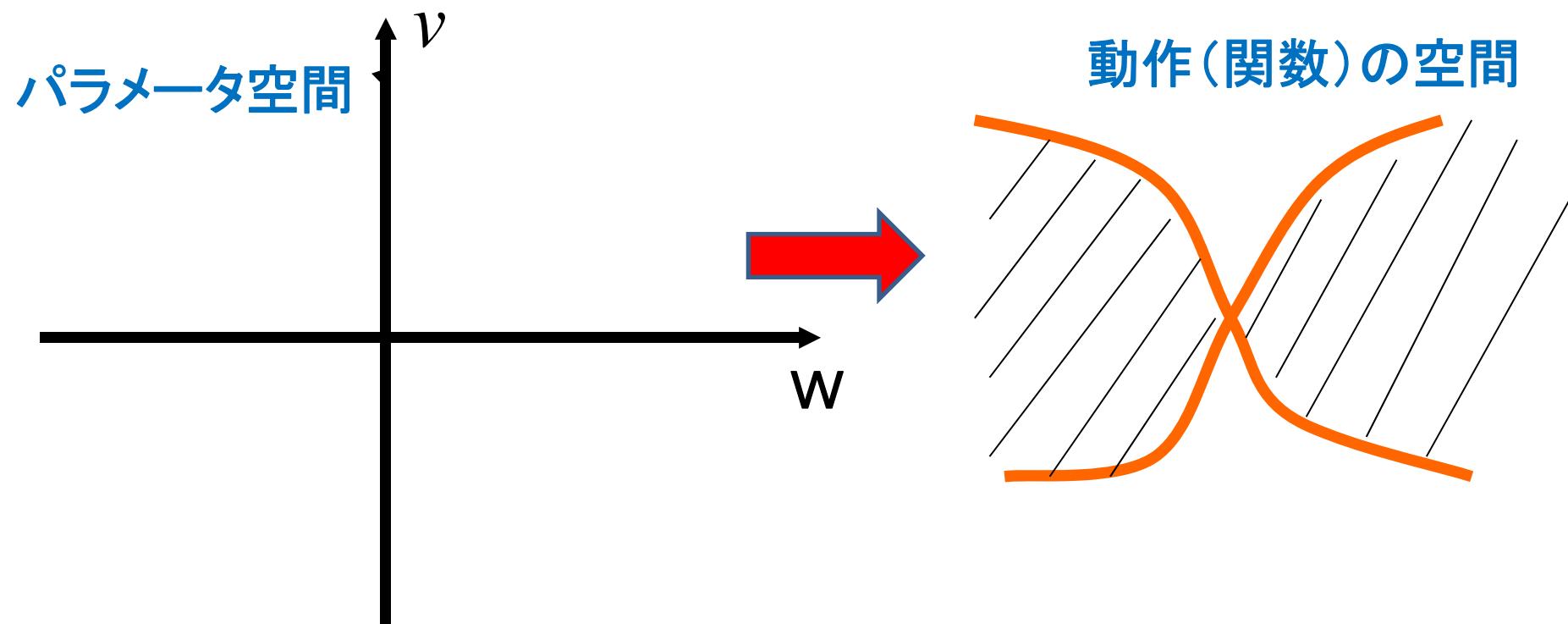


特異点・特異構造



特異モデルの例

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

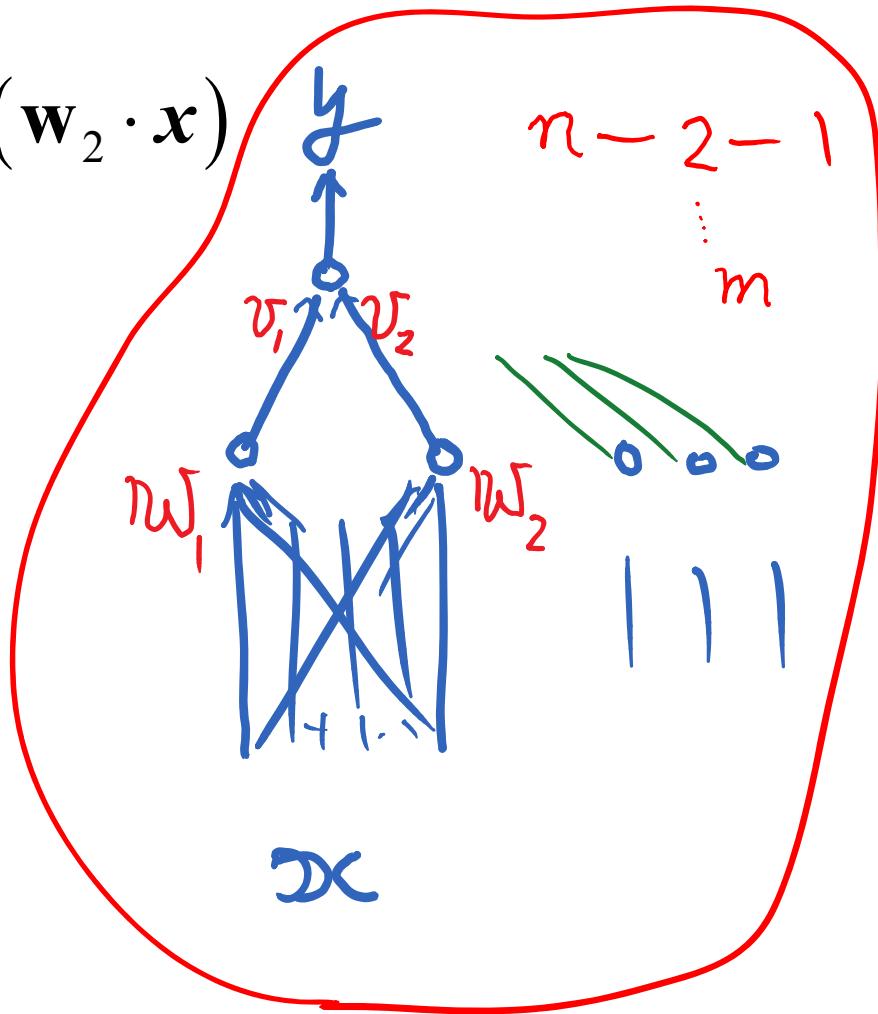


典型モデルの解析

$$f(\mathbf{x}, \theta) = v_1 \varphi(\mathbf{w}_1 \cdot \mathbf{x}) + v_2 \varphi(\mathbf{w}_2 \cdot \mathbf{x})$$

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

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



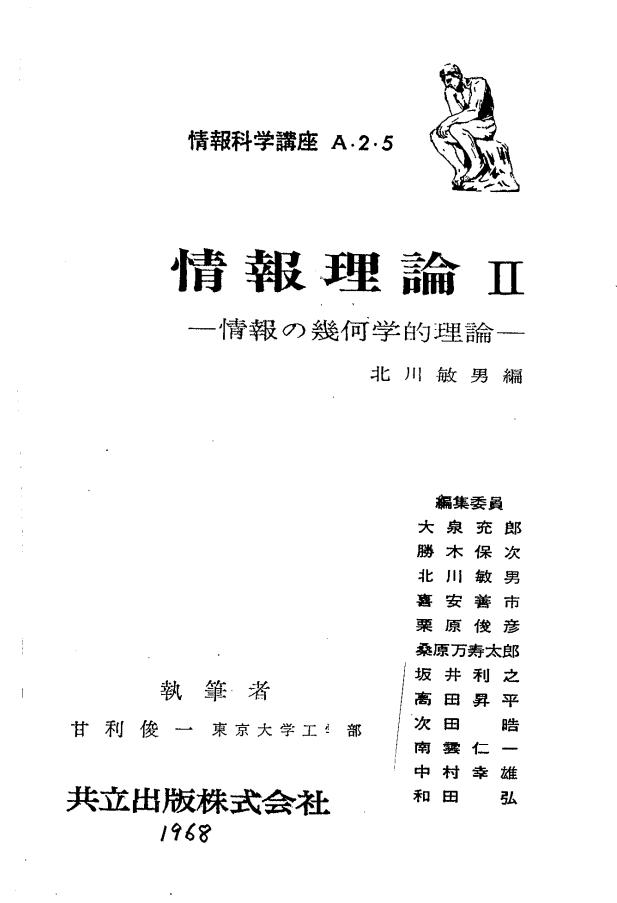
loss function: $l(x, y; \theta) = \frac{1}{2} \{y - f(x, \theta)\}^2$

y : teacher signal : θ_0 **stochastic descent learning**

$$\dot{\theta} = -\eta \left\langle \frac{\partial l(x_t, y_t, \theta_t)}{\partial \theta} \right\rangle$$

backprop : vanilla gradient

First stochastic descent learning of MLP (1967;1968)



