

人工知能学会全国大会

脳、心、人工

付録一統計神経力学の構想

理化学研究所 甘利俊一

# 人工知能の衝撃

人間の知的能力を超えるのか？  
第4次産業革命

人工知能の歴史：

記号と論理ー並列分散（ニューラルネット）

# 脳—人間とは何か

思考・言語・意識  
人間・社会・文明



# 宇宙誌と脳 — 脳ができるまで

ビッグバン

(138億年前) 物理学・化学

生命

(36億年前) 生命科学

脳・神経系

(5億年前) 神経科学・情報科学

文明・社会

(20万年前?) 脳科学・情報科学・人間科学

物質の法則： 宇宙

生命の法則： 情報 + 物質： 進化

文明の法則： 心ころ + 情報 + 物質：  
社会・文化

# 人工知能と脳のモデル：一歴史の要約

第一次ブーム：記号と論理 VS パターンと学習

1956~

AI

Dartmouth 会議

記号と論理

知的推論、ゲーム

脳モデル

Perceptron

学習する普遍計算機構

実用的でない!!

暗黒期 (1965後半~1970's)

## 第2次ブーム

1970~ AI  
エキスパートシステム  
(MYCIN, DENDRAL)  
ミックス

1980~ BT (神経回路)  
MLP (backprop)  
連想記憶モデル、ダイナ

一兆円産業か？

沈静化  
確率Bayes推論  
chess (1997)

# 第3次ブーム 2010~ 脳型の人工知能（融合）

## 深層学習 Deep learning

（畳み込み多層回路（福島） + 確率勾配降下：  
日本でなぜ実現しな かったか？）

## 確率推論

深層学習の勝利 — 人間以上の識別能力

パターン認識: vision, auditory, sentence analysis

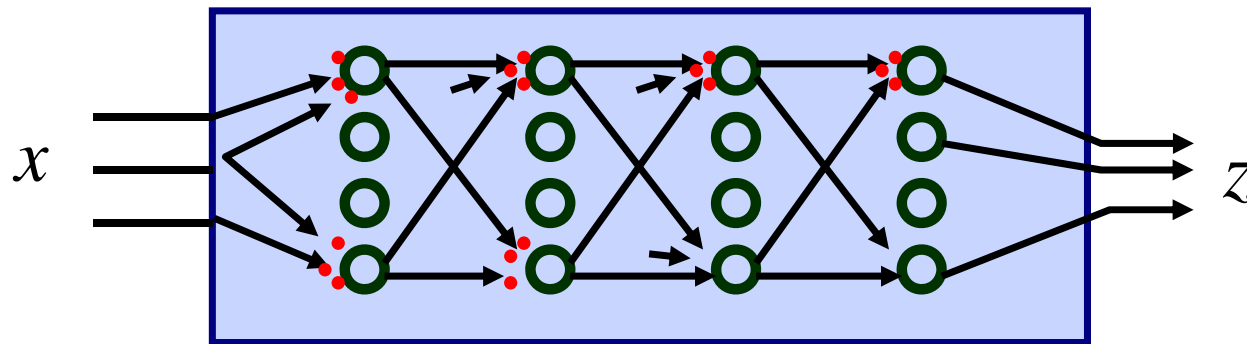
囲碁：強化学習

時系列とダイナミックス、動的パターン；言語処理

記号と論理 VS パターンとダイナミックス、学習 — 融合



# 層状学習回路網 multilayer perceptron



パーセプトロン Perceptron  
バックプロパゲーション Backpropagation

$$L(\mathbf{x}, W) = |y - \mathbf{g}(\mathbf{x}, W)|^2$$

$$\mathbf{w} \rightarrow \mathbf{w} + \Delta \mathbf{w}, \quad \Delta \mathbf{w} = -c \frac{\delta L(\mathbf{x}, W)}{\delta W}$$

# 最初のMLPの確率勾配降下学習法 (1967;1968)

情報科学講座 A・2・5



## 情報理論 II

—情報の幾何学的理論—

北川敏男編

編集委員

大泉充郎  
勝木保次  
北川敏男  
喜安善市  
栗原俊彦  
桑原万寿太郎  
坂井利之  
高田昇平  
次田皓一  
南雲仁雄  
中村幸弘  
和田弘

執筆者

甘利俊一 東京大学工学部

共立出版株式会社

1968

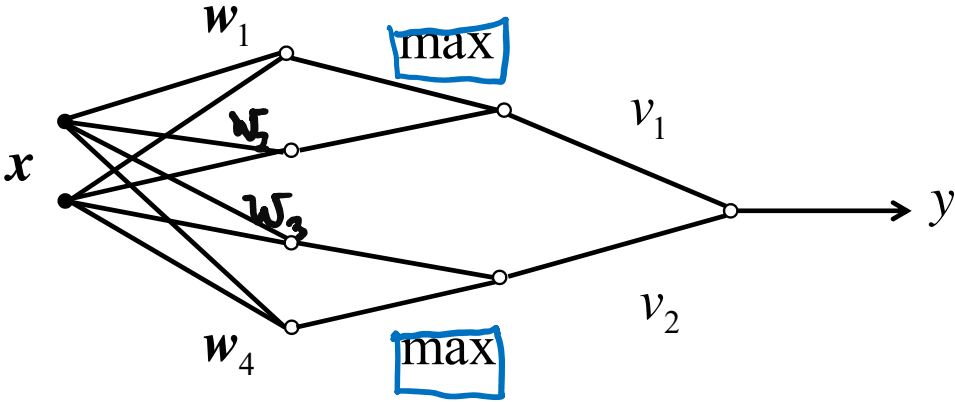
## Information Theory II --Geometrical Theory of Information

Shun-ichi Amari  
University of Tokyo

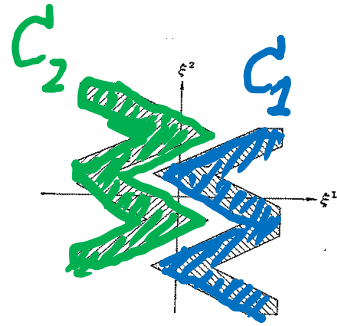
Kyouritu Press, Tokyo, 1968

$$f(x, \theta) = v_1 \max\{w_1 \cdot x, w_2 \cdot x\} + v_2 \min\{w_3 \cdot x, w_4 \cdot x\}$$

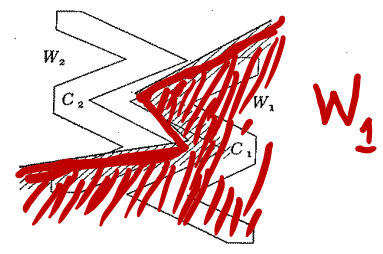
アナログニューロン、シグモイド関数



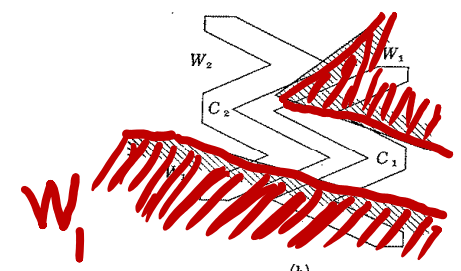
# 線形分離不可能 パターン分類



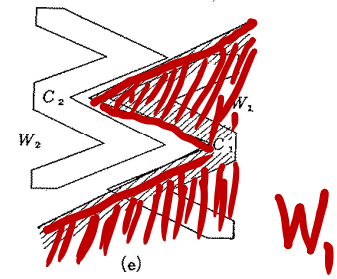
(a)



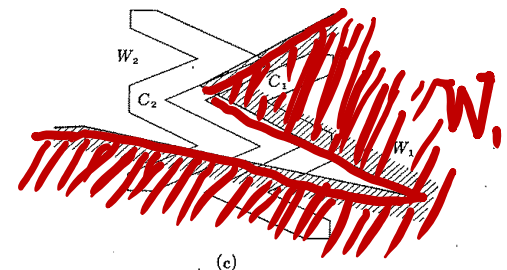
(d)



(b)



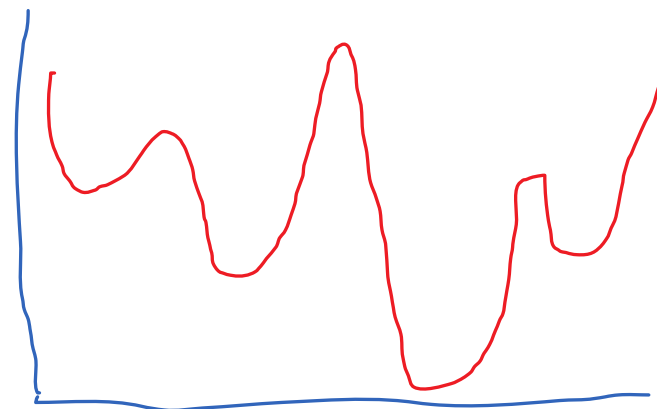
(e)



(c)

