

統数研記念式典

人工知能の歴史、発展、 社会への影響

甘利俊一

理化学研究所 荣誉研究員
東京大学 名誉教授

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—並列分散(ニューラルネット)

動力技術

機械技術

材料技術

生命技術—遺伝子編集

情報技術—知能・心

技術は止まらない！

脳一人間とは何か

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

・ビッグバン(138億年前)

物質の法則

エネルギー・物質－天体－分子(秩序)
物理学、化学



46億年前：
地球の誕生

火の玉、全球凍結

生命の誕生（36億年前）

進化の法則

生命 = 情報 + 物質

=自己を複製し次世代に伝える物質

生命科学

遺伝、分子機構、自己保存



DNA



多細胞生物

環境の情報を利用、
記憶・学習、判断・行動

脳=情報

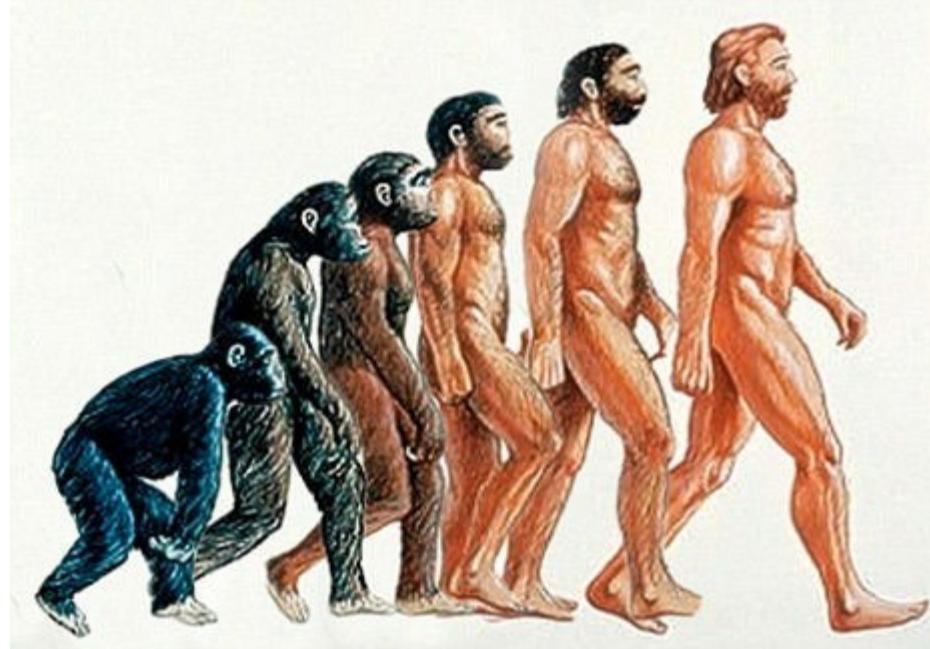
脳・神経系(5億年前)：神経科学



類人猿、そして人間 社会に生きる生命；文化と社会



人類の登場(700万年前)



心、意識
文明

猿人、原人

旧人(ネアンデルタール、50万年-3万年前)

新人(ホモ・サピエンス、20万年前-現在)

心の理論 ロボットに意識はあるか



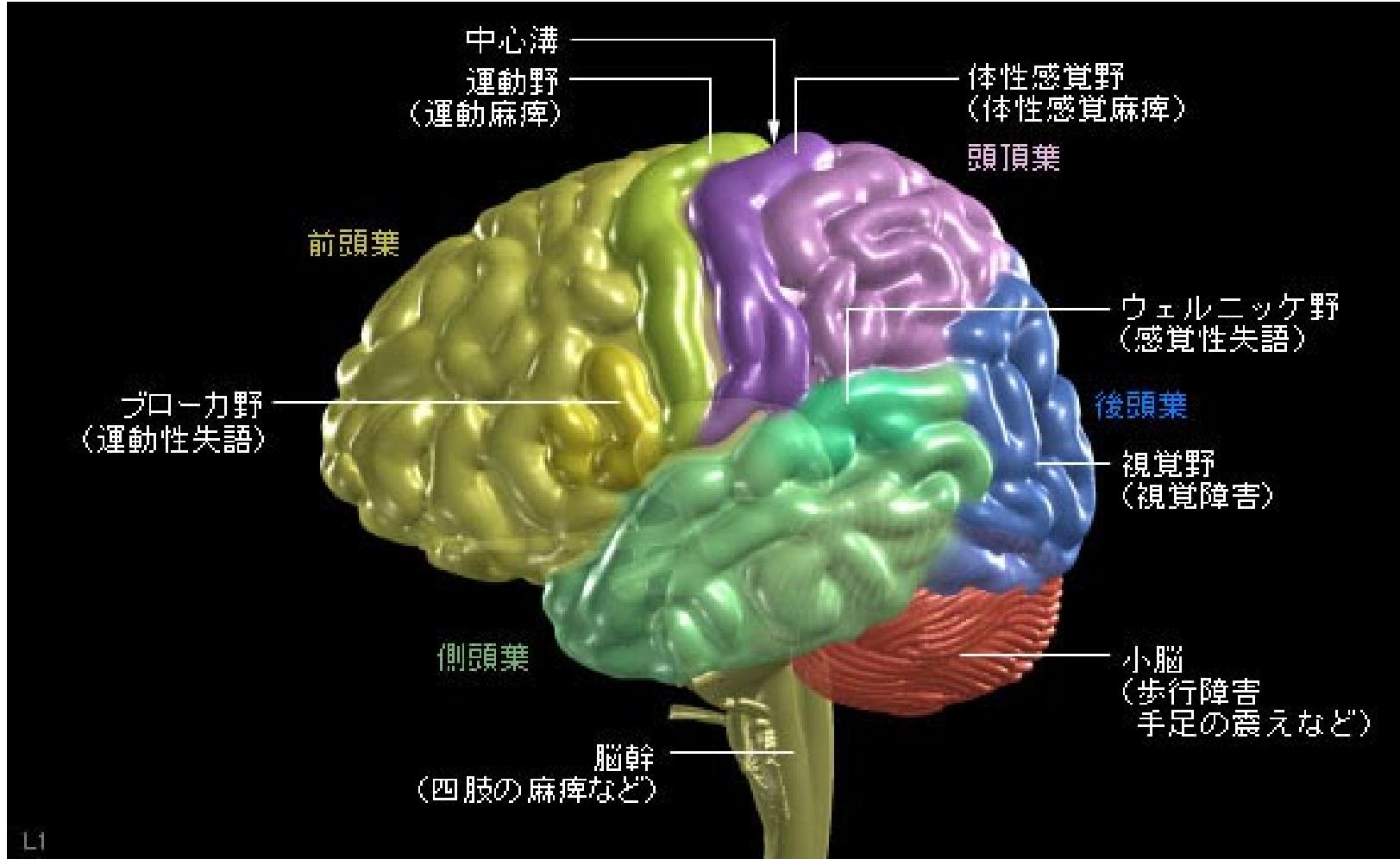
物質の法則： 宇宙

生命の法則： 情報+物質： 進化

文明の法則： こころ+情報+物質：
社会・文化

脳：大脳、海馬、小脳、脳幹

脳科学：ミクローマacro、理論：神経回路網



人工知能と脳のモデル： —歴史の要約

第一次ブーム

1956~ AI
Dartmouth 会議
記号と論理
知的推論、ゲーム

実用的でない!!