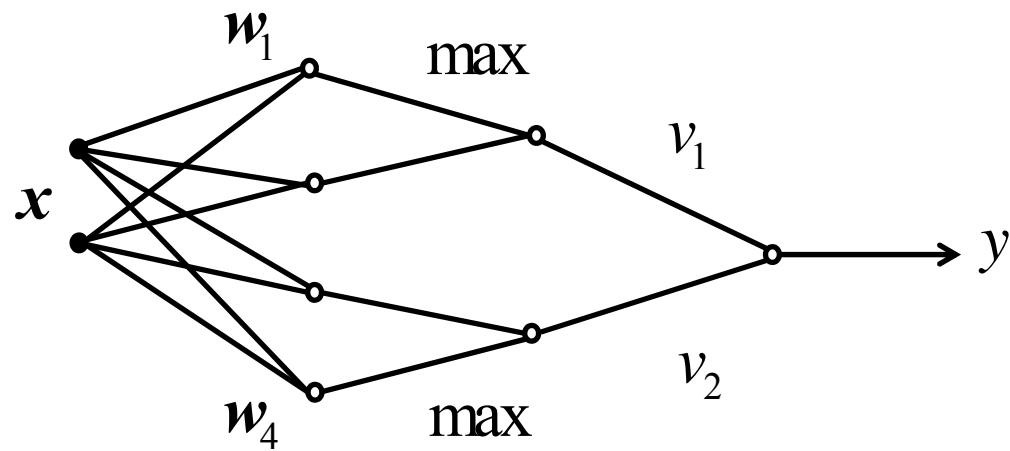
Information Theory II
--Geometrical Theory of Information

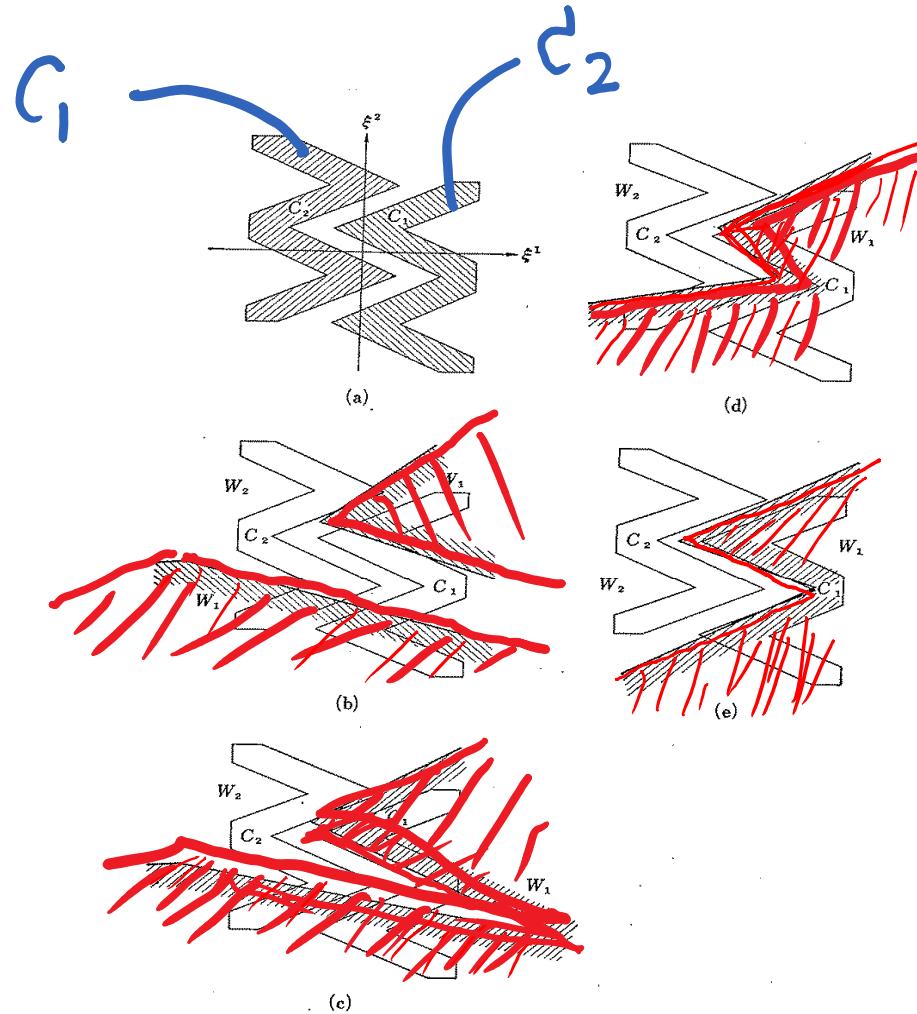
Shun-ichi Amari
University of Tokyo

Kyouritu Press, Tokyo, 1968

$$f(\mathbf{x}, \theta) = v_1 \max \{ \mathbf{w}_1 \cdot \mathbf{x}, \mathbf{w}_2 \cdot \mathbf{x} \} + v_2 \max \{ \mathbf{w}_3 \cdot \mathbf{x}, \mathbf{w}_4 \cdot \mathbf{x} \}$$

$$v_1 = 1; v_2 = -1$$





自然勾配學習法

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

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

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

不变； 最急降下

自然勾配 (Riemannian)

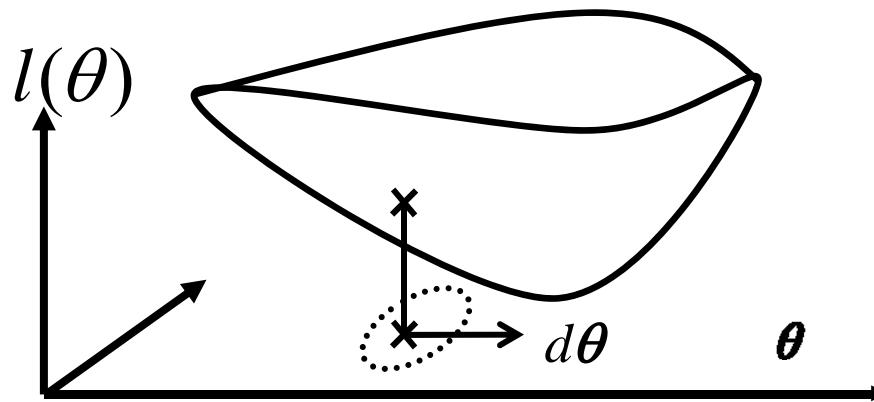
$$\max \quad dl = l(\theta + d\theta) - l(\theta)$$

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

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

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

最急降下---Natural Gradient



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

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

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

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

自然勾配の長所

Steepest descent; invariant Yan Ollivier

Fisher-efficient

Natural gradient is non-vanishing even in multiple layers

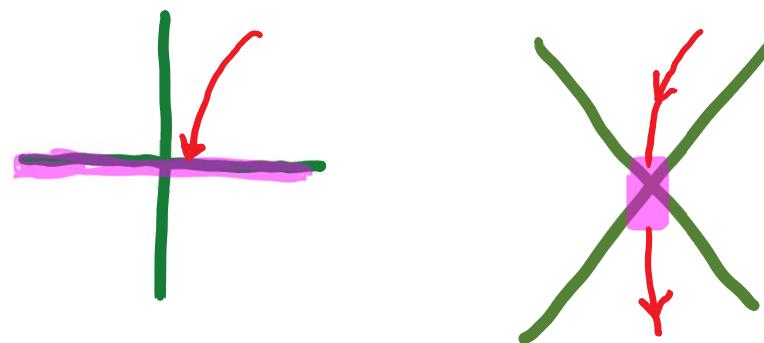
Good at singular regions (avoid plateaus: Milnor attractor)

適応自然勾配

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

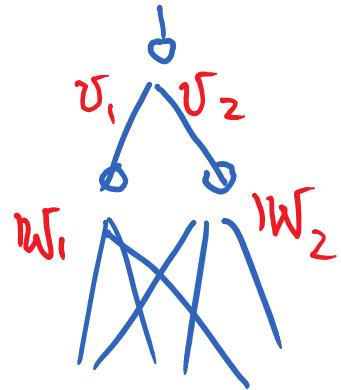
$G^{-1} \rightarrow \infty$, $\nabla \ell \rightarrow 0$ at singularities