# 深層学習

自己組織化学習（教師なし学習）  
確率降下学習（教師あり学習）



## 情報表現の獲得

大域解と局所解：高次元

生成モデル GAN 囲碁、言語：何でもあり

## 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

# 深層学習

大量のデータ、計算力

## 入力を基に正解

現象の予測

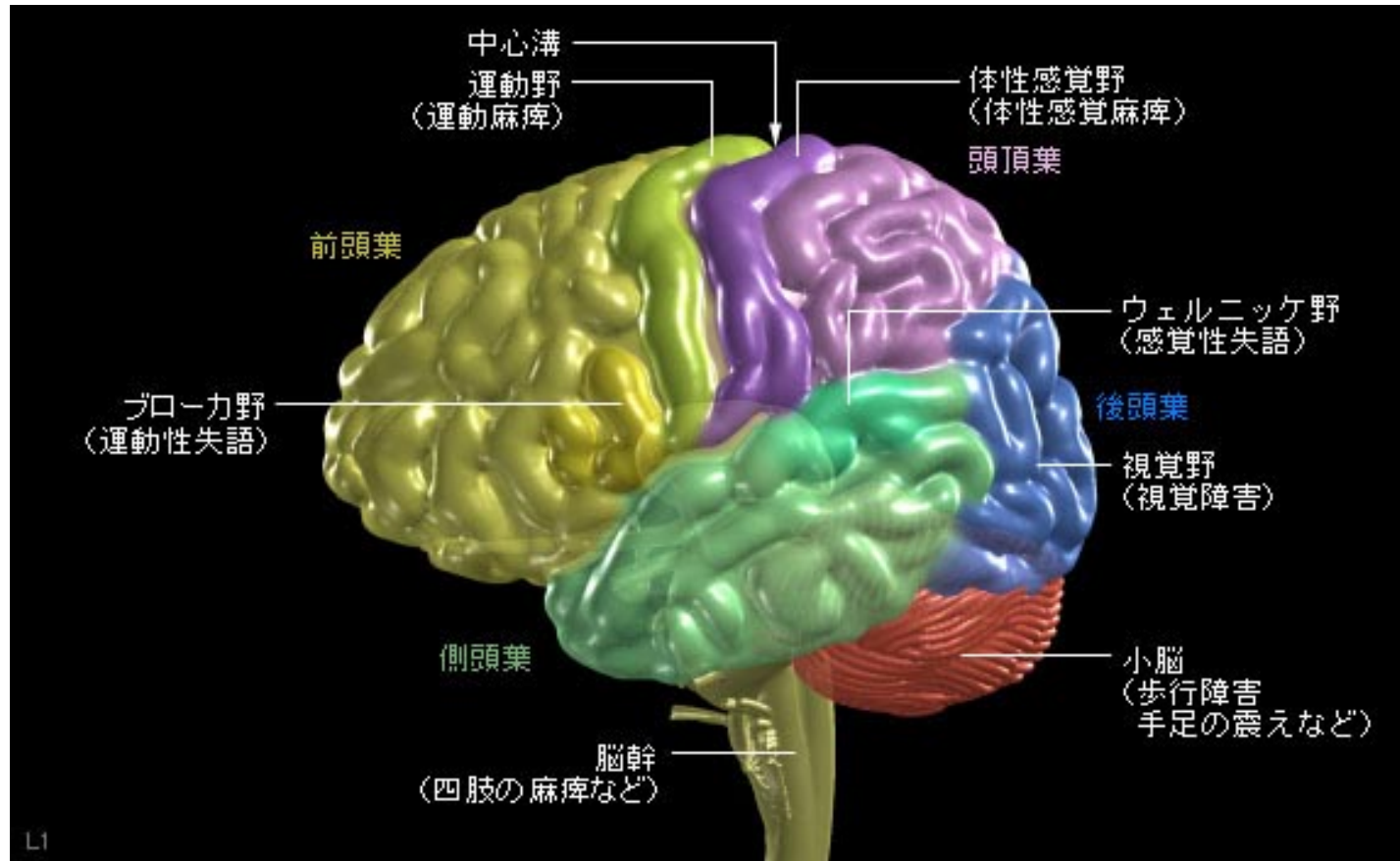
日蝕の予測

ケプラーの法則、ニュートン力学

原理の創出・理解一人間

# 脳：大脳、海馬、小脳、脳幹

脳科学：ミクローマクロ、理論：神経回路網





# 数理脳科学は脳の基本原理を探求する

単純な基本モデルを用いる：数理的探索（現実とは違う）

→ 計算論的神経科学

（脳はいかにこの原理を実現したか）

→ AI：技術による原理の実現（脳とは違う）

→ 神経科学（ありのままの脳）

# 脳は基本原理をどう実現したか

進化によるランダムサーチ

使える材料の制約

歴史的な制約

ごたごたの設計の中で精妙な実現：超複雑

# 人工知能は何をどう実現するか？

# 人工知能は脳に何を学ぶのか： 心 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



# 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

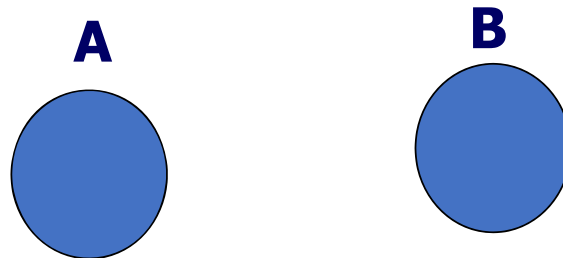
# 心の理論



# 葛藤する心

究極のゲーム  
(ultimatum game)

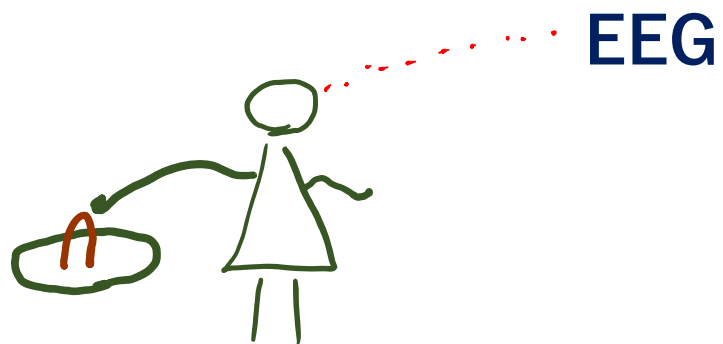
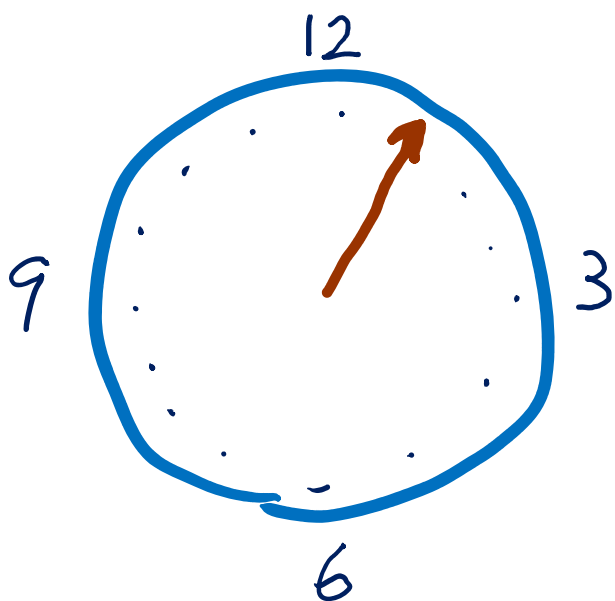
10万円



**A:** 配分を決定する----- **7万円-3万円**  
**B:** 同意または拒否

この時の脳の働き、各領野の確執  
利益、公正、その後のこと、評判  
二人か社会ゲームか

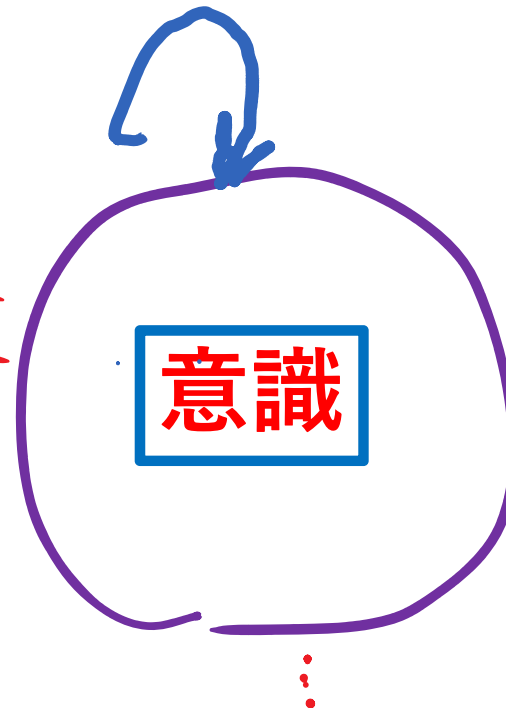
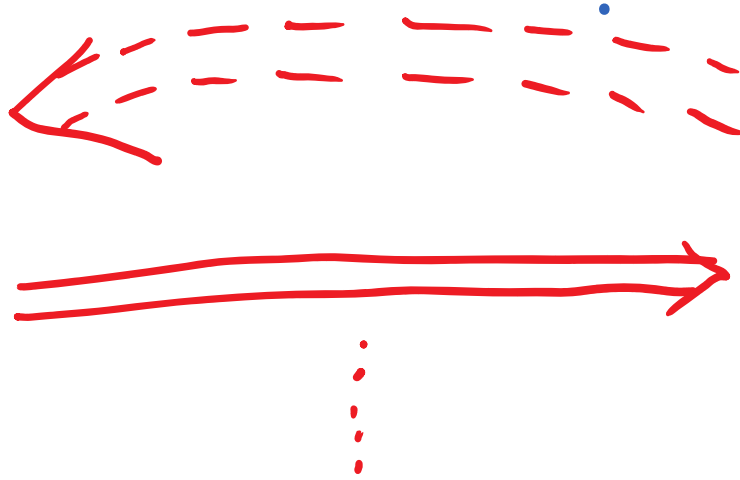
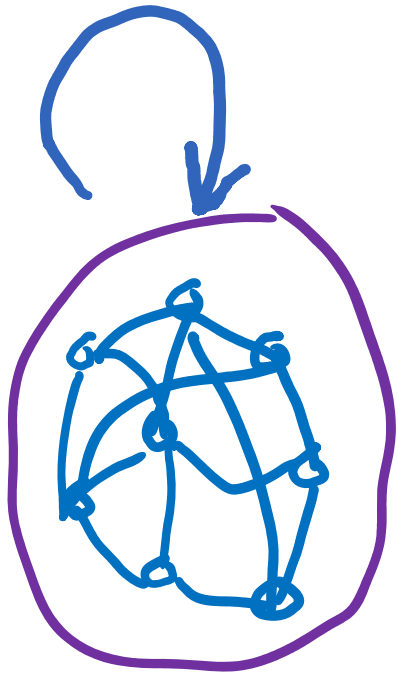
# Libet の実験：自由意志