暗黒期 (1965後半~1970's)

stochastic descent learning (1967) for MLP

脳モデル
Perceptron
学習する普遍計算機構
線形分離可能

第2次ブーム

1970~

AI

エキスパートシステム
(MYCIN, DENDRAL)

沈静化

確率Bayes推論
chess (1997)

1980~ BT (神経回路)

MLP (backprop)

連想記憶モデル、ダイナミックス
—兆円産業か？

第3次ブーム 2010~ 脳型の人工知能(融合)

深層学習 Deep learning

(畠み込み多層回路(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 (graphical model; Bayesian; WATSON)

深層学習の勝利 ——人間以上の識別能力

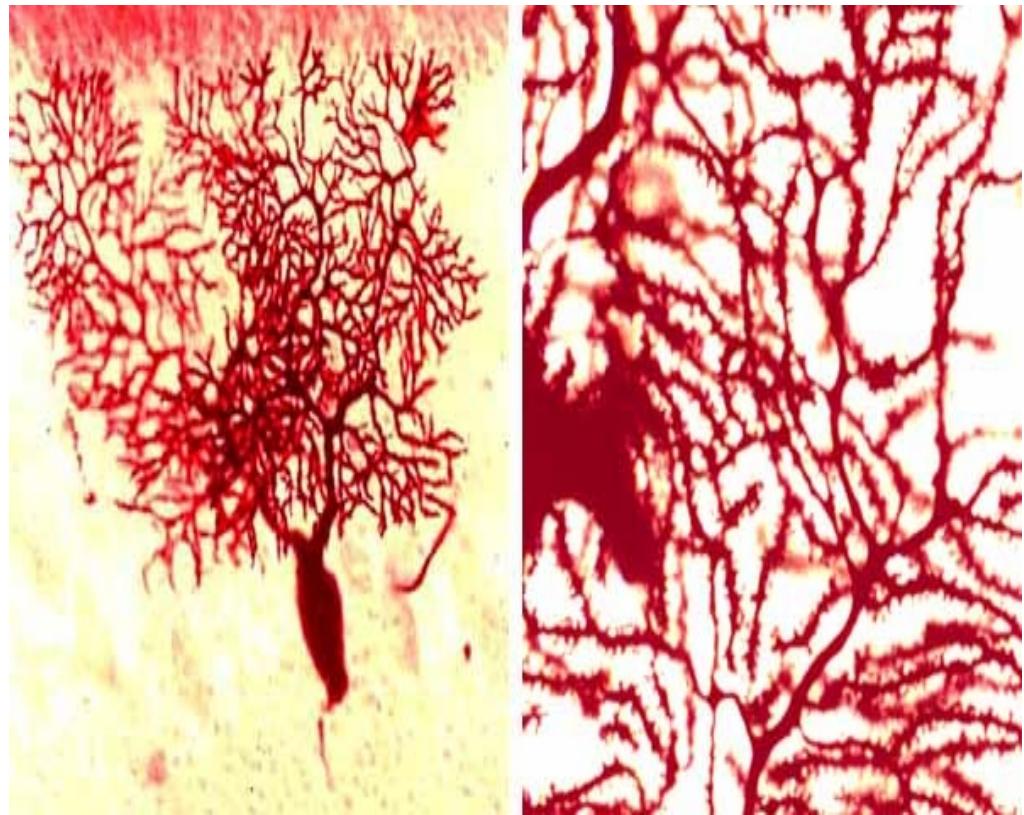
パターン認識: vision, auditory, sentence analysis

囲碁: 強化学習

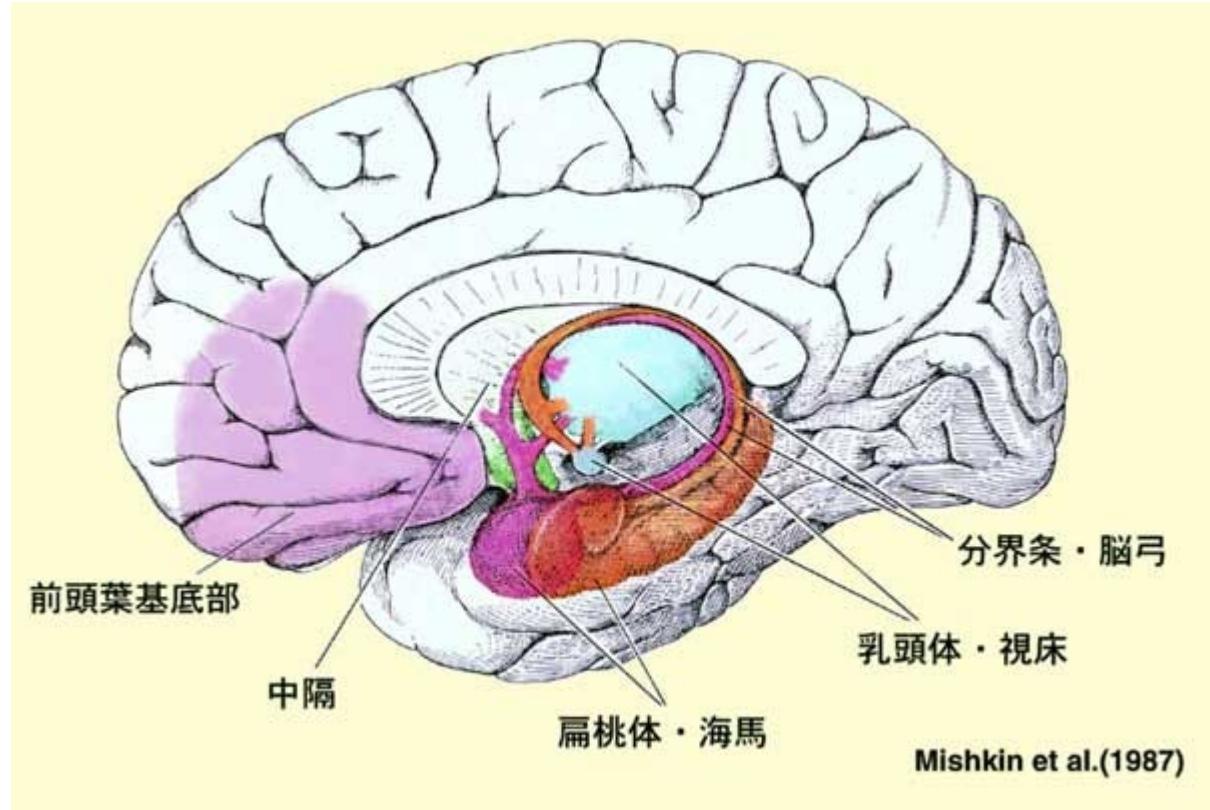
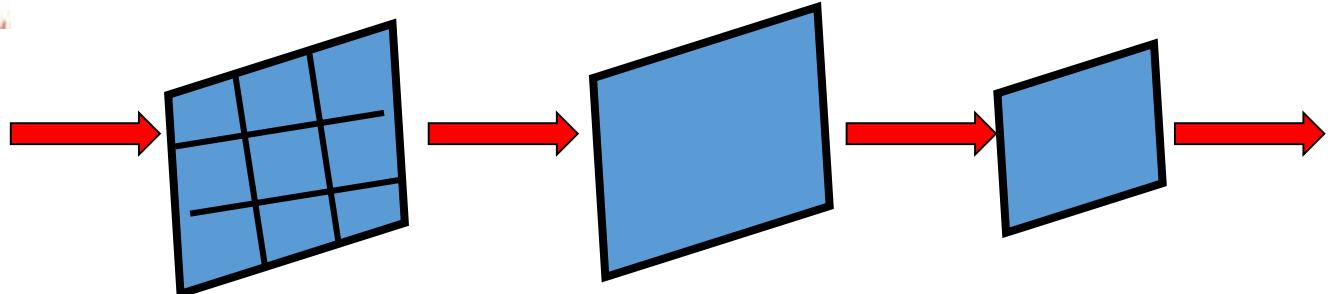
時系列とダイナミックス、動的パターン; 言語処理