$G^{-1} \nabla \ell$



特異領域

$$R(v, w) = \{\theta | w_1 = w_2 = w, v_1 + v_2 = v\}$$



$$\cup \{\theta | v_1 = 0, v_2 = v, w_2 = w\}$$

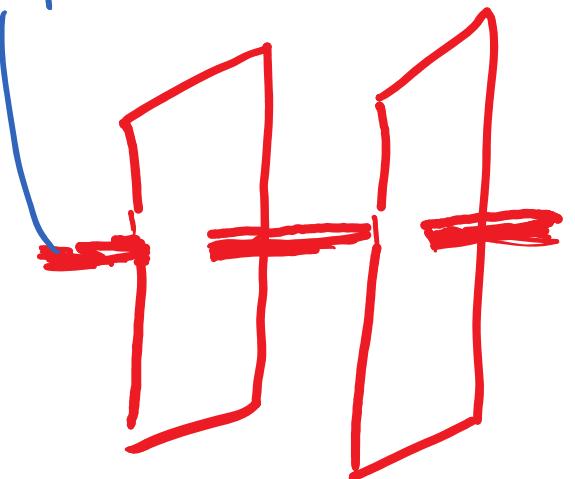
$$w_1 = w_2$$

$$\cup \{\theta | v_1 = v, v_2 = 0, w_1 = w\}$$

$$f(x, \theta) = w_1 \varphi(J_1 \cdot x) + w_2 \varphi(J_2 \cdot x)$$

$$v_1 = 0$$

$$v_2 = 0$$



座標変換

$$\nu = \frac{w_1 J_1 + w_2 J_2}{w_1 + w_2},$$

$$w = w_1 + w_2,$$

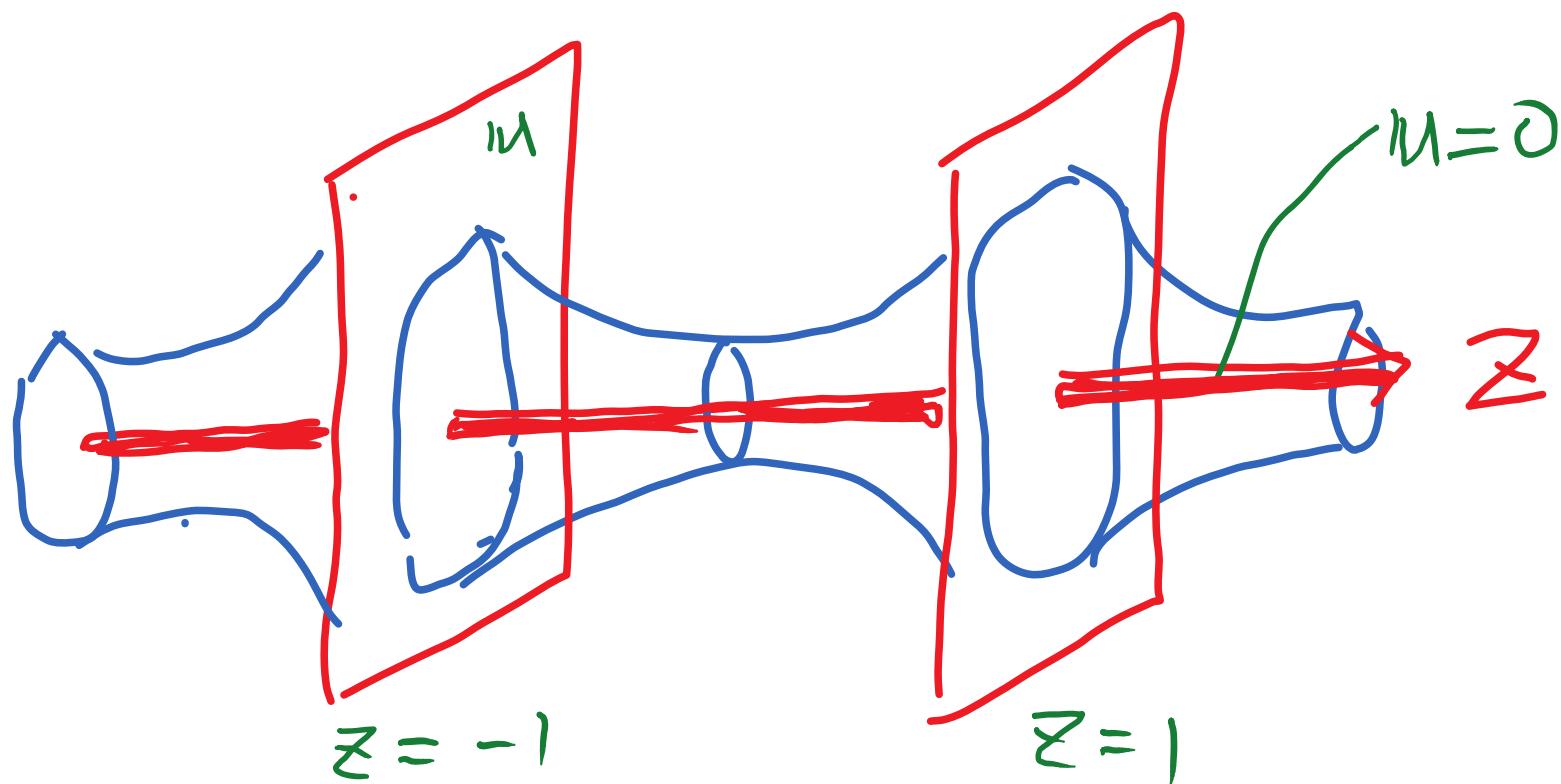
$$u = J_2 - J_1,$$

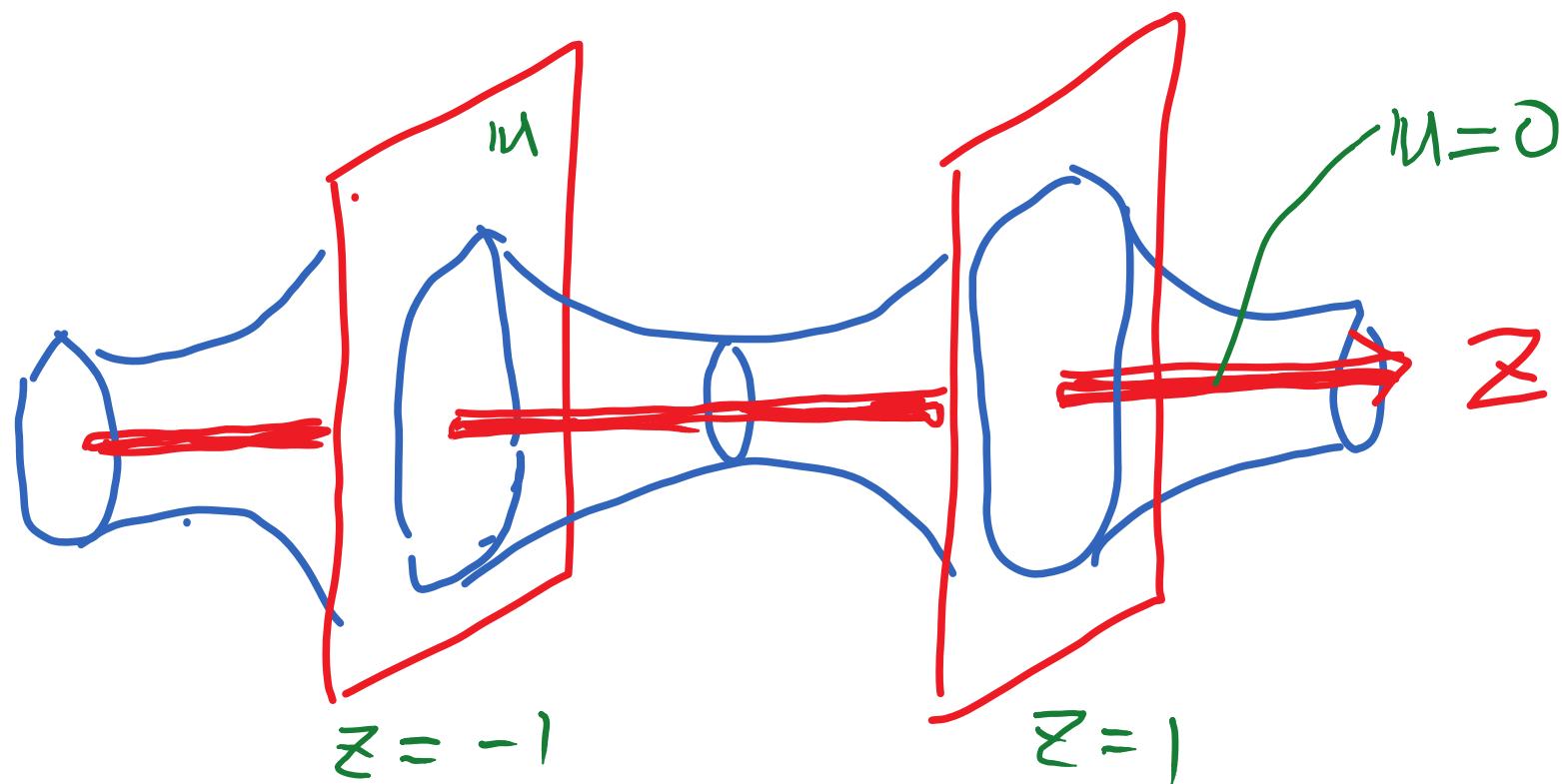
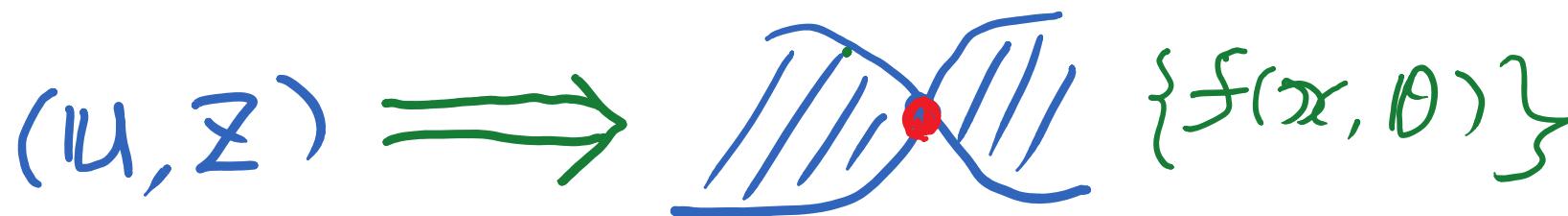
$$z = \frac{w_2 - w_1}{w_1 + w_2}$$

$$\xi = (\nu, w, u, z)$$

特異領域

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





Taylor expansion u : small

$$f(x, \xi) = w\varphi(v \cdot x) + \frac{w}{8}\varphi''(v \cdot x)(1 - z^2)(u \cdot x)^2$$

$$-\frac{w}{24}\varphi'''(v \cdot x)z(1 - z^2)(u \cdot x)^3 + \dots$$

fast dynamics $\rightarrow w, v$: stability

slow dynamics $\rightarrow u, z$

Rの近傍でのダイナミックス

$$\dot{\mathbf{u}} = -\frac{\eta w^*}{2} (1 - z^2) \langle \varphi'' e(\mathbf{u} \cdot \mathbf{x}) \mathbf{x} \rangle$$

$$\dot{z} = -\frac{\eta}{4w^*} z (3 + z^2) \langle \varphi'' e(\mathbf{u} \cdot \mathbf{x})^2 \rangle$$