**When!**

# 予測(先付け)と後付け prediction and postdiction

dual dynamics



ダイナミクス

意思決定と行動

反省、正当化、論理



**人工知能が脳に学ぶべきこと： 数理的な理解**

**判断；制御；認知；記憶**

**意識と心の役割；後付け**

**連想式記憶システム；知識体系**

心の理論

ロボットに意識はあるか



# 心を持ったロボットがつくれるか？

人の心の動きを理解する

ロボットが心を持つように見える  
(感情移入)

# ロボットが心を持てるのか？

## 人間の心：進化の産物

### 意識、意図、論理、感情

種の生存と個の不合理

人間は不合理； 芸術、喜び、愛、苦悩：使命感

ただ一度の、かけがいのない人生

ロボットは合理的

# 人工知能と倫理

人工知能の安全性、制御可能性

人工知能と戦争；人工知能の金融支配；  
支配の道具、格差

暴走： 人間の暴走を範として

# 社会への影響

失業問題：人口減　：AIは仕事を奪うか？　より高度な仕事

格差の拡大：

ベーシックインカムと人類の家畜化：働く喜び

# 人工知能と技術的特異点

2045

人工知能が人間を超えるとき  
人工知能が研究し、技術を進める

人間は素晴らしいが、愚かである。

人間はどんな知能システムを作るのか？  
社会の進化と支配

# 人工知能と未来社会の設計

深層学習を超えて

科学研究、技術開発

社会、文明

その脆弱性・崩壊



# 統計神経力学

Rozonoer (1969)

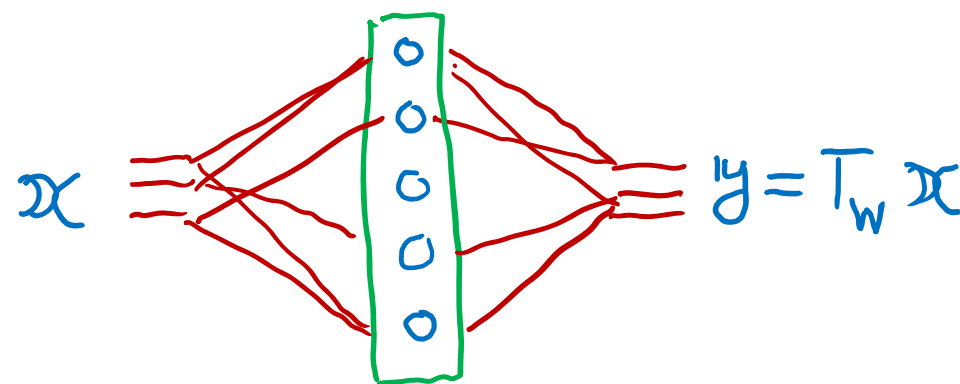
Amari (1971; 1974)

Amari et al (2013)

Toyoizumi et al (2015)

Poole, ..., Ganguli (2016)

Schoenholz et al (2017)



$$w_{ij} \sim N(0, 1)$$

## 巨視的振舞い

ほとんどすべての (典型的) 回路に共通

# 巨視變數

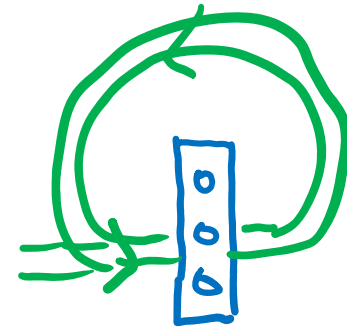
$$\text{活動度: } A = \frac{1}{n} \sum x_i^2$$

$$\text{距離・計量: } D = D[\mathbf{x} : \mathbf{x}']$$

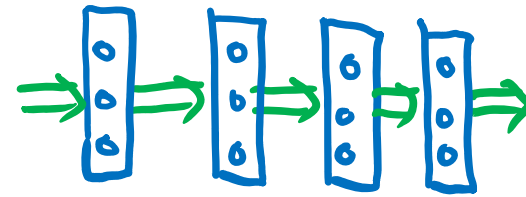
曲率:

$$A_{l+1} = F(A_l)$$

$$D_{l+1} = K(D_l)$$



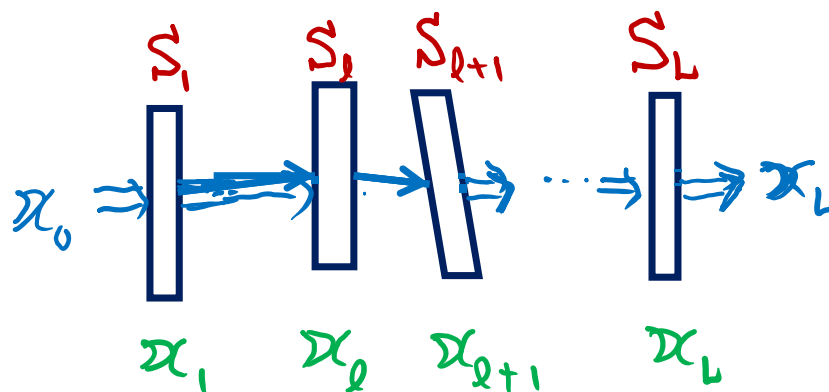
recurrent dynamics



multilayer dynamics

# 深層回路

$$x_{l+1} = \varphi\left(\sum_l w_{ij} x_l + w_{0i}\right)$$



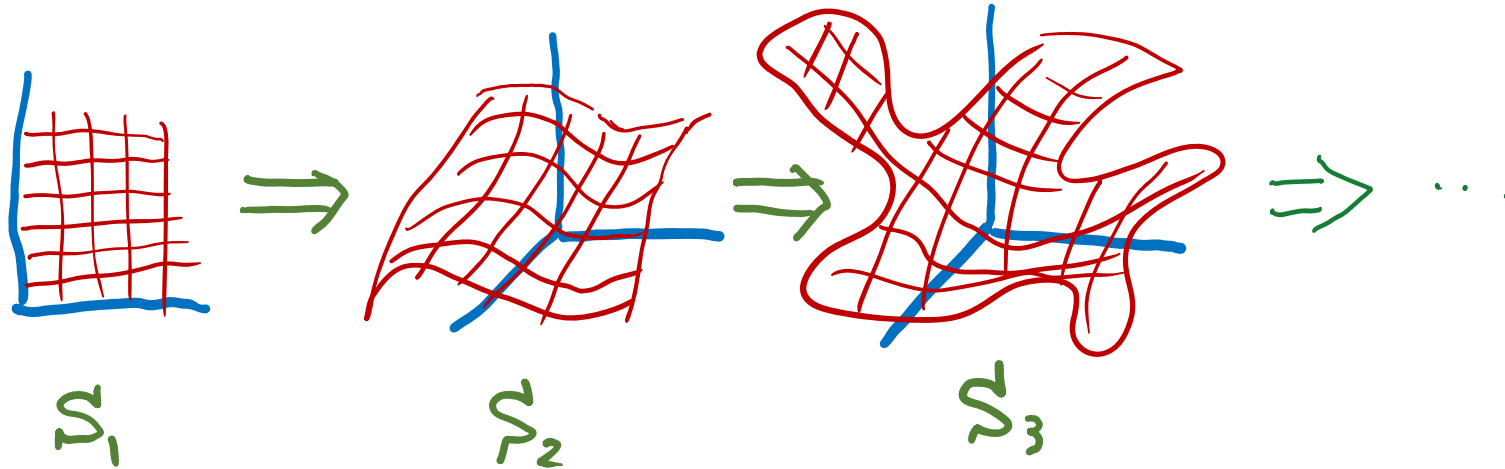
$$A_l = \frac{1}{n_l} \sum_l x_l^2$$

$$w_{ij} \sim N(0, \sigma^2 / \sqrt{n})$$

$$A_{l+1} = F(A_l)$$

$$w_{0i} = b \sim N(0, \sigma_b^2)$$

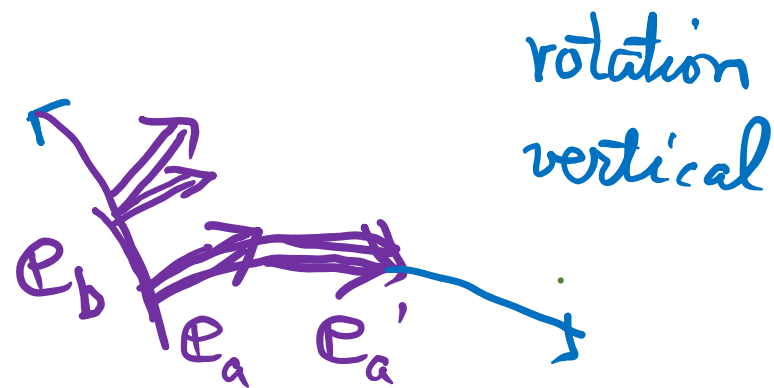
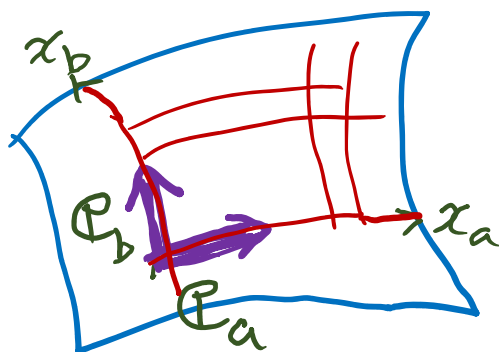
# 引き戻し計量 (リーマン計量・距離・曲率)