記号と論理 VS パターンとダイナミックス、学習—融合

ニューラルネットと脳

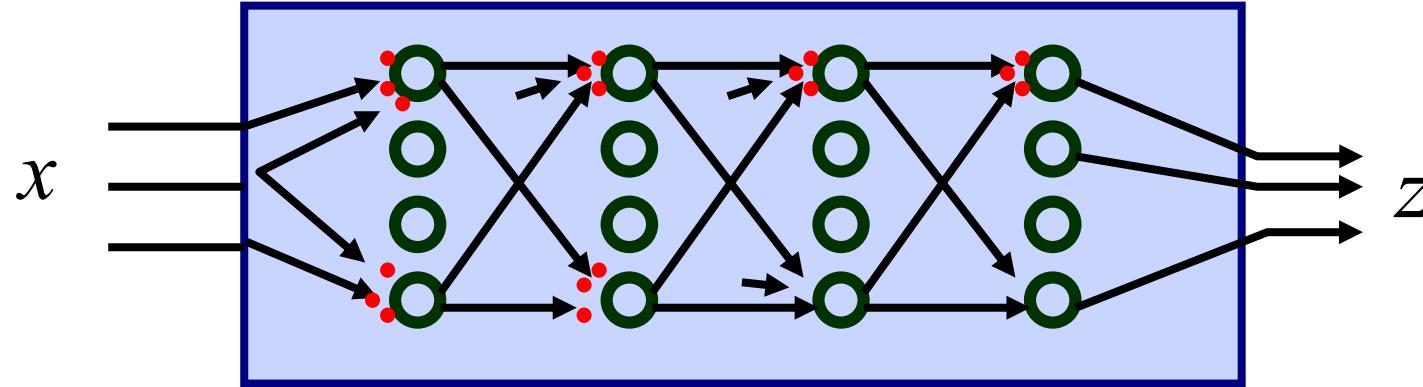


多層パーセプトロン
表現 ダイナミックス 万能性



層状学習回路網

multilayer perceptron



パーセプトロン Perceptron

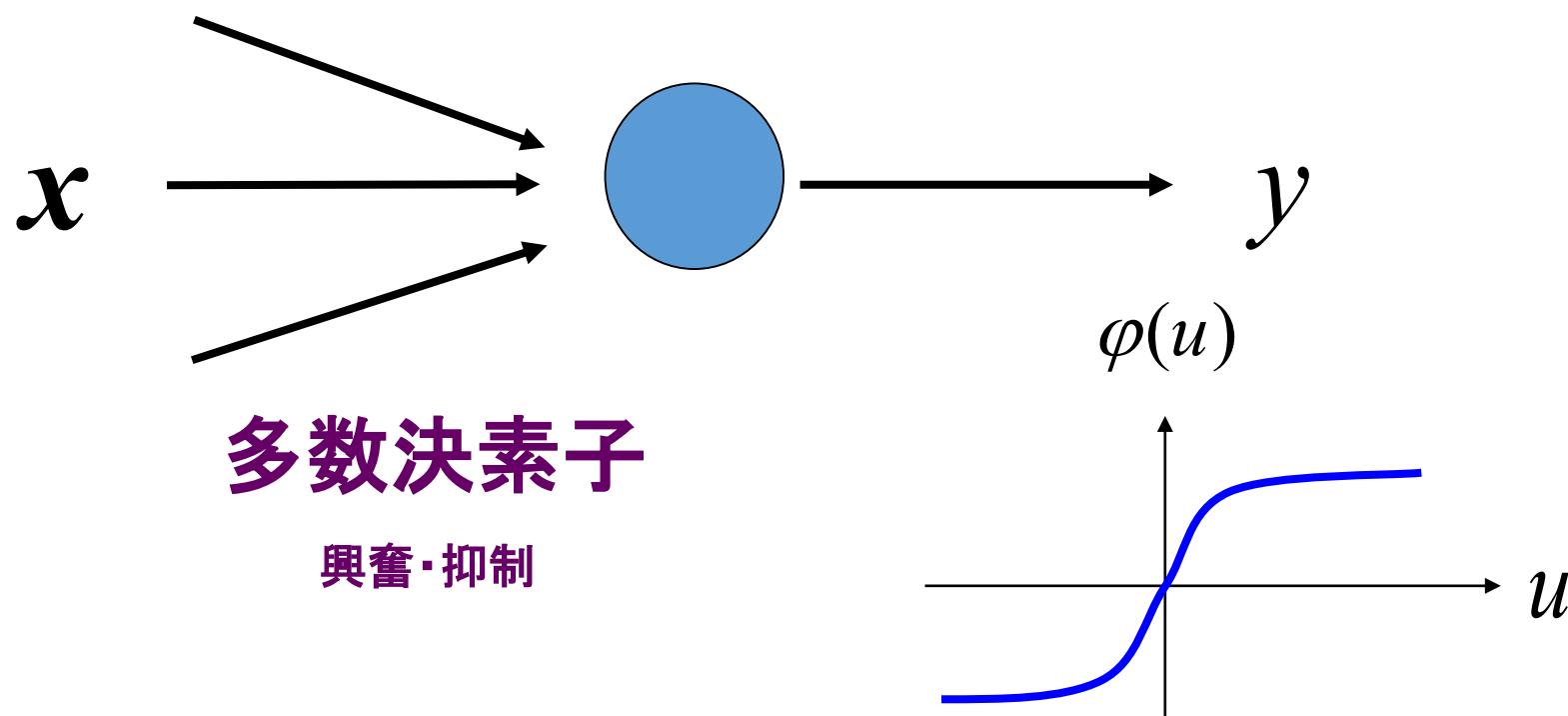
バックプロパゲーション Backpropagation

$$L(x, W) = |y - g(x, W)|^2$$

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

ニューロンの数理モデル

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



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 晃一

南雲 仁一

中村 幸雄

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

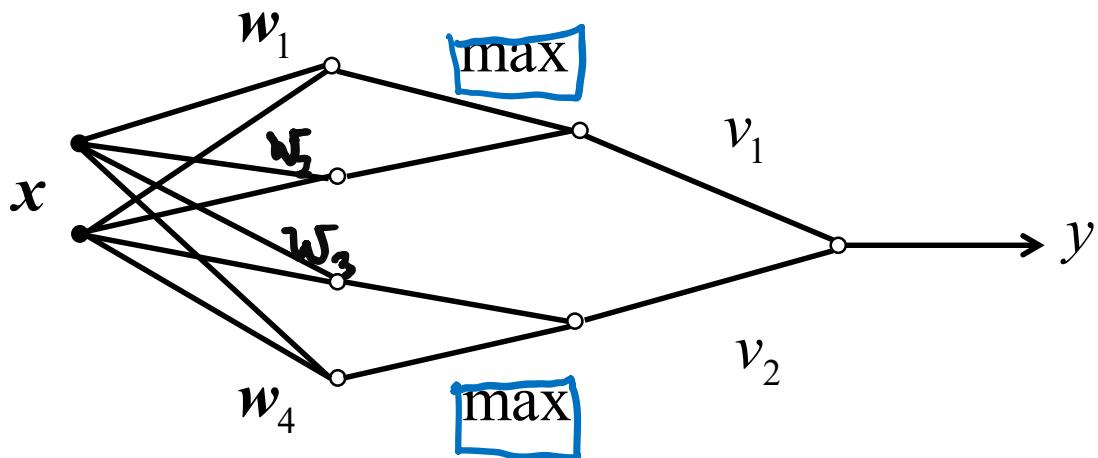
Shun-ichi Amari
University of Tokyo

Kyouritu Press, Tokyo, 1968

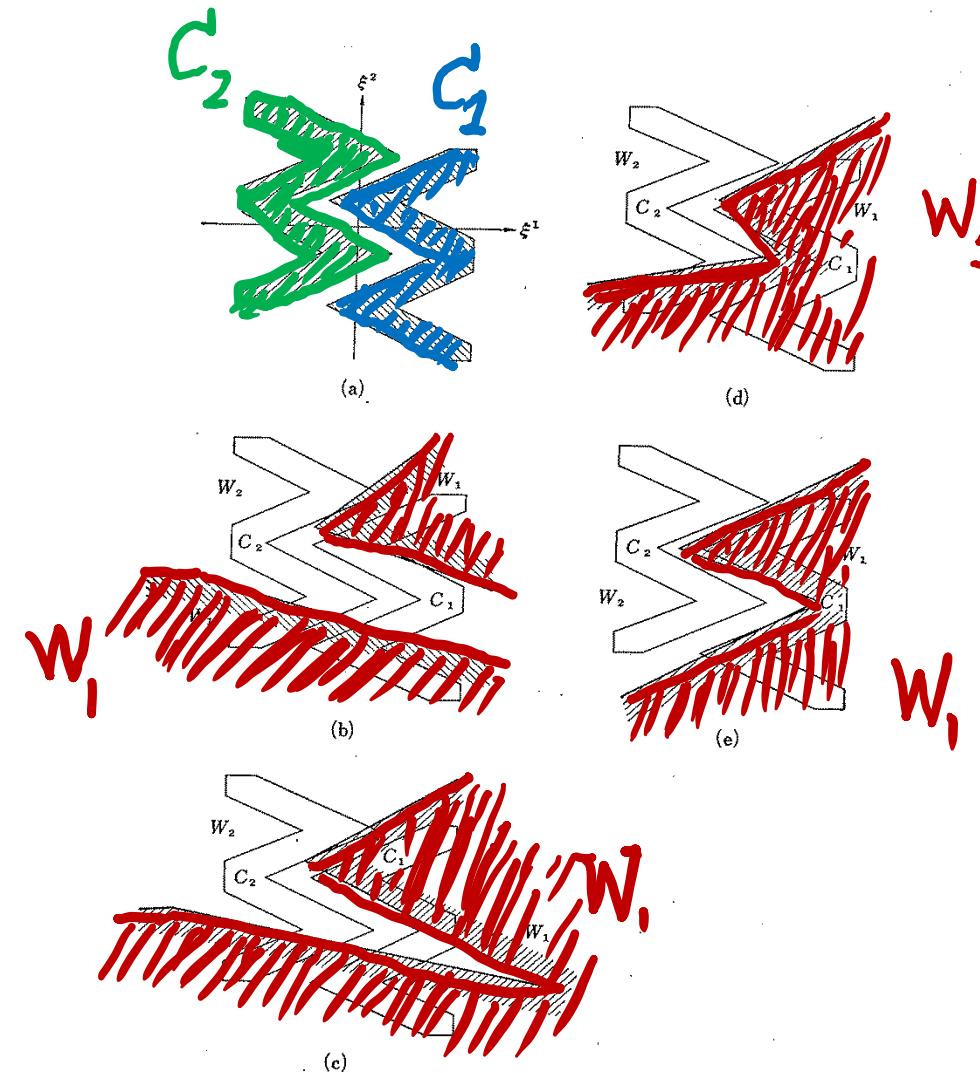
$$f(x, \theta) = v_1 \max \{w_1 \cdot x, w_2 \cdot x\} + v_2 \min \{w_3 \cdot x, w_4 \cdot x\}$$

1

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



深層学習の勝利

パターン認識(画像、音声,...)

囲碁・ゲーム

テキスト生成

言語翻訳

深層学習の問題点 1

大量のデータ、計算力： 実験式：原理を発見しない

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(新概念)