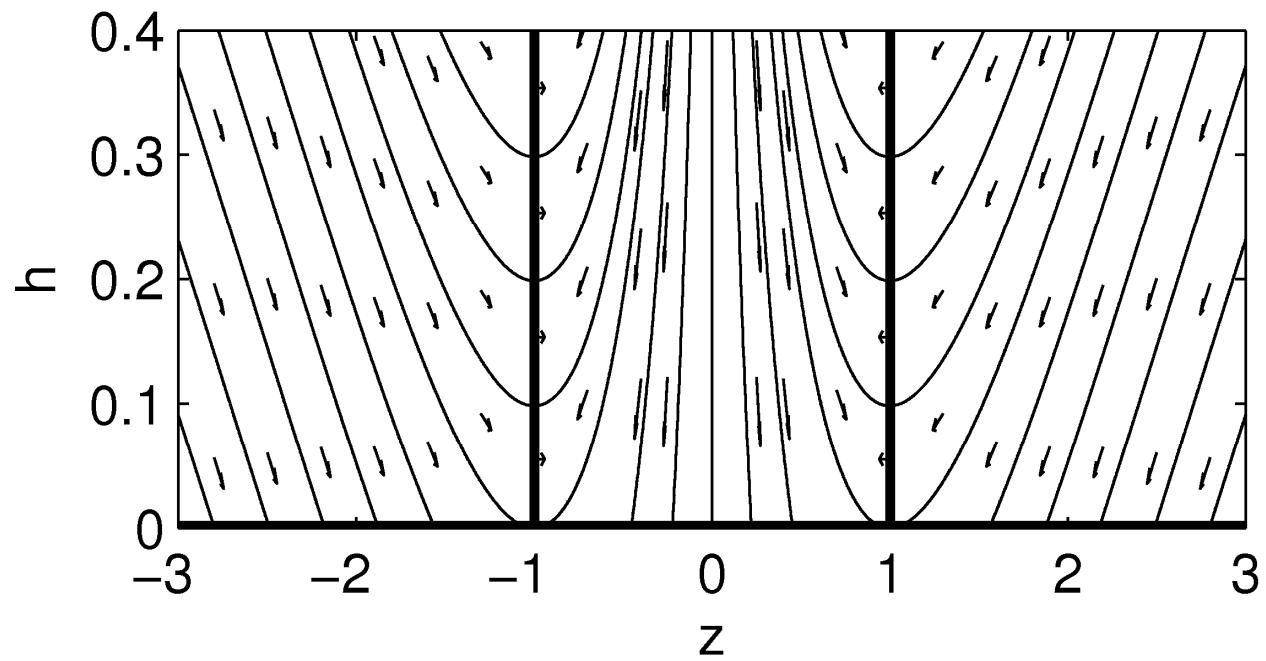
解の軌道

$$\frac{1}{2} |\mathbf{u}(t)|^2 = \frac{2}{3} w^* \log \frac{(z^2 + 3)^2}{|z(t)|} + c$$

安定論 1

true solution is in R :

$$R = \{u = 0 \text{ or } z = \pm 1\} : \text{stable}$$



Dynamic vector fields: Redundant case

安定論 2 : true solution is outside R

$$H = \langle \varphi'' e(x) \mathbf{x} \mathbf{x}^T \rangle$$

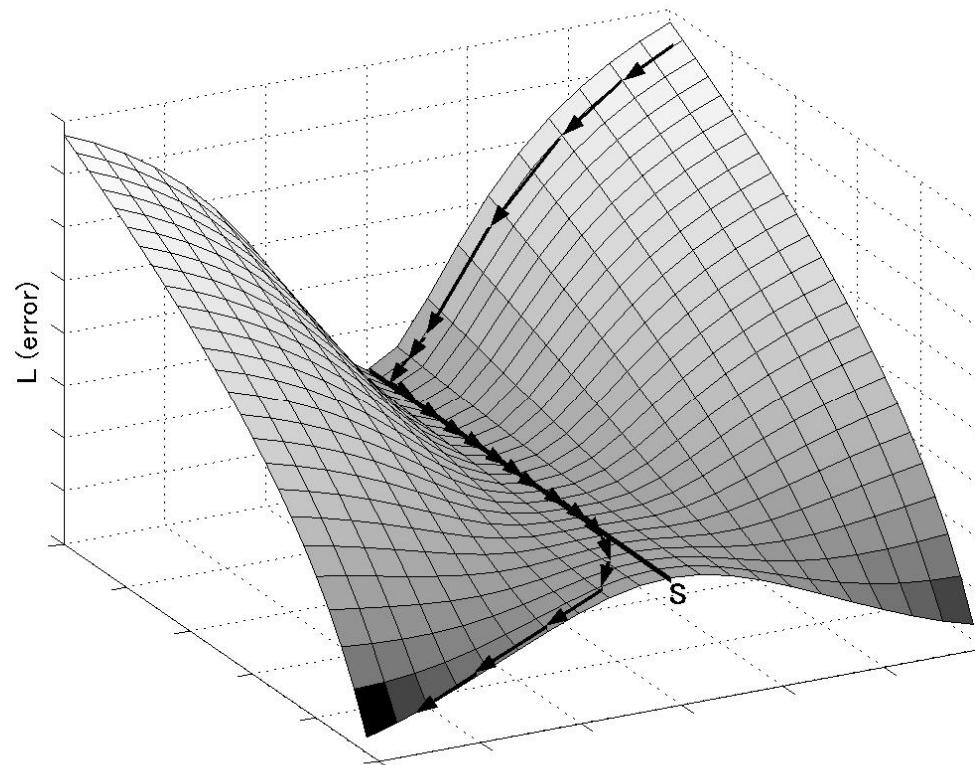
wH : positive-definite

$$\Leftrightarrow |z| < 1 \text{ stable} ; |z| > 1 \text{ unstable}$$

wH : negative-definite

$$\Leftrightarrow |z| > 1 \text{ stable} ; |z| < 1 \text{ unstable}$$

特異領域における解の軌道



Milnor attractor

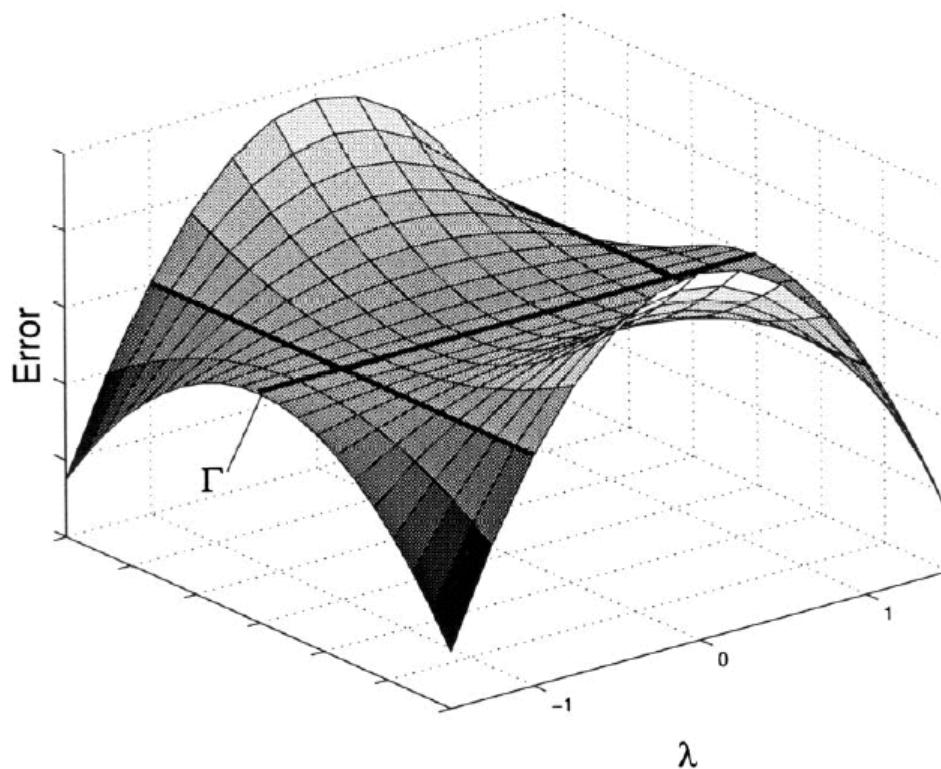


Fig. 5. Critical set with local minima and plateaus.