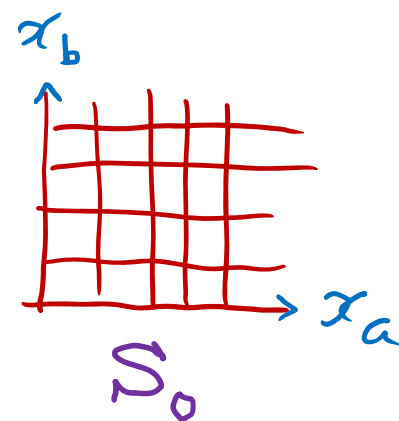
$$ds^2 = \sum g^l_{ab} dx^a dx^b = \frac{1}{n_l} d\mathbf{x}^l \cdot d\mathbf{x}^l$$

$$g^l_{ab} = \mathbf{e}^l_a \cdot \mathbf{e}^l_b$$

# 曲率

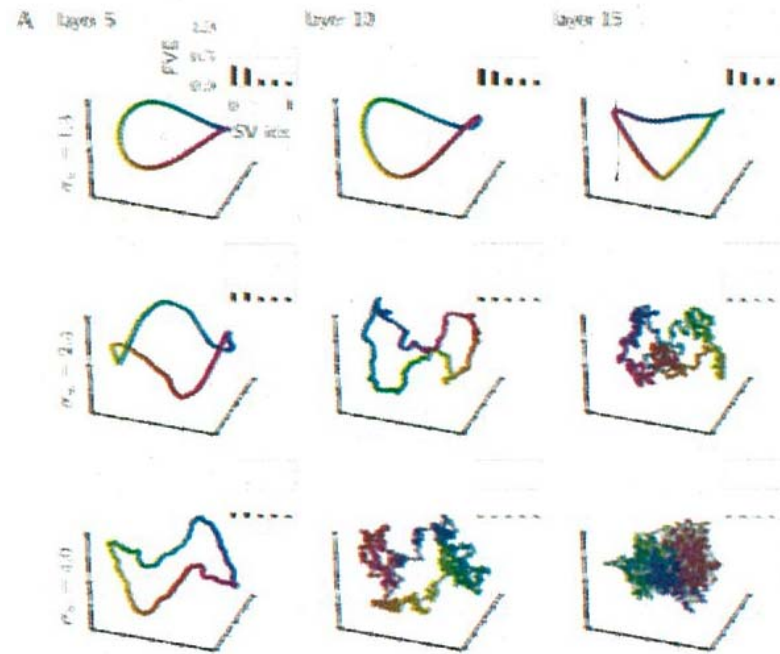
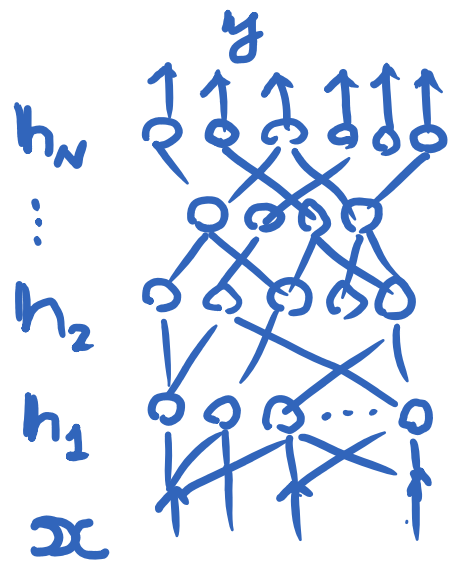


$$H_{abi}^{\ell} = \nabla_a e_b^{\ell}$$

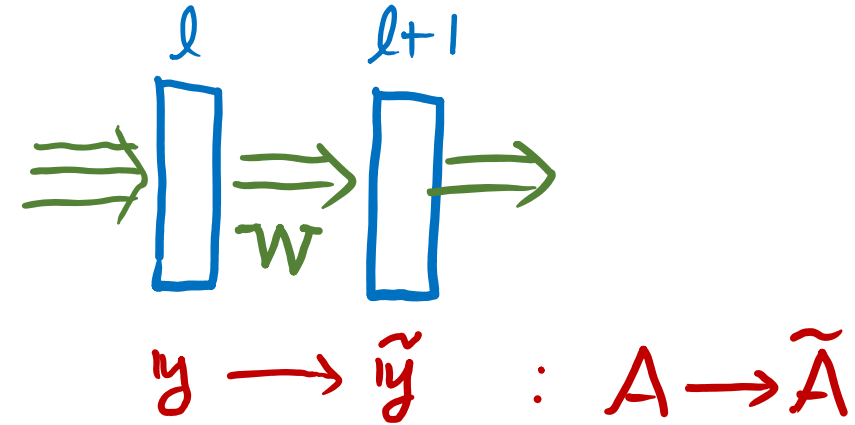


# Poole et al (2016)

## Random deep neural networks



# 活動度の力学



$$\tilde{y}_\alpha = \varphi\left(\sum w_{\alpha k} y_k + b_\alpha\right) = \varphi(u_\alpha)$$

$$u_\alpha \sim N(0, \sigma_A^2); \quad \sigma_A^2 = \sigma^2 A + \sigma_b^2$$

$$\tilde{A} = \frac{1}{n_{l+1}} \sum (\tilde{y}_\alpha)^2 = E[\varphi(u_\alpha)^2] = \chi_0(A) = \frac{1}{2\pi} \cos^{-1}\left(-\frac{\sigma_A^2}{1 + \sigma_A^2}\right)$$

$$\chi_0(A) = \int \varphi^2(\sqrt{A}v) Dv \quad v \sim N(0, 1)$$

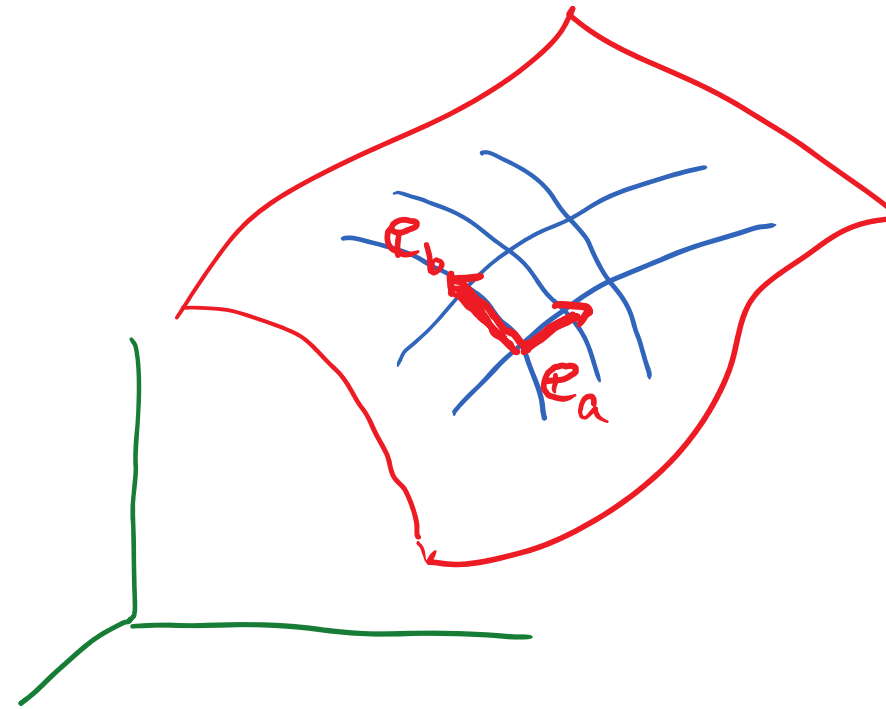
# Basis vectors

$$dx_{i_l} = \sum \varphi'(u_{i_l}) W_{i_{l-1}}^{i_l} dx_{i_{l-1}} = \sum B_{i_{l-1}}^{i_l} dx_{i_{l-1}}$$

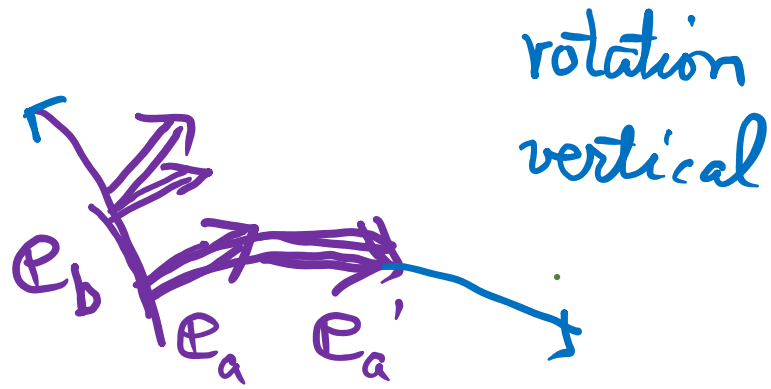
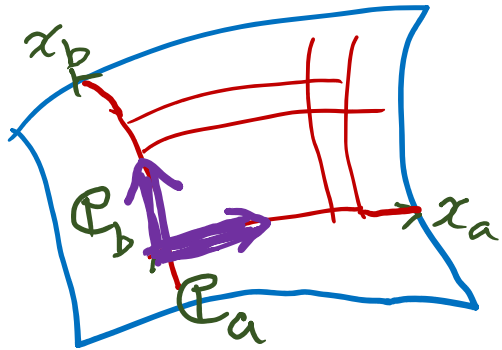
$$d\mathbf{x} = B d\mathbf{x} \quad (= B \dots B d\mathbf{x})$$