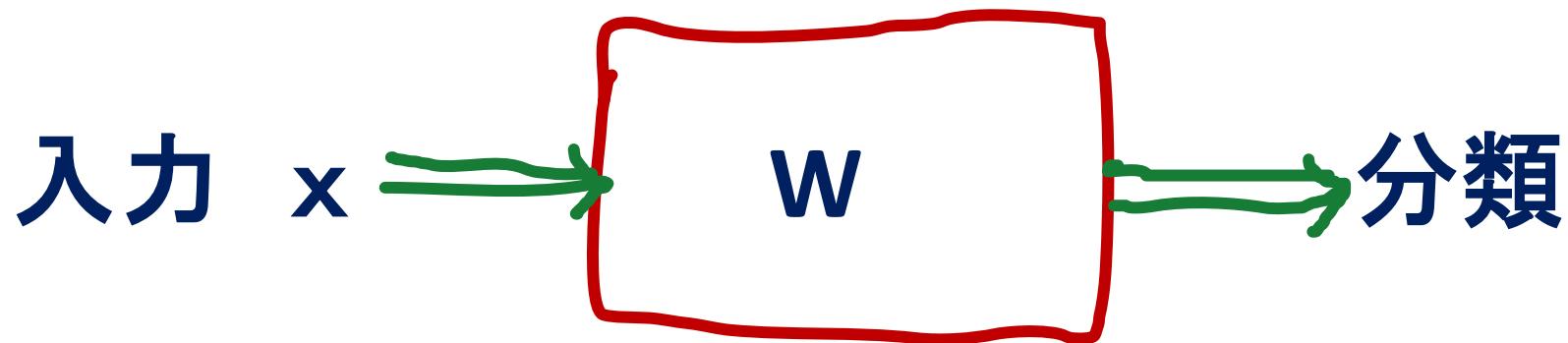
原理の創出・理解一人間

科学

深層学習の問題点 2

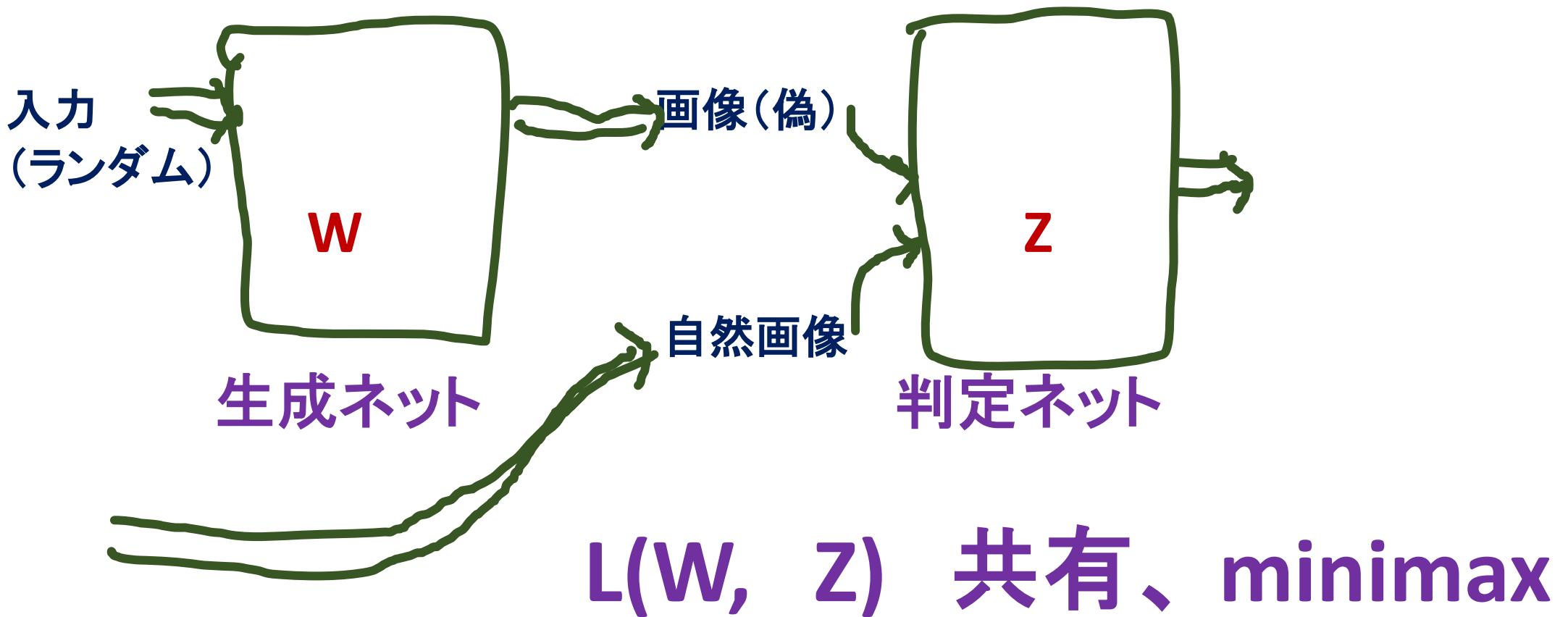
何故1000層も必要か 情報表現
敵対的例題と脆弱性

敵対的例題



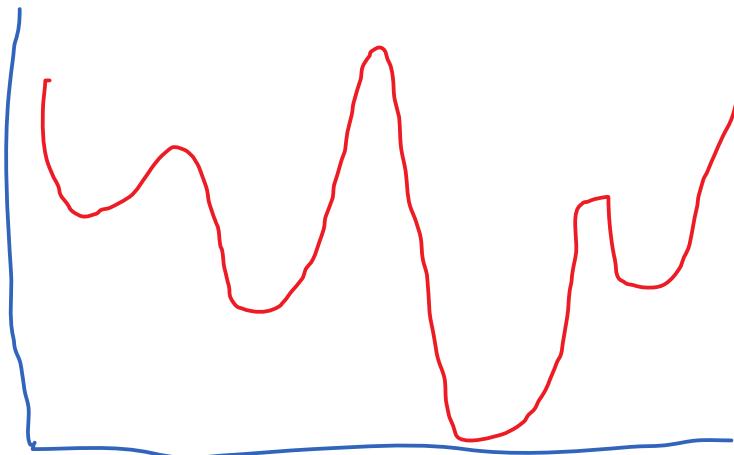
x の学習

敵対的ネットワーク：真偽判定回路



深層学習の問題点 3

局所解と大域解



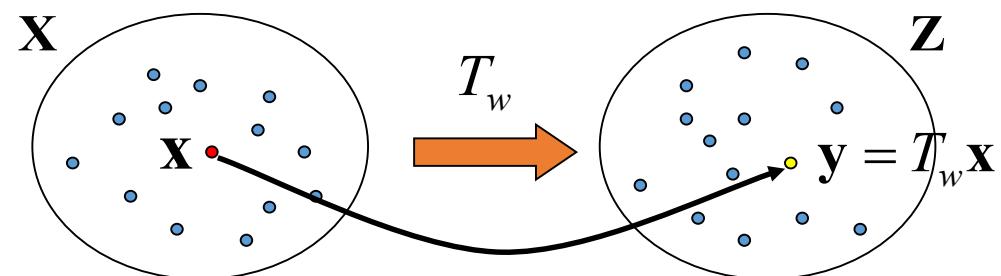
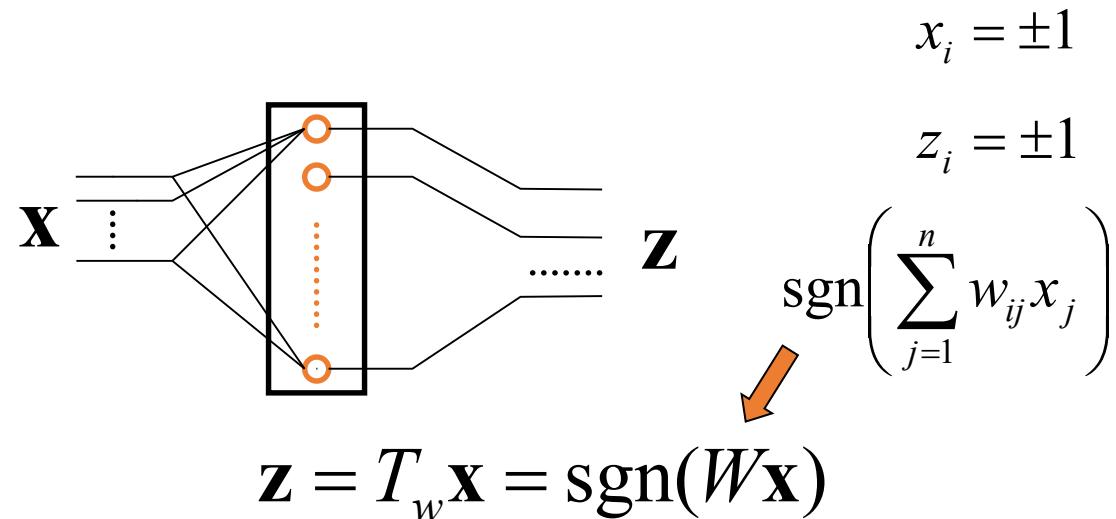