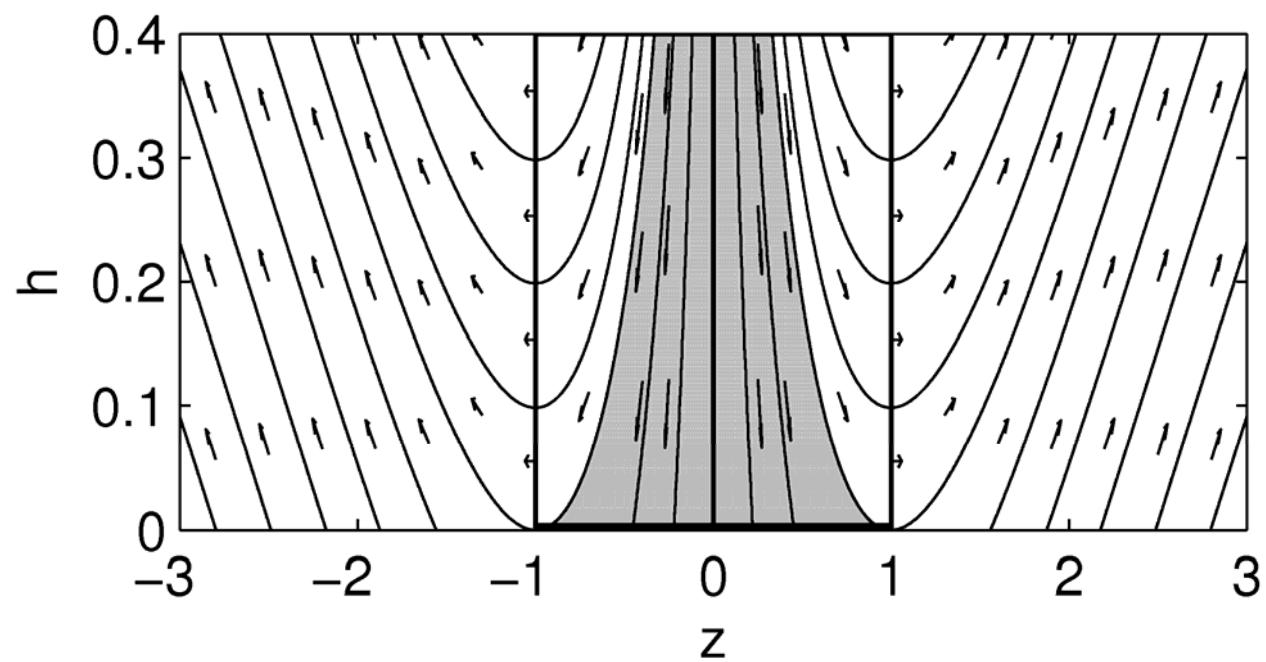
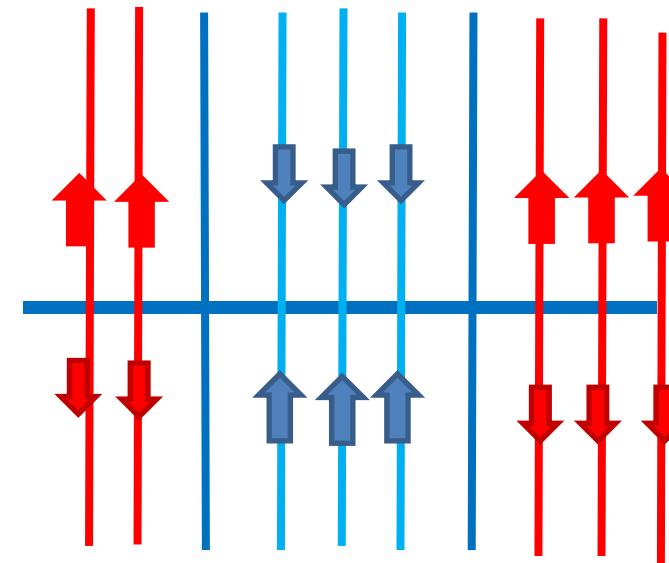
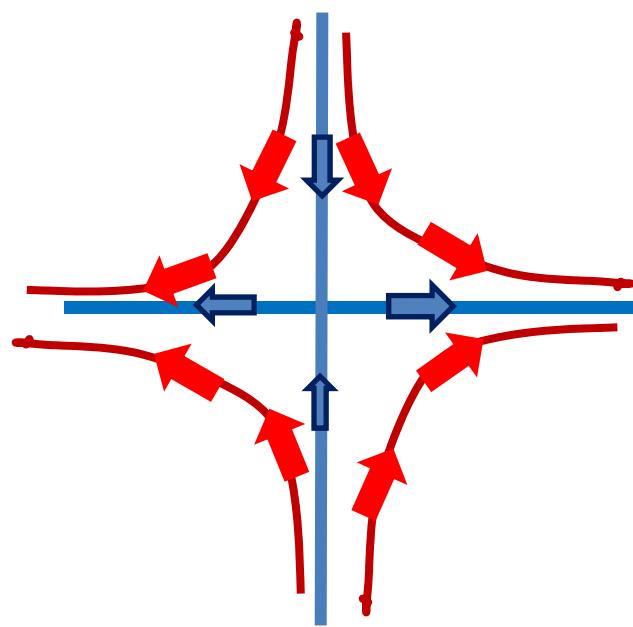


Fig. 2: trajectories

鞍点とプラトー



学習の遅滞

$$E = \frac{1}{2} \langle e^2 \rangle$$

$$\dot{E} = O(u^5)$$

$$\dot{E} = O(u^2)$$

特異領域Rのトポロジー

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

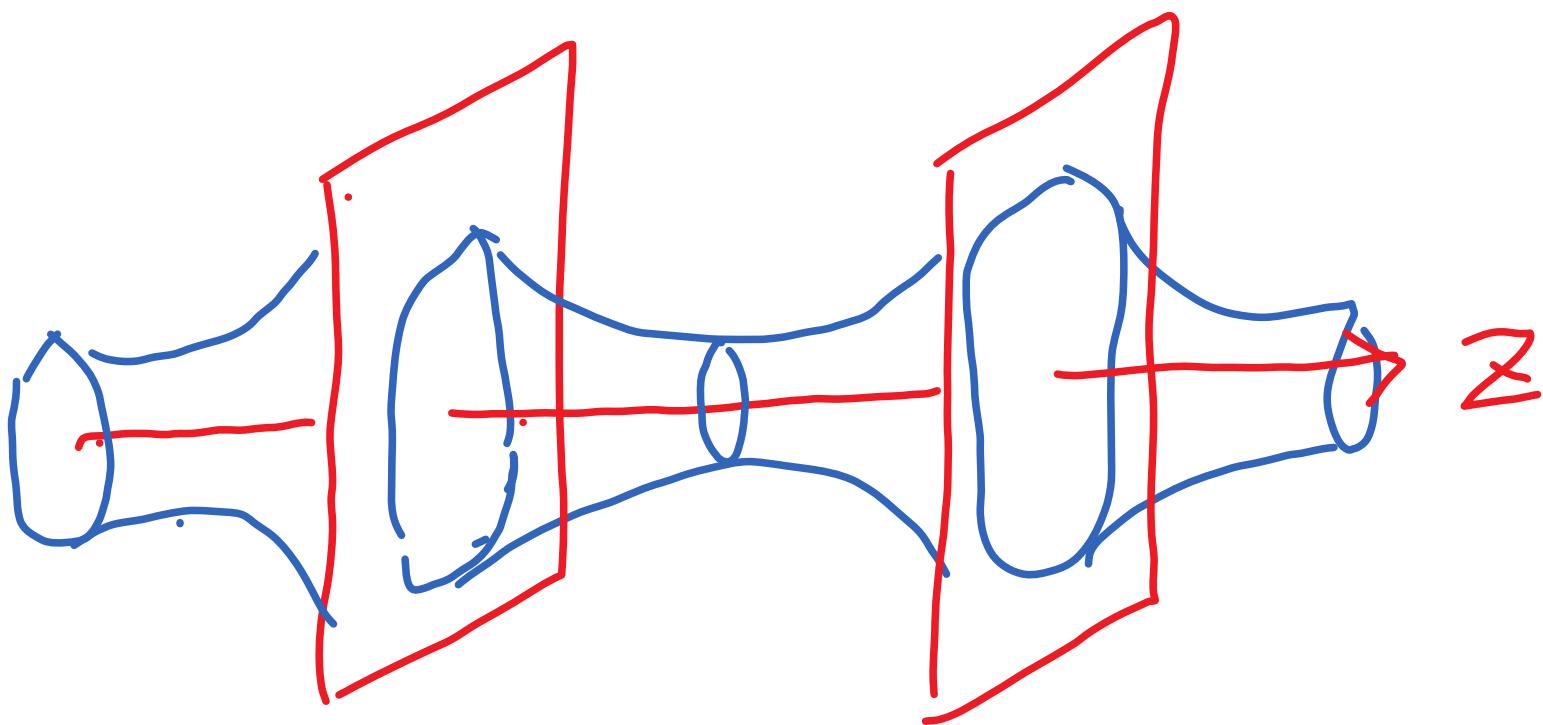
$$\tau = c_1 (1 - z^2) u^2, \quad u = |\mathbf{u}|$$