$$B_{i_{l-1}}^{i_l} = \varphi'(u_{i_l}) W_{i_{l-1}}^{i_l}$$

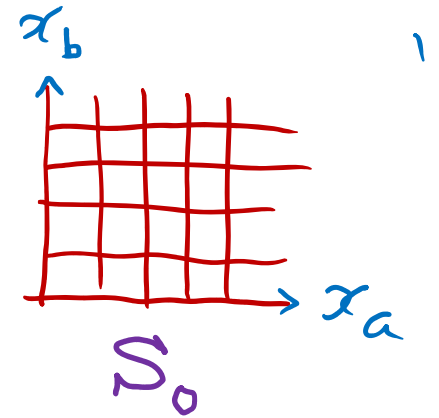
$$\mathbf{e}_a = B \mathbf{e}_a = B \dots B \mathbf{e}_a$$







$$g_{ab} = \frac{1}{n_l} \mathbf{e}_a \cdot \mathbf{e}_b$$



• • •

# Metric

$${}^l g_{ab} = \left\langle {}^l \mathbf{e}_a, {}^l \mathbf{e}_b \right\rangle = BB {}^{l-1} g_{ab}$$

$$ds^2 = \sum {}^l g_{ab} d {}^l x_a d {}^l x_b$$

$$BB = \sum_{i_l} W_{i_{l-1}}^{i_l} W_{i'_{l-1}}^{i_l} \varphi'(u_{i_l})^2 \approx \sigma_l^2 E[\varphi'^2] \delta_{i_{l-1} i'_{l-1}}$$

$$\chi_1 = \sigma_l^2 E\left[\varphi'(u_{i_l})^2\right]$$

# リーマン計量の力学

$$\tilde{y}_\alpha = \varphi\left(\sum w_{\alpha k} y_k + b_\alpha\right) = \varphi(u_\alpha)$$

$$d\tilde{y}_\alpha = \sum B_k^\alpha dy_k \quad \tilde{\mathbf{e}}_a = B\mathbf{e}_a$$

$$ds^2 = \sum g_{ij} dy^i dy^j = \langle d\mathbf{y}, d\mathbf{y} \rangle$$

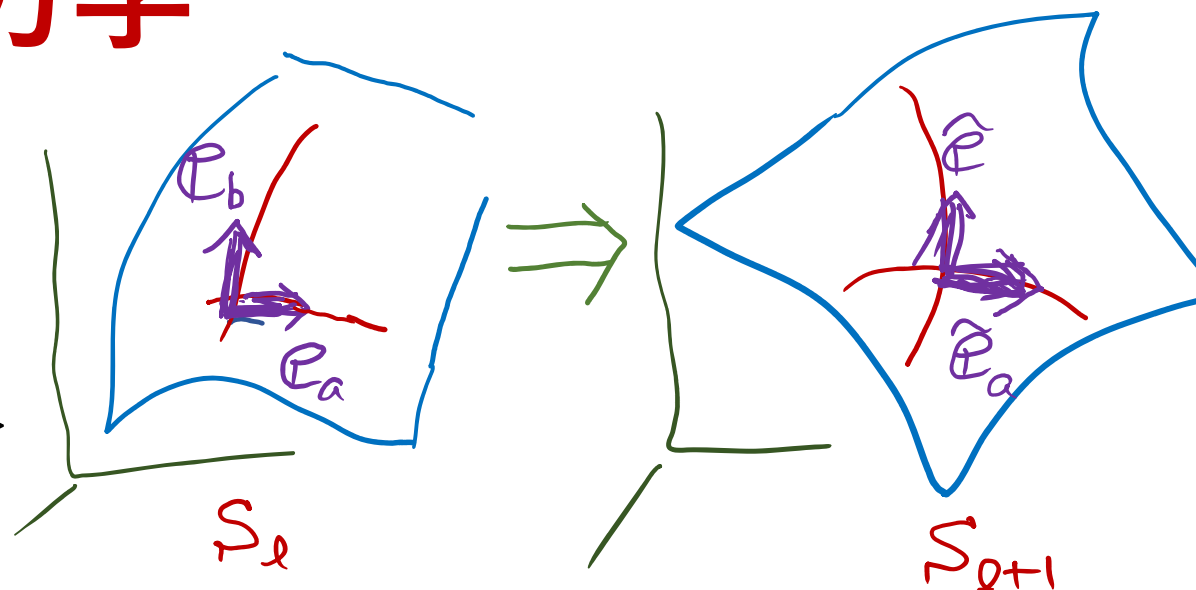
$$B = (B_k^\alpha) = (\varphi'(u_\alpha) w_k^\alpha)$$

$$\langle \tilde{\mathbf{e}}_a, \tilde{\mathbf{e}}_b \rangle = \tilde{g}_{ab} = \sum B_k^\alpha B_j^\alpha \langle \mathbf{e}_k, \mathbf{e}_j \rangle$$

$$E[(\varphi'(u_\alpha))^2 w_k^\alpha w_j^\alpha] = E[(\varphi'(u_\alpha))^2] E[w_k^\alpha w_j^\alpha]$$

平均場近似

$$\chi_1(A) = \int \sigma^2 \{\varphi'(\sqrt{A}v)\}^2 Dv = \frac{1}{2\pi} \frac{\sigma^2 A + \sigma_b^2}{\sqrt{1 + 2(\sigma^2 A + \sigma_b^2)}}$$



# law of large numbers

$$d s^2 = \sum (dx^i)^2$$

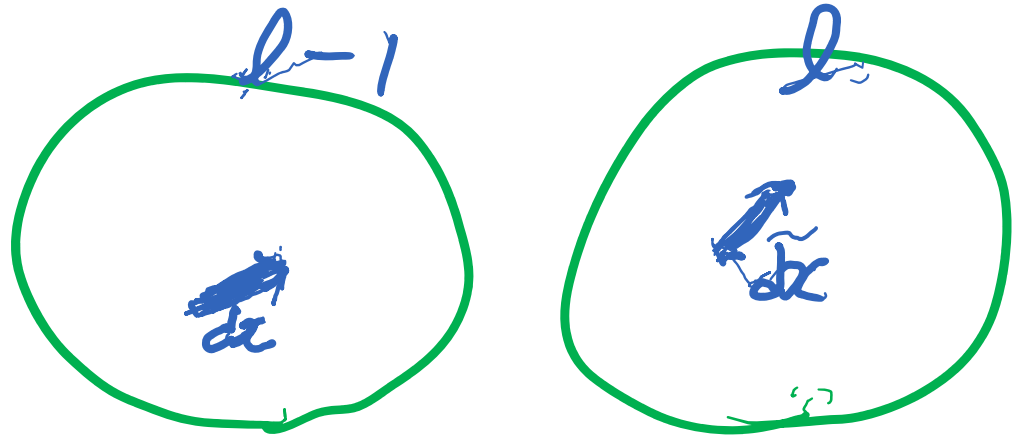
$$d s^2 = \sum_{i=1}^l B_{i-1}^i B_{i-1}^i dx^i dx^i$$

$$\sum_{i=1}^l B_{i-1}^i B_{i-1}^i = \sum [\varphi'(u_{i-1})^2 W_{i-1}^i W_{i-1}^i]$$

$$= n E[\varphi'(u_{i-1})^2 W_{i-1}^i W_{i-1}^i] + O_p(1/n)$$

$$= E[\varphi'(u_{i-1})^2] E[W_{i-1}^i W_{i-1}^i] = \chi \delta_{i-1}^i + O_p(1/n)$$

$$\chi = \sigma^2 \int \varphi'(u_i)^2 Dv$$



$$\tilde{g}_{ab} = \chi_1(A) g_{ab}$$

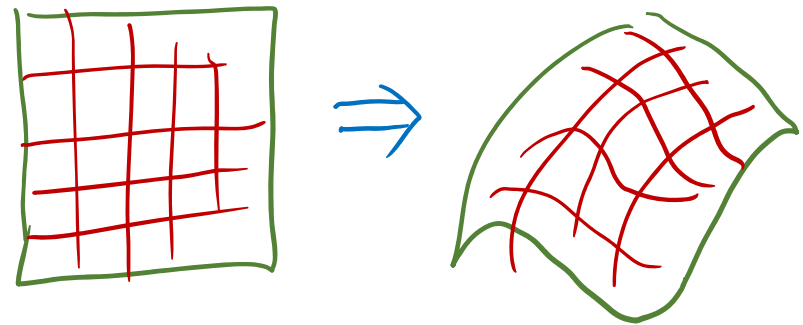
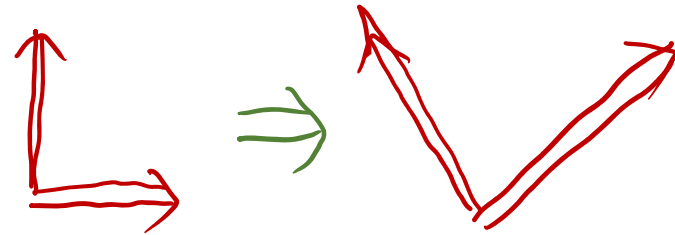
conformal transformation!