大規模系の特徴

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

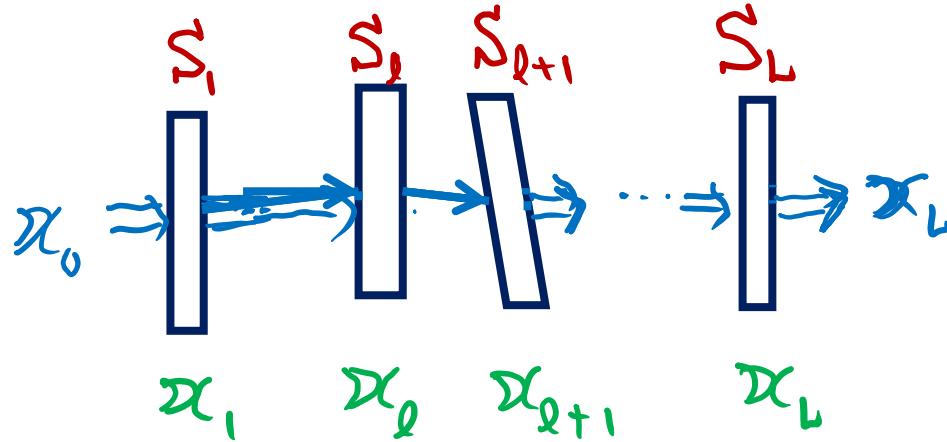
統計神経力学

1-layer network



深層回路

$$x_i = \varphi \left(\sum_{l+1} w_{ij} x_i + w_{0i} \right)$$



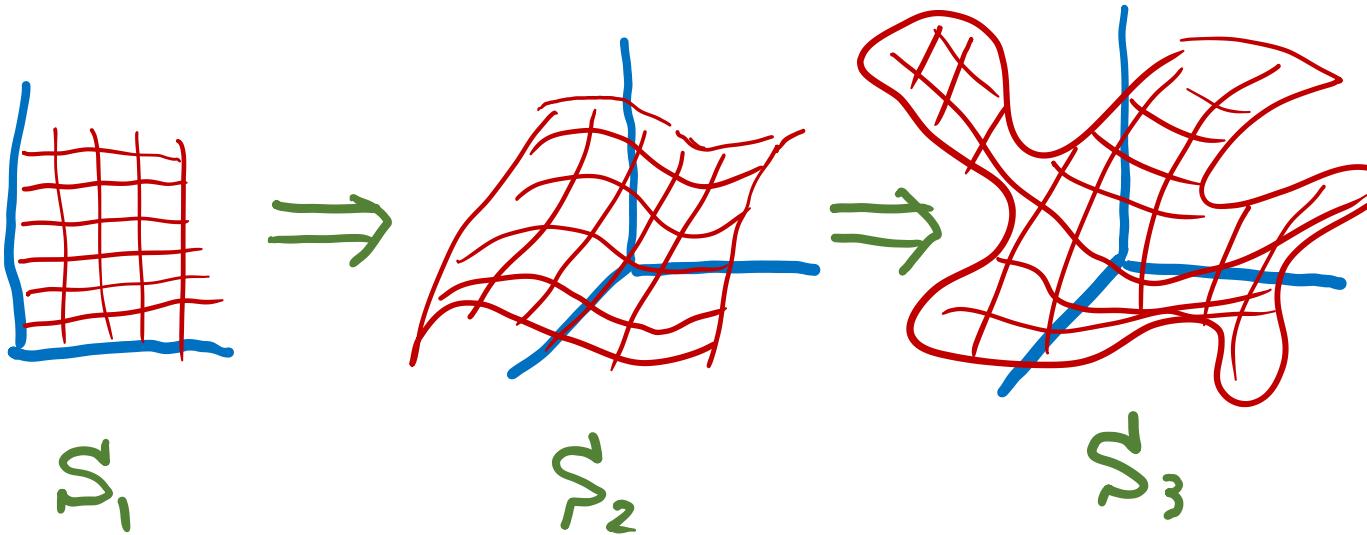
$$A_l = \frac{1}{n_l} \sum_i {x_i}_l^2$$