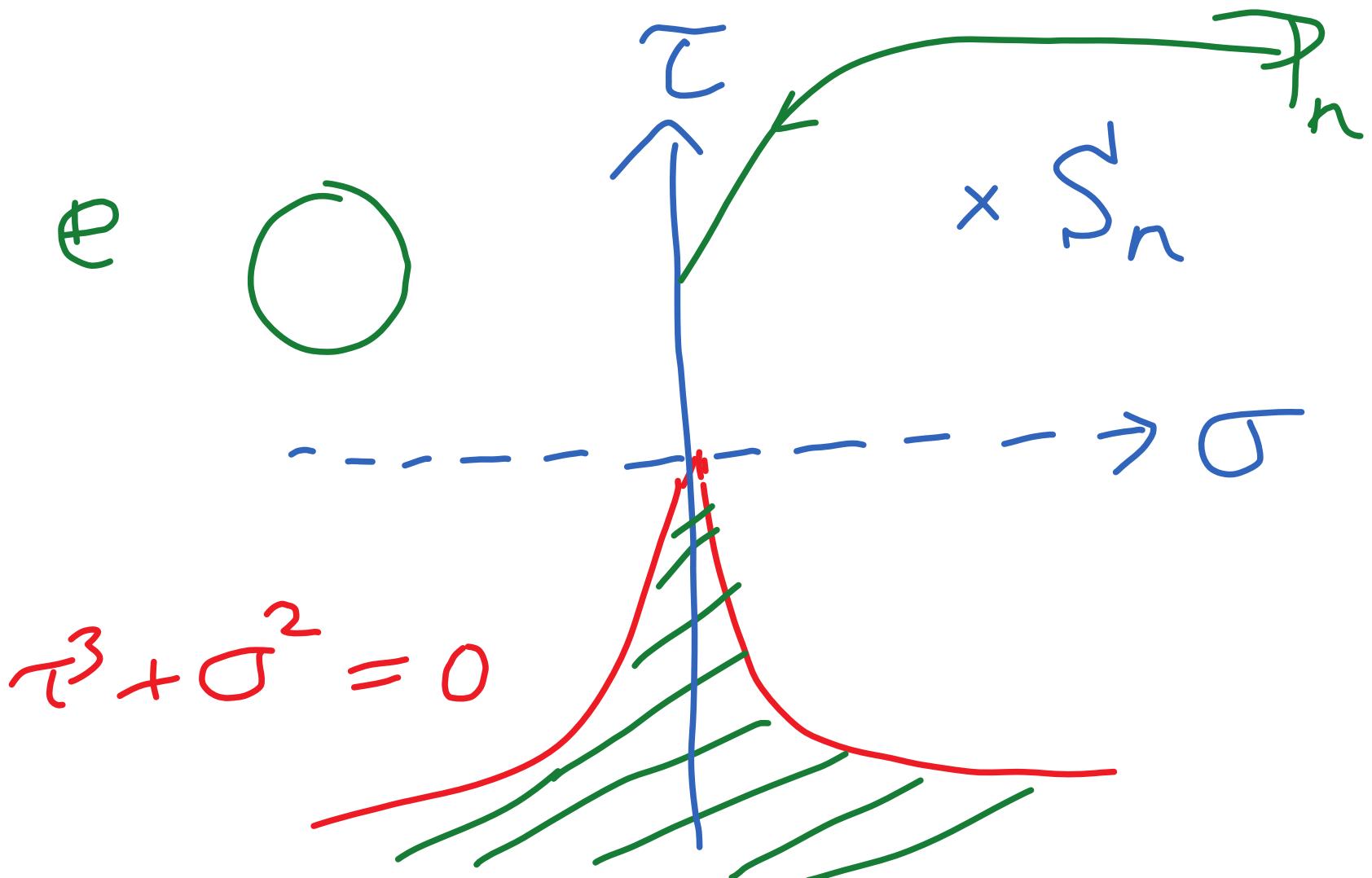
$$\sigma = c_2 z (1 - z^2) u^3,$$

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

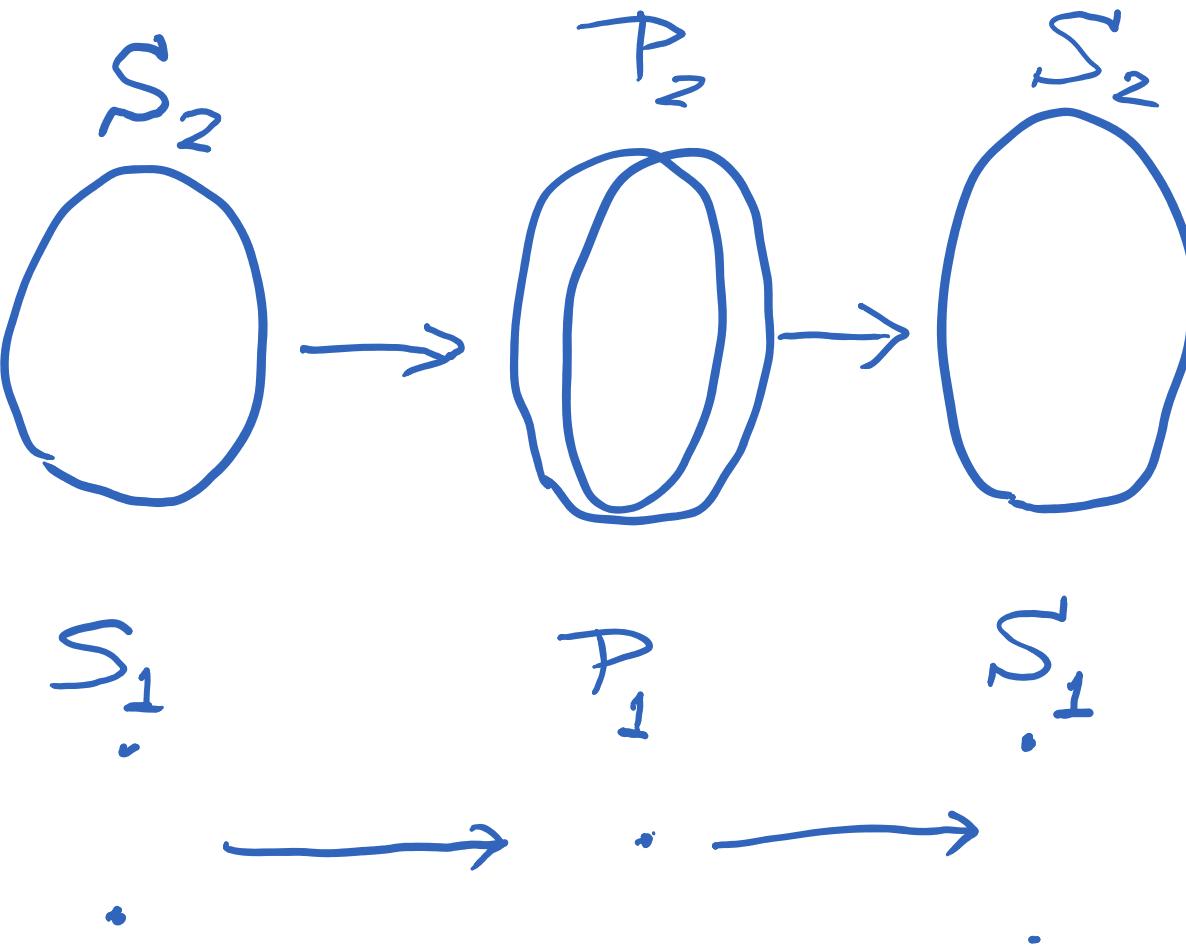
Singular Region

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





Sphere S_n and Projective space P_n



特異点における自然勾配学習

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

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