$$\bar{\chi}_1 = \bar{\chi}_1(\bar{A}) > 1:$$

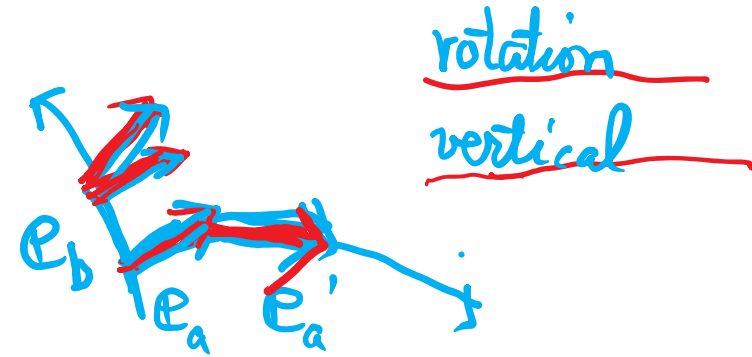
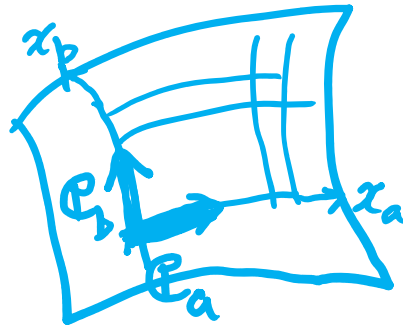
拡大 (カオス、Lyapunov指数)

$$\Rightarrow g^l_{ab} = \prod \chi_1(A^s) \delta_{ab}$$

rotation, expansion



# 曲率の力学



$$\tilde{H}_{ab}^{\alpha} = \nabla_a \tilde{\mathbf{e}}_b^{\alpha} = \partial_a \partial_b \tilde{y}^{\alpha}$$

$$= \varphi''(u_{\alpha})(\mathbf{w} \cdot \mathbf{e}_a)(\mathbf{w} \cdot \mathbf{e}_b) + \varphi'(\mathbf{w} \cdot \partial_a \mathbf{e}_b)$$

$$\tilde{\mathbf{H}}_{ab} = \mathbf{H}_{ab}^{\perp} + \mathbf{H}_{ab}^{\parallel}$$

Euler-Schouten曲率  
Affine connection

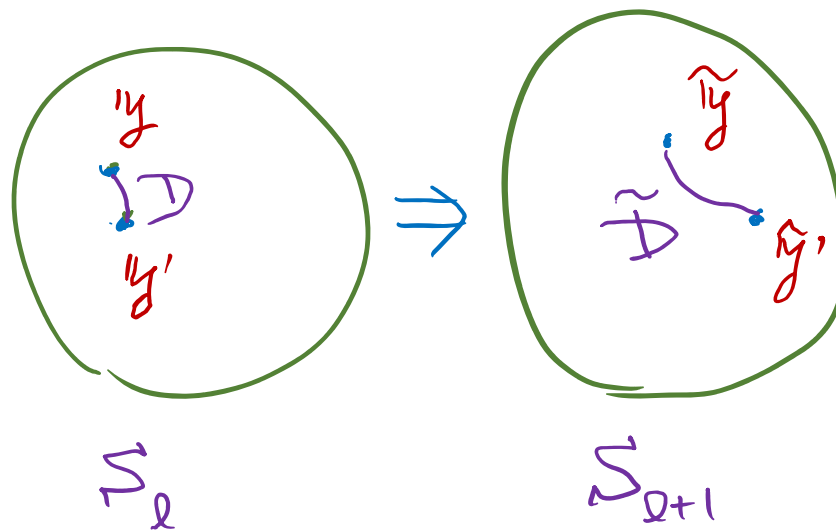
$$\tilde{H}_{ab}^2 = |\tilde{\mathbf{H}}_{ab}|^2$$

# 距離法則 (Amari, 1974)

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

$$C(x, x') = \frac{1}{n} x \cdot x' = \sum x_i x_i'$$

$$D = A + A' - 2C$$



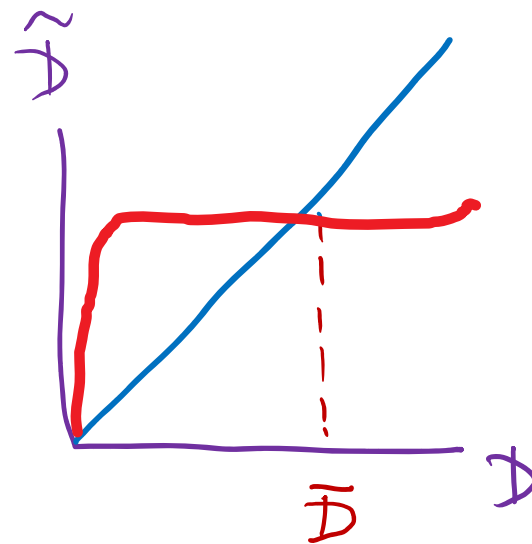
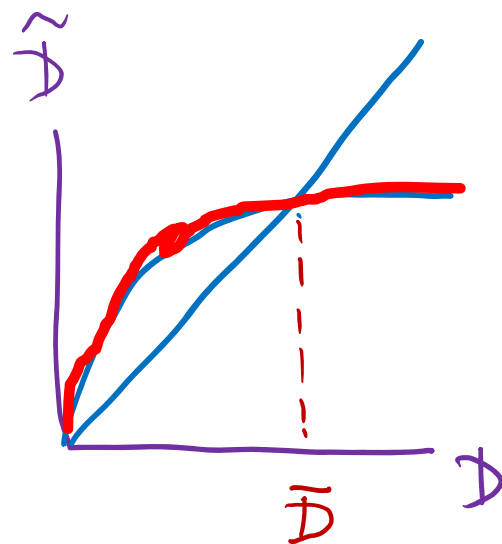
$$D_{l+1} = \xi(D_l)$$

$$\frac{d\tilde{D}}{dD} = \lambda_1 > 1$$

$$D(\mathbf{x}, \mathbf{x}') \Rightarrow \bar{D}$$

$$\bar{D} = \xi(\bar{D})$$

$$D_l = \xi * \xi * \dots * \xi(D_0)$$





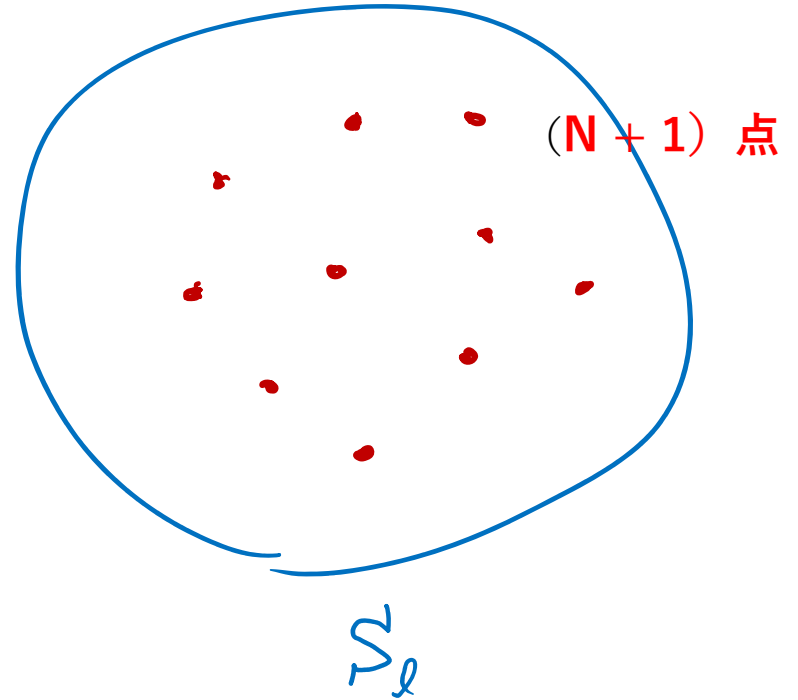
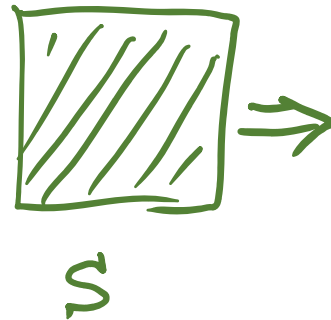
本当か!  $n \rightarrow \infty; l \rightarrow \infty$

*equidistance*

$$D(\mathbf{x}_l, \mathbf{x}'_l) \rightarrow \bar{D}$$

$$\bar{D} = \xi(\bar{D})$$

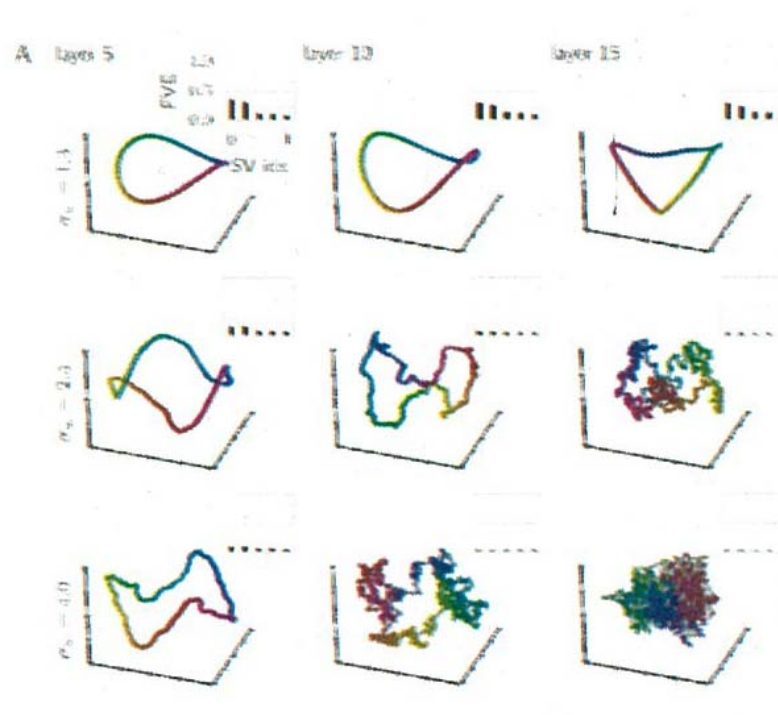
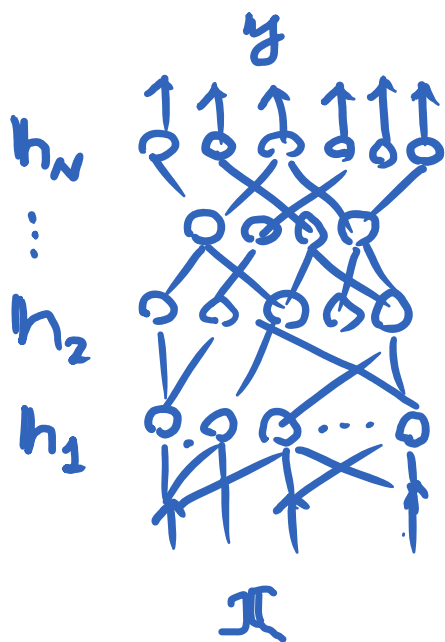
等距離