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

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

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

引き戻し計量(リーマン計量・距離)



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

$$g^l{}_{ab} = \mathbf{e}^l{}_a \cdot \mathbf{e}^l{}_b$$

リーマン計量の力学

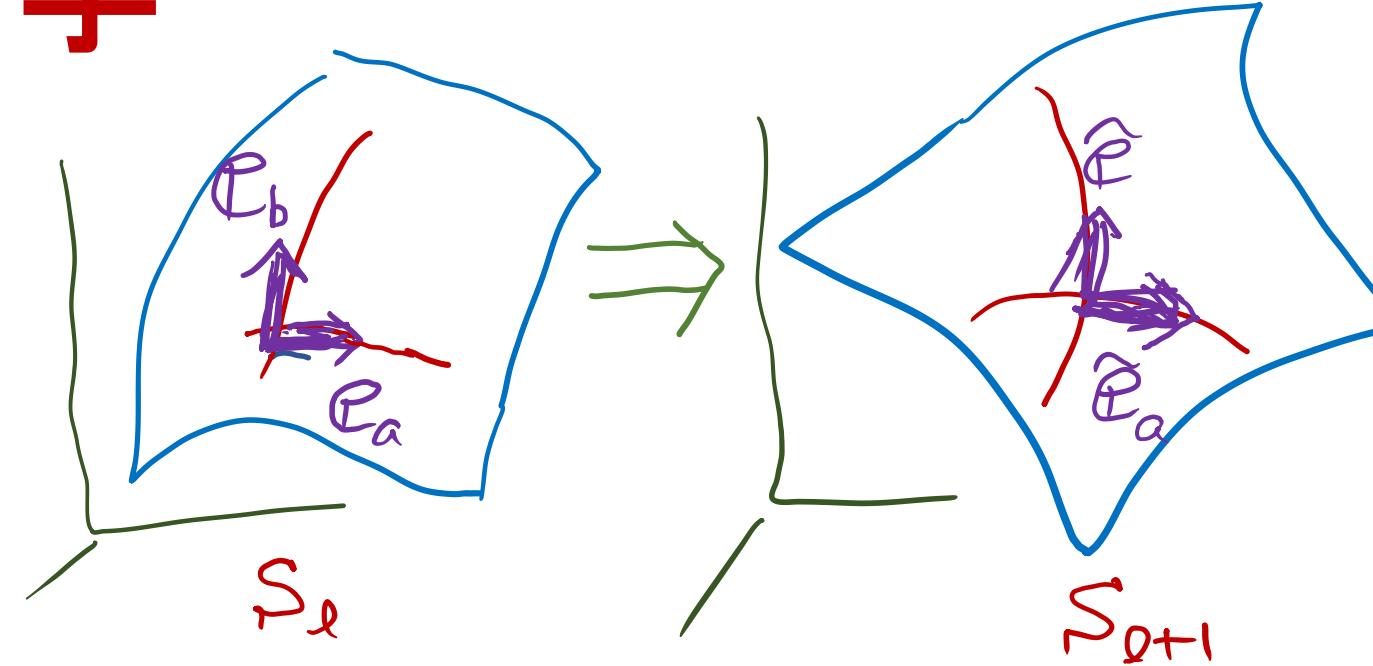
$$\tilde{y}_\alpha = \varphi(\sum w_{\alpha k} y_k + b_\alpha) = \varphi(u_\alpha)$$

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

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

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

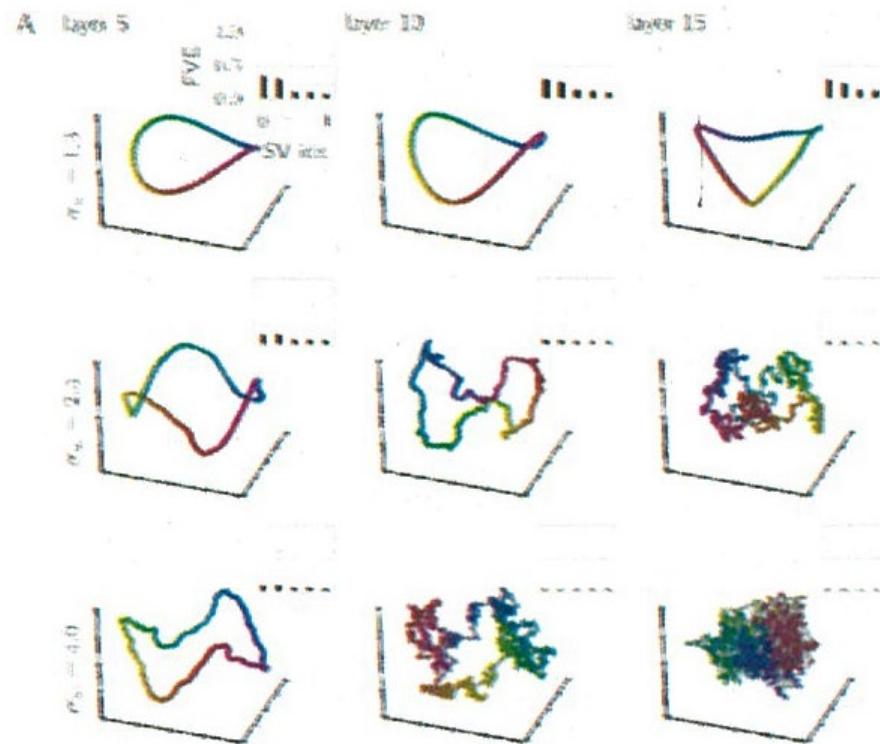
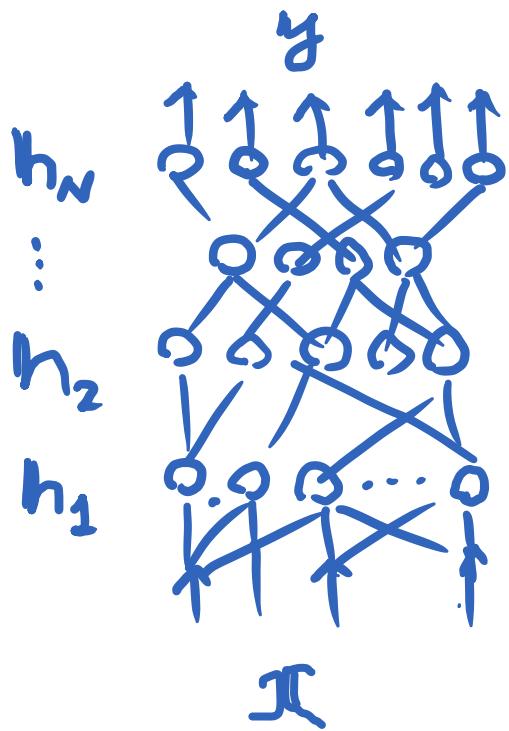
$$\text{E}[\varphi'(u_\alpha))^2 w_k^\alpha w_j^\alpha] = \text{E}[\varphi'(u_\alpha))^2] \text{E}[w_k^\alpha w_j^\alpha]$$