Milnor attractor

自然勾配学習の近似実現法

adaptive natural gradient

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

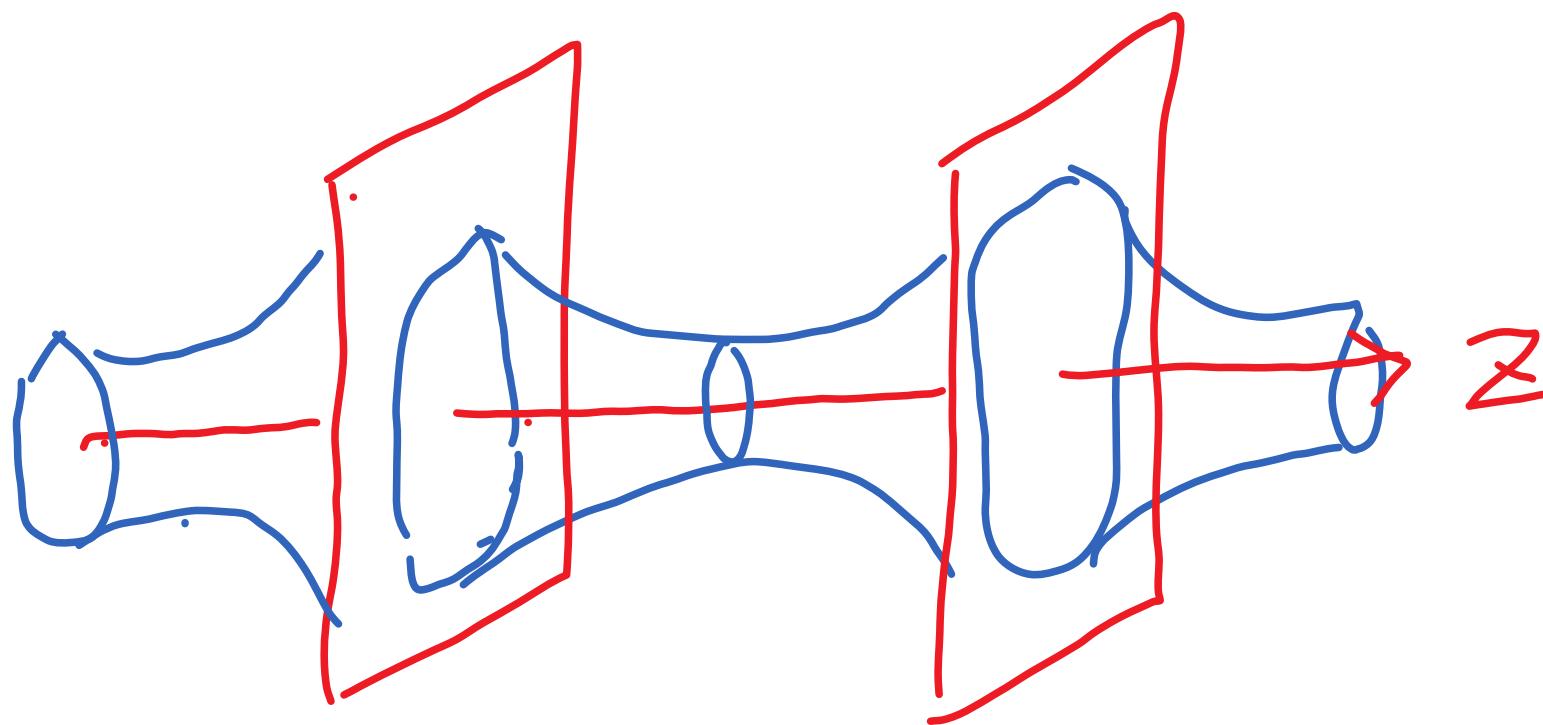
Unitwise diagonalization of G : Yan Olliver

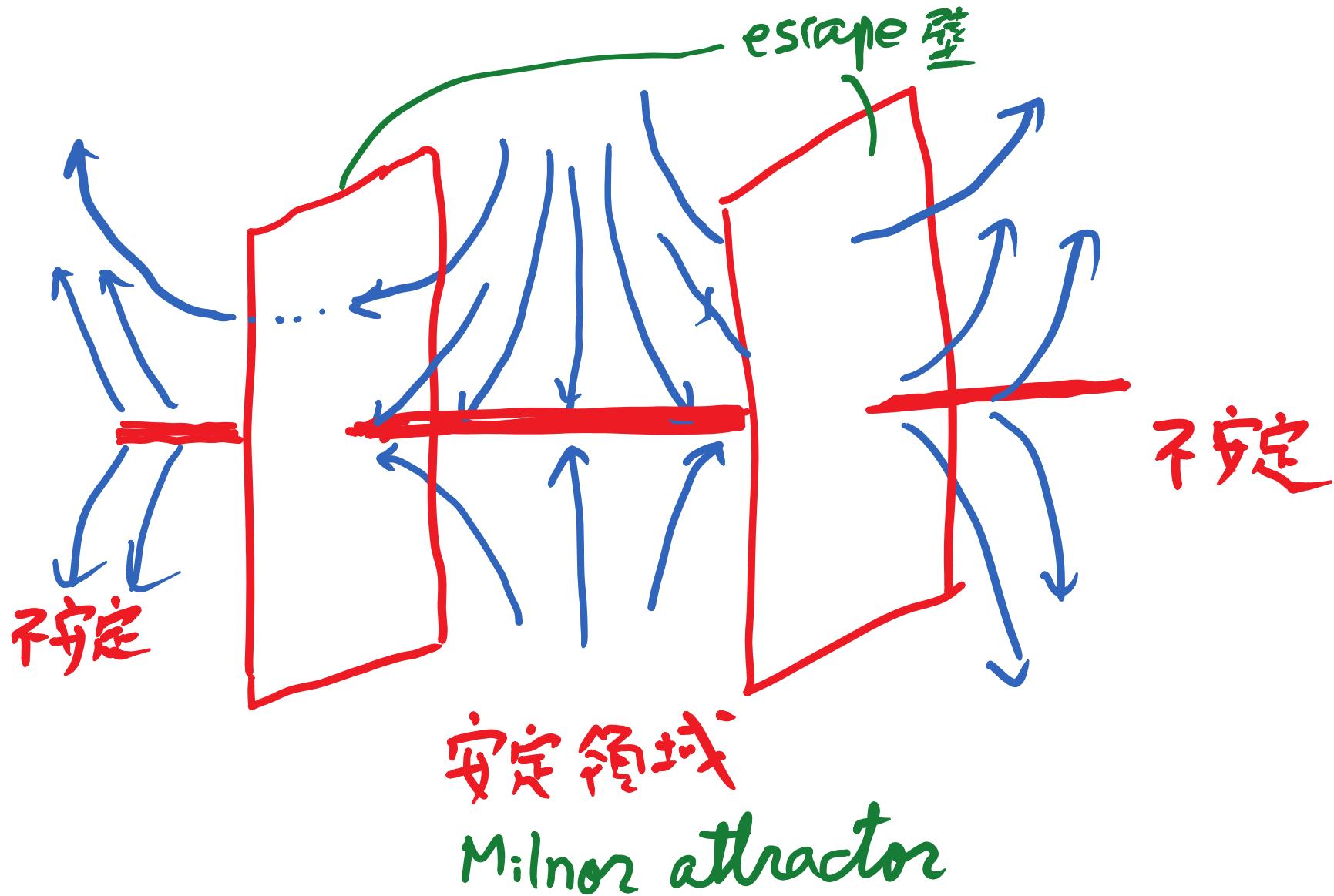
$G^{-1} \nabla l$: non-singular

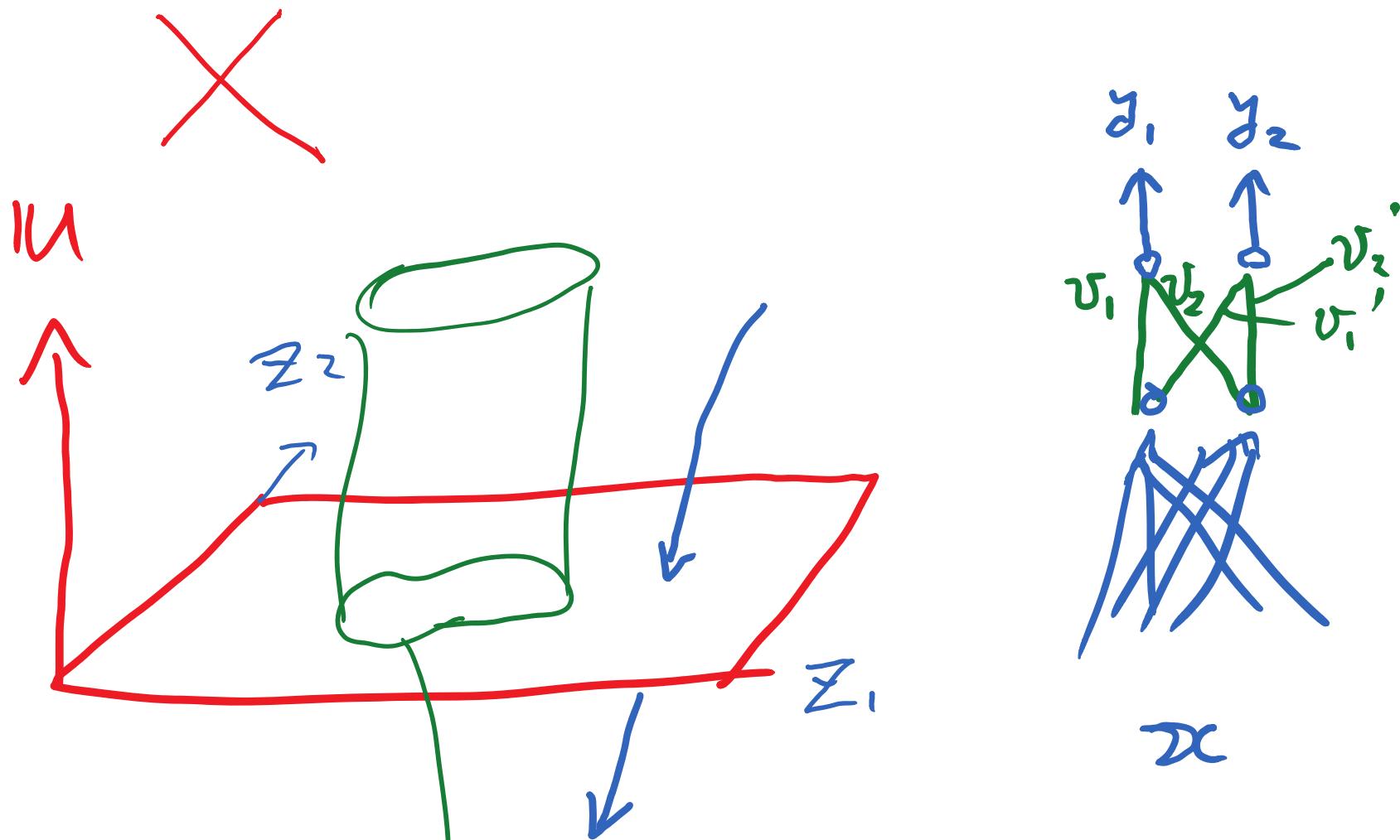
G: unitwise-diagonalization is OK (Ollvier)