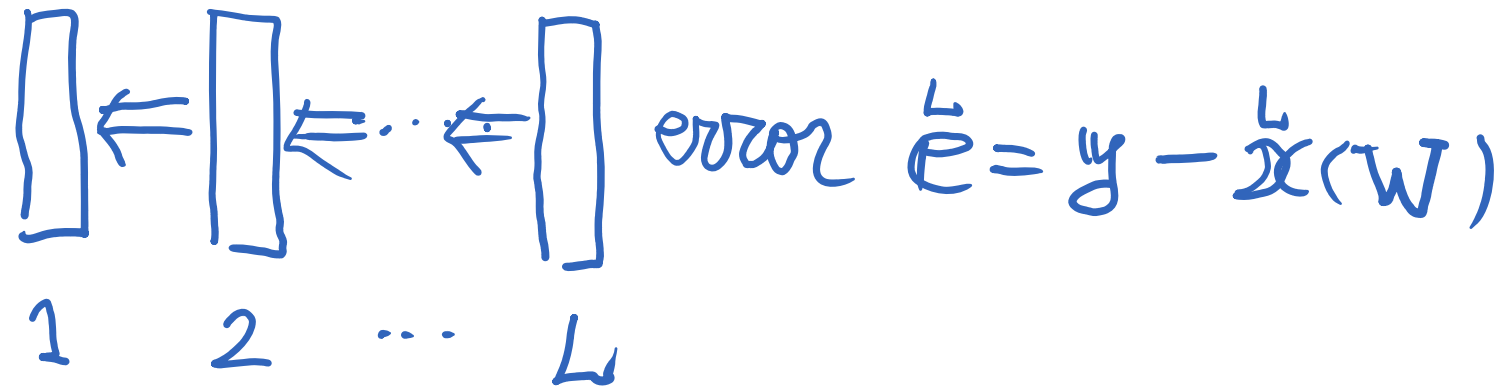
フラストレーション  
フラクタル

Poole et al (2016)  
Deep neural networks

# フラクタル 敵対的例題



# Backward path; error back-propagation Fisher information matrix



$$l(x, W) = \frac{1}{2} |y - \varphi(x; W)|^2 = |e(x, y)|^2$$

# Stochastic model : 深層回路の多様体

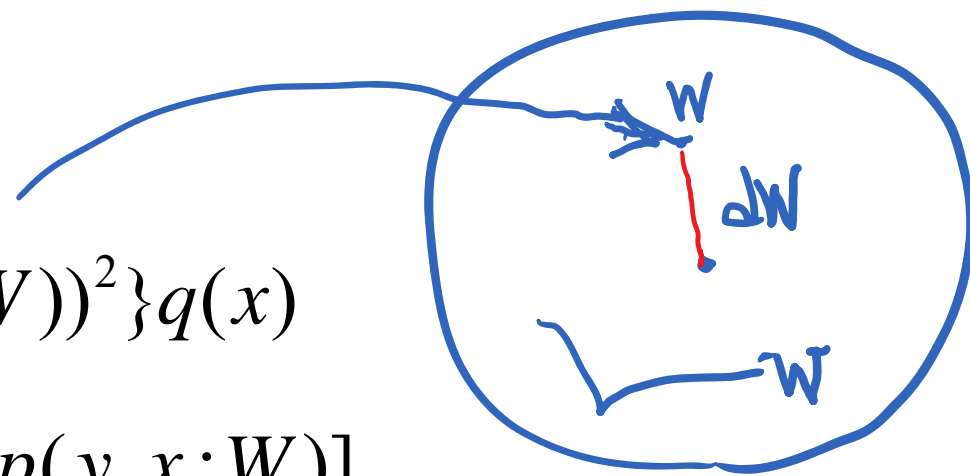
$$y = \varphi(u) + \varepsilon; \quad \varepsilon \sim N(0,1)$$

$$p(y, x : W) = c \exp\left\{-\frac{1}{2}(y - \varphi(x; W))^2\right\} q(x)$$

$$G = E_x[\nabla_W \log p(y, x : W) \nabla_W \log p(y, x : W)]$$

$$\underline{ds^2 = dW G dW}$$

Fisher information



Riemannian

# Fisher information

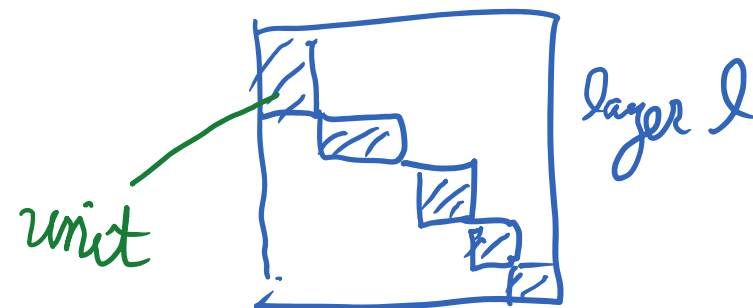
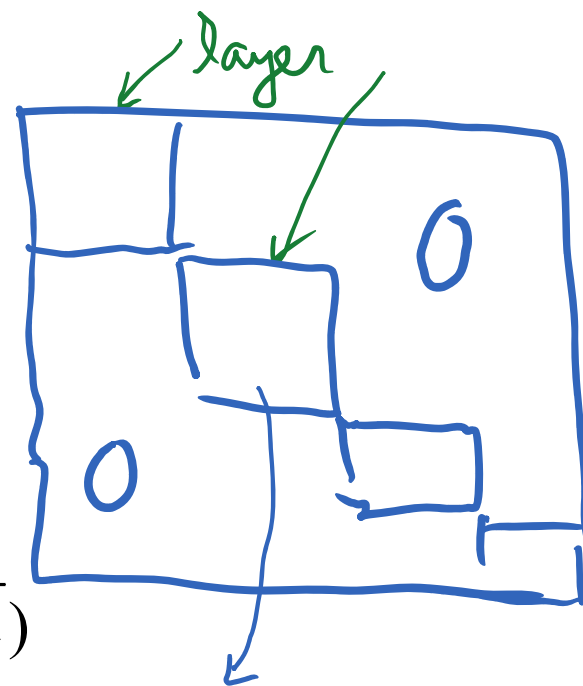
$$G = E_x \left[ \frac{\partial \phi}{\partial W_m} \frac{\partial \phi}{\partial W_l} \right]$$

$$\frac{\partial \phi^l}{\partial W_m} = \underbrace{\phi' W}_{\mathcal{B}} \frac{\partial \phi^{l-1}}{\partial W_m} = B \frac{\partial \phi^{l-1}}{\partial W_m} = \underbrace{BB \dots B}_{\mathcal{B}} \frac{\partial \phi^{m+1}}{\partial W_m}$$

$$G(W_l, W_m) = \prod \chi_1 E_x \left[ \begin{matrix} \phi' \left( \begin{matrix} l \\ \mathbf{w}_i \end{matrix} \right)^2 & & \\ & \mathbf{x} & \mathbf{x} \end{matrix} \right] + O_p(1/\sqrt{n})$$

$$G(W_l, W_m) = 0 \sim O_p(1/\sqrt{n}), \quad l \neq m$$

$$G \left( \begin{matrix} l & l \\ \mathbf{w}_i & \mathbf{w}_j \end{matrix} \right) = 0 \sim O_p(1/\sqrt{n}), \quad i \neq j$$

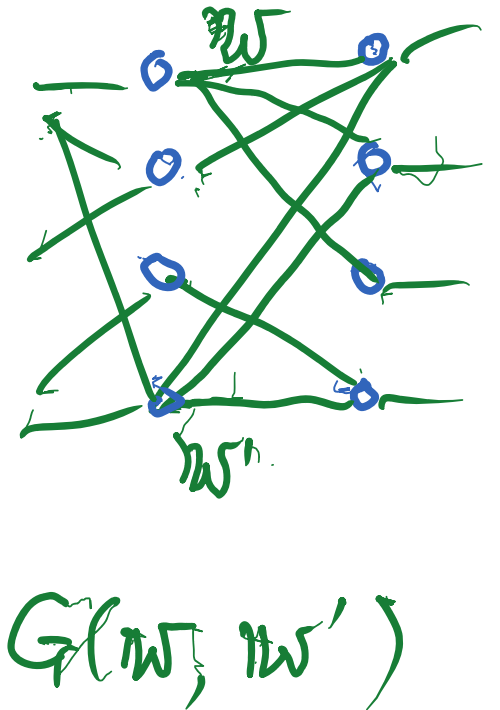


$$\mathbf{B}_{i_m}^{i_l} = \sum B_{i_{l-1}}^{i_l} B_{i_2}^{i_{l-1}} \dots B_{i_{m-1}}^{i_m} = O_p\left(\frac{1}{n}\right)$$

$$\sum \mathbf{B}_{i_m}^{i_l} \mathbf{B}_{i'_m}^{i_l} = \chi \chi \dots \chi \delta_{i_m i'_m} + O_p\left(\frac{1}{n}\right)$$

$$G(W_l, W_m) = E_x \left[ \sum \mathbf{B}_{i_l}^{i_L} \mathbf{B}_{i_m}^{i_L} \varphi' \varphi' \mathbf{xx} \right] \\ + o_p(1/n); \quad l \neq m$$

$$G(W_l, W_l) = E_x \left[ \sum \mathbf{B}_{i_l}^{i_L} \mathbf{B}_{i'_l}^{i_L} \varphi' \varphi' \mathbf{xx} \right] \\ = \delta_{i_l i'_l} E_x [\varphi' \varphi' \mathbf{xx}]_{i_{l-1} i'_{l-1}}$$