平均場近似不使用

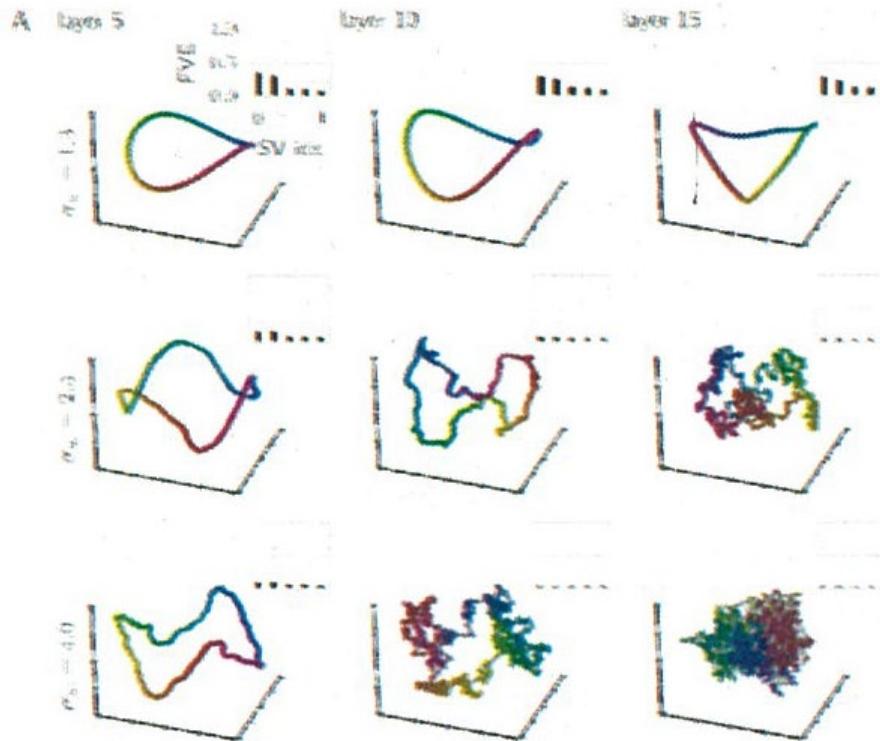
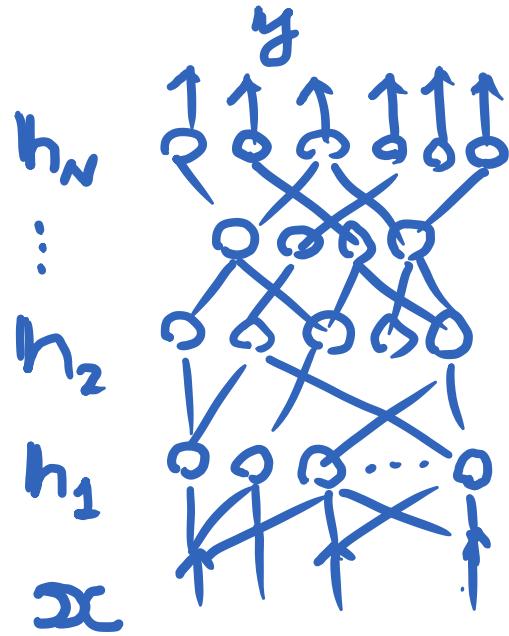
$$\chi_1(A) = \int \sigma^2 \{\varphi'(\sqrt{A}v)\}^2 Dv = \frac{1}{2\pi} \frac{\sigma^2 A + \sigma_b^2}{\sqrt{1 + 2(\sigma^2 A + \sigma_b^2)}}$$

フラクタル 敵対的例題



Poole et al (2016)

Random deep neural networks



Natural Gradient

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

$$|d\theta|^2 = \varepsilon \quad \text{KL}[p(x, \theta) : p(x, \theta + d\theta)] = \varepsilon$$

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

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

Fisher information

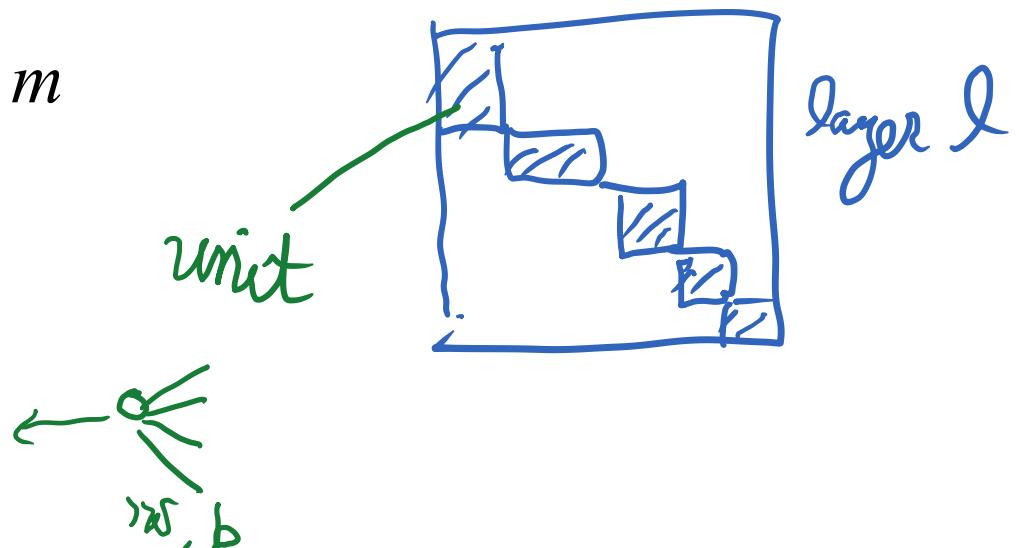
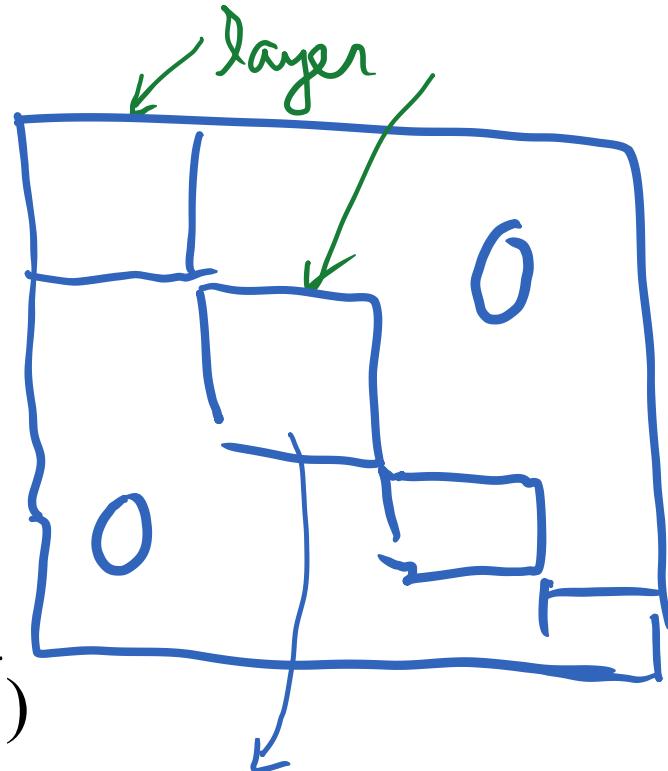
$$G = E_x \left[\frac{\partial \varphi}{\partial W_m} \frac{\partial \varphi}{\partial W_l} \right]$$

$$\frac{\partial \varphi^l}{\partial W_m} = \varphi' W \frac{\partial \varphi^{l-1}}{\partial W_m} = B \frac{\partial \varphi^{l-1}}{\partial W_m} = BB...B \frac{\partial \varphi^{m+1}}{\partial W_m}$$