Singular Region

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







安定領域の消失：通過点

脳に何を学ぶのか 意識と無意識のダイナミックス

記号 --- 興奮パターン

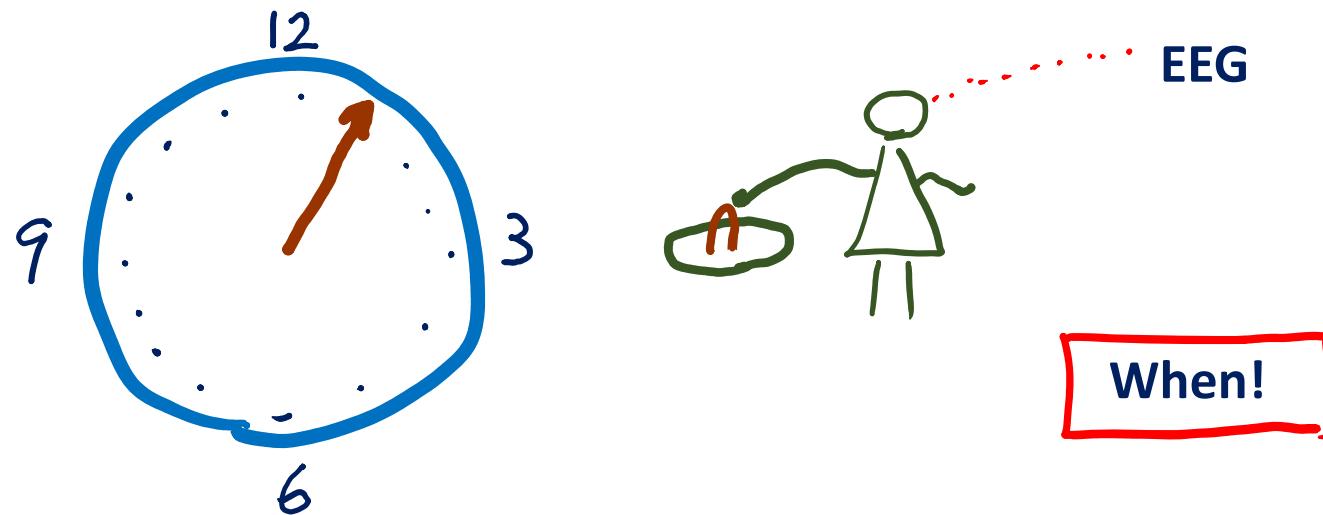
論理的推論 --- 並列ダイナミックス

AI

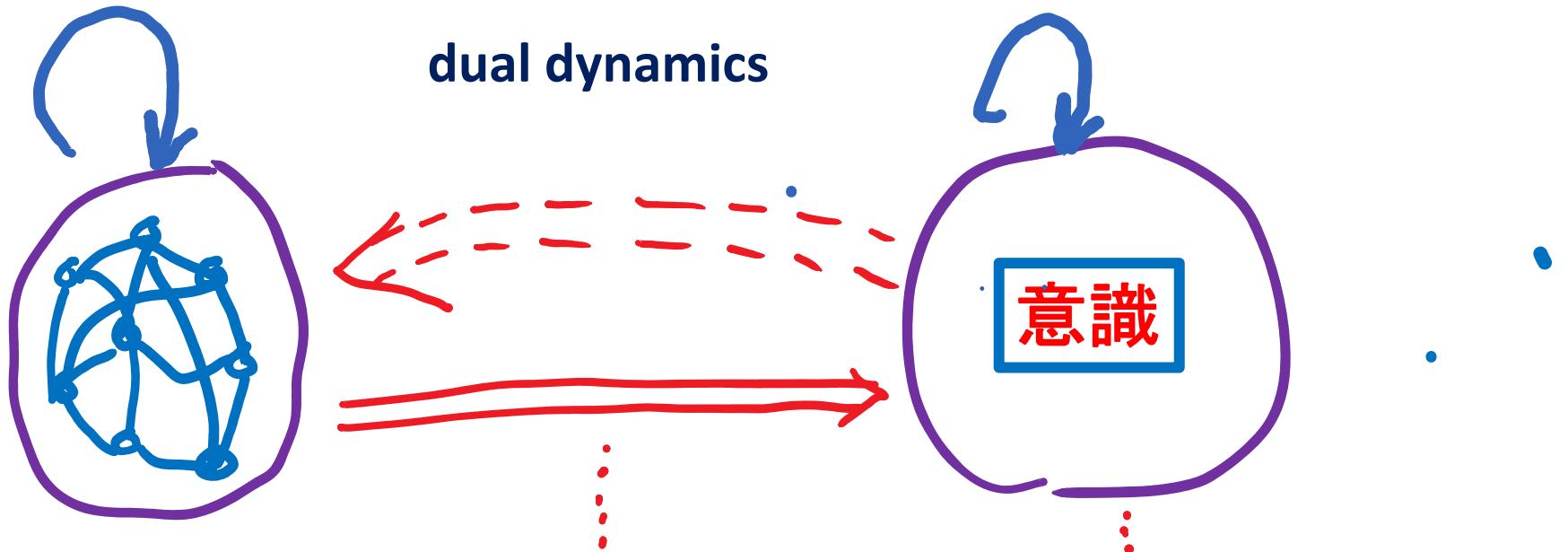
NN



Libet の実験：自由意志



予測(先付け)と後付け Prediction and Postdiction



ダイナミックス

意思決定と行動

反省、正当化、論理

意識の発生

共同作業、自分の意図を自分で知る

言語：論理的思考、数学

AI と脳科学の将来

Postdiction(後付け): 記号と論理を駆使
→パターン力学への介入

脳型の記憶方式 (連想記憶)

原理の共有; 存在証明

意識の役割

情報統合理論と情報幾何：意識の定量化

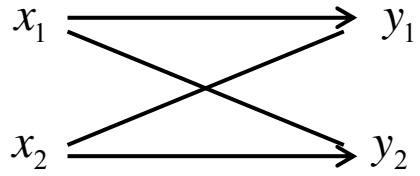
多層パーセプトロンのトポロジー

甘利俊一 理化学研究所脳科学総合研究センター

規範ダイバージェンスは定義できるか

多端子統計推論

情報統合理論と情報幾何: 意識の定量化



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



分離モデル: $S_S = \{q(\mathbf{x}, \mathbf{y})\}$

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

Tononi, Barrett and Seth, Ay

Split Model S_G

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

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

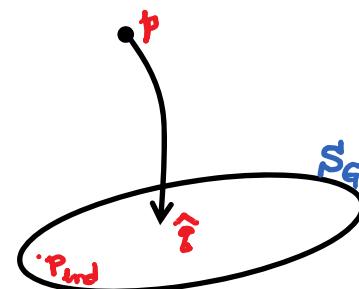
$$q(x_1, y_2 | x_2, y_1) = q(x_1 | x_2, y_1) q(y_2 | x_2, y_1)$$

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

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

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

graphical model



要請 1: $\Phi = D[p : q] \quad q : \text{split}$

要請 2: Markov条件を満たす最小の split モデル

$$X_1 \rightarrow X_2 \rightarrow Y_2$$

$$X_2 \rightarrow X_1 \rightarrow Y_1$$

要請 3 D は KL-ダイバージェンス

$S_H \subset S_G, S_D : S_G, S_H$ dually flat

$S_M \subset S_G : S_M, S_D$ not flat

$\Phi_G \leq \Phi_H, \Phi^*; \Phi_D \leq \Phi_H$