# Domino Theorem

$$\frac{\partial \mathbf{x}}{\partial \mathbf{x}} \stackrel{l}{=} B \frac{\partial \mathbf{x}}{\partial \mathbf{x}} \stackrel{l-1}{=} BB \frac{\partial \mathbf{x}}{\partial \mathbf{x}} \stackrel{l-2}{=} BBB \dots$$

$$\frac{\partial \mathbf{x}}{\partial W} \stackrel{l}{=} B \frac{\partial \mathbf{x}}{\partial W} \stackrel{l-1}{=} BB \frac{\partial \mathbf{x}}{\partial W} \stackrel{l-2}{=} BBB \dots$$

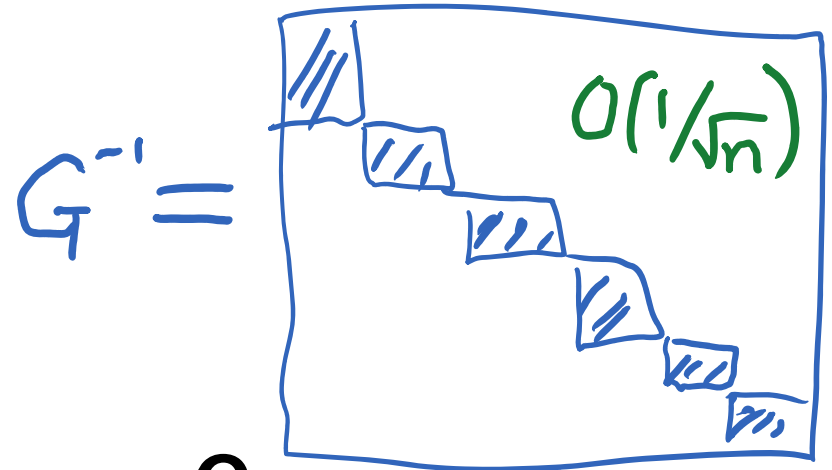
$$\Sigma \delta_{i_L i'_L} B_{i_{L-1}}^{i_L} B_{i'_{L-1}}^{i'_L} = \chi_1 \delta_{i_{L-1} i'_{L-1}} + O_p\left(\frac{1}{\sqrt{n}}\right)$$

$$\Sigma \delta_{i_L i'_L} B_{i_{L-1}}^{i_L} B_{i'_{L-1}}^{i'_L} BBBB = \chi_1 \chi_1 \chi_1 \chi_1 \delta_{i_{L-1} i'_{L-1}} + O_p\left(\frac{1}{\sqrt{n}}\right)$$

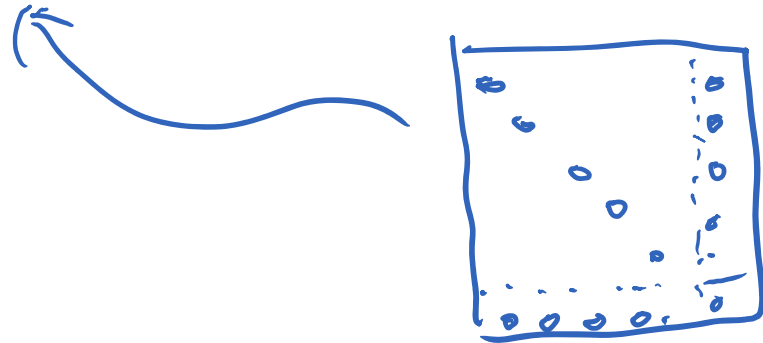


# Unitwise natural gradient

$$\Delta W = -\eta G^{-1} \nabla_W l$$



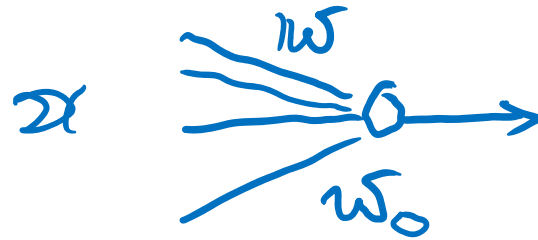
Y. Ollivier; Marceau-Caron



# Fisher information of unit $\tilde{\mathbf{w}} = (w, w_0)$

$$y = \varphi(\mathbf{w} \cdot \mathbf{x} + w_0)$$

$$G = E_x [\partial_w \varphi \partial_w \varphi] = E_x [(\varphi')^2 \mathbf{xx}]$$



$$G(\tilde{\mathbf{w}}, \tilde{\mathbf{w}})$$

$\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\} \Rightarrow \{\mathbf{e}_1^*, \mathbf{e}_2^*, \dots, \mathbf{e}_n^*\}$ : ortho-normal basis

$$\mathbf{e}_n^* = \frac{\mathbf{w}}{w}$$

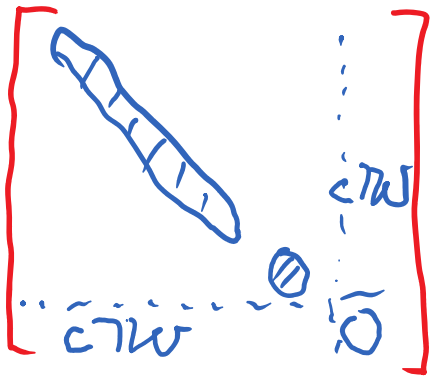
Input  $x$ : independent and identically distributed, 0-mean

$$G = A\mathbf{I} + \frac{B}{w^2} \mathbf{w}\mathbf{w} + \frac{C}{w} (\mathbf{w}\mathbf{b} + \mathbf{b}\mathbf{w})$$

$$\mathbf{b} = [0 \ 0 \ \dots \ 0 \ 1]'$$

$G^{-1}$  : similar form

$A, B, C$   
                      
 $(w, b)$



$$\Delta \mathbf{w} = -\eta e G^{-1} \mathbf{x}$$

TANGO  
QD自然勾配法?

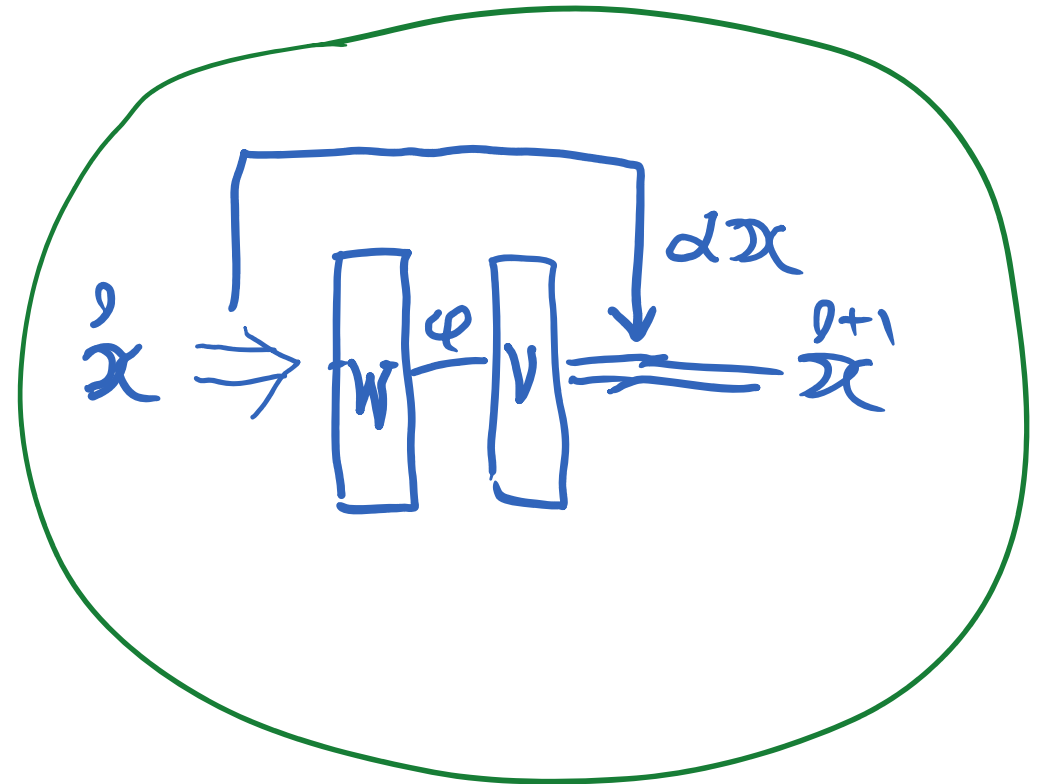
$$= -\eta e \left[ \bar{A} \mathbf{x} + \left( \frac{\bar{B}}{w^2} \mathbf{w} \cdot \mathbf{x} + c \right) \mathbf{w} \right]$$

$$\Delta w_0 = -\eta e D$$

# Lesnet

$$\mathbf{x}^l = V \varphi \left( W \mathbf{x}^{l-1} \right) + \alpha \mathbf{x}^{l-1}$$

$$\chi_1 \rightarrow \sigma_v^2 \chi_1 + \alpha^2$$



# Karakida theory

eigenvalues of  $G$

$$\frac{1}{P} \sum \lambda_i = \frac{1}{n}, \quad \frac{1}{P} \sum \lambda_i^2 = O(1)$$

distorted Riemannian metric