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

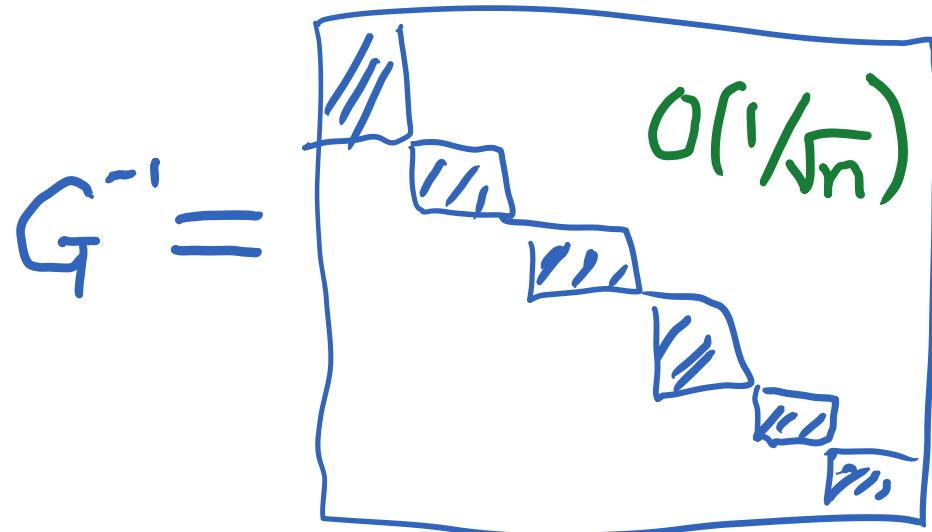
$$G(W_l, W_m) = 0 \sim O_p(1/\sqrt{n}), \quad l \neq m$$

$$G(\mathbf{w}_i, \mathbf{w}_j) = 0 \sim O_p(1/\sqrt{n}), \quad i \neq j$$

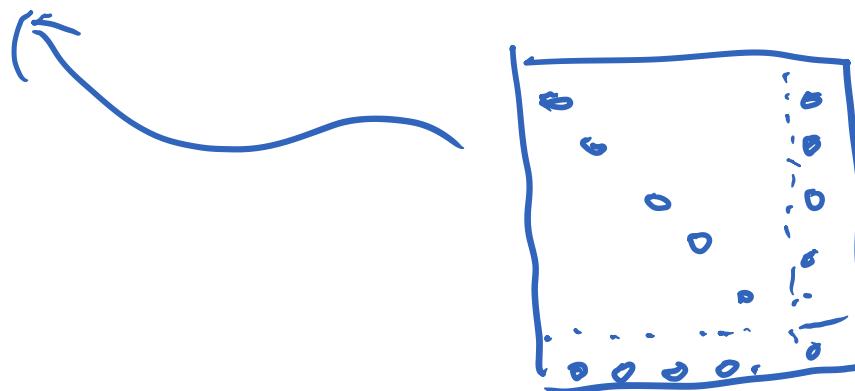


Unitwise natural gradient

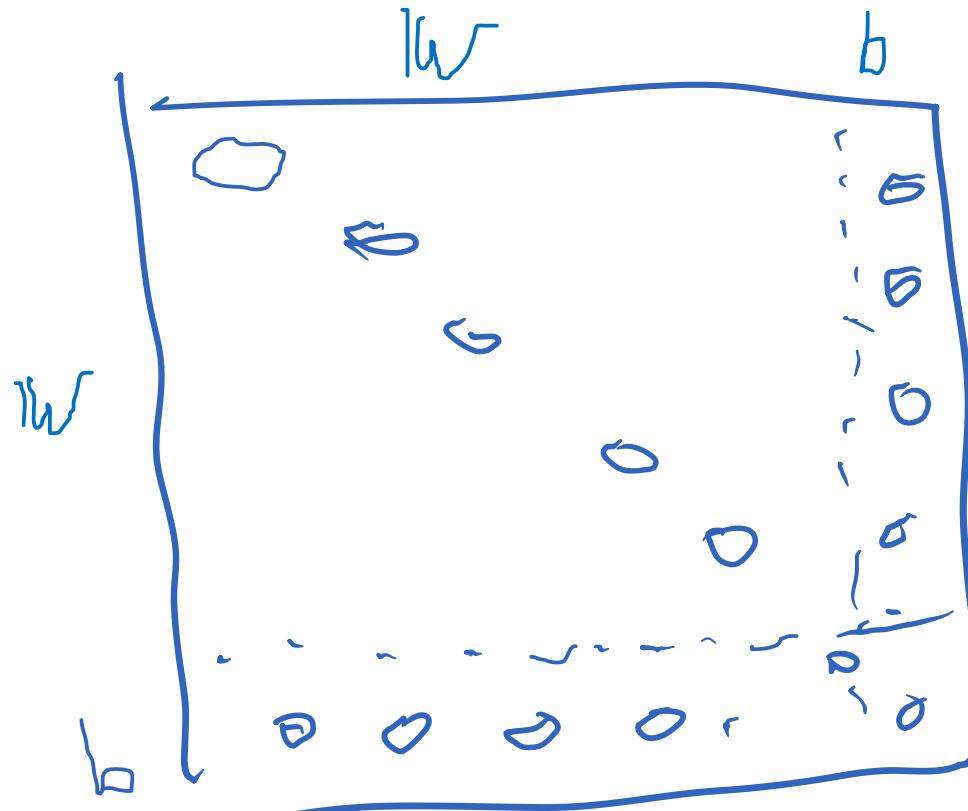
$$\Delta W = -\eta G^{-1} \nabla_W l$$



Y. Ollivier; Marceau-Caron



G=



+WW

Ollivier, Marceau-Caron : quasi-diagonal natural gradient

No simulations: Oh No!!

Natural gradient is no more dream!

TANGO: Ollivier
recursive formula for

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

Neural tangent kernel

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)

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

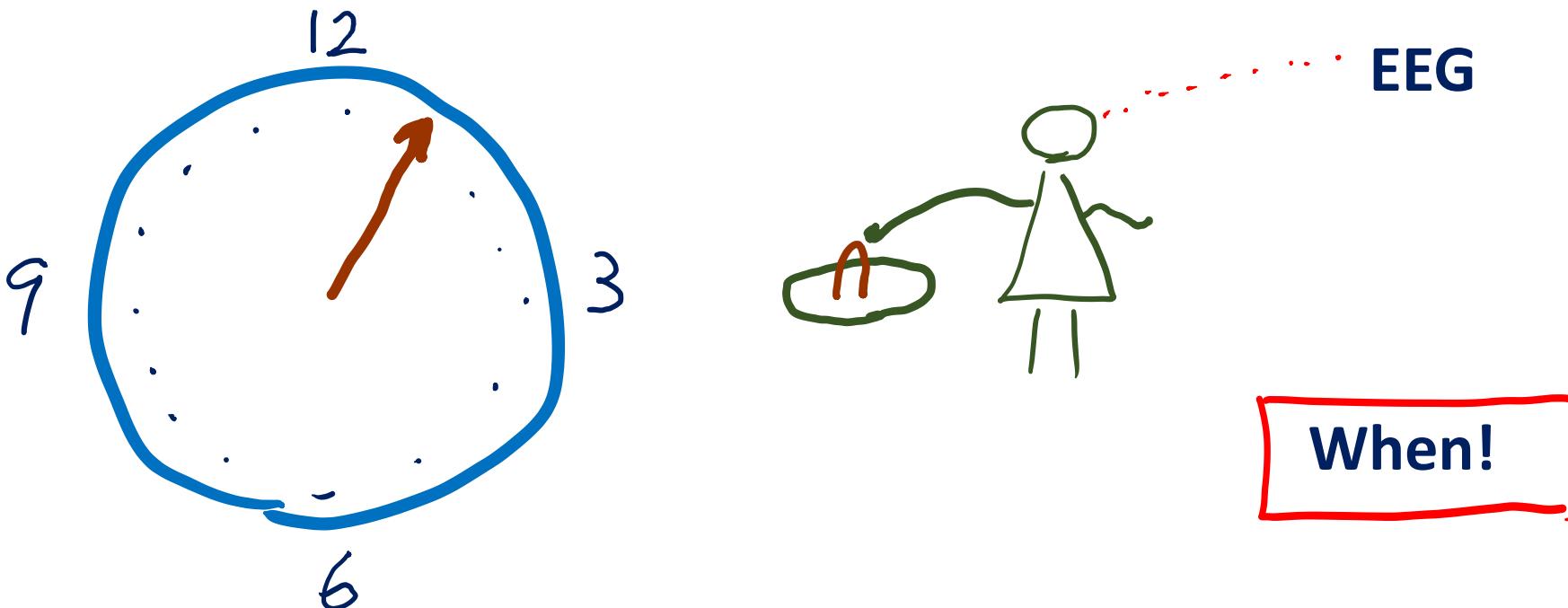


意識の発生

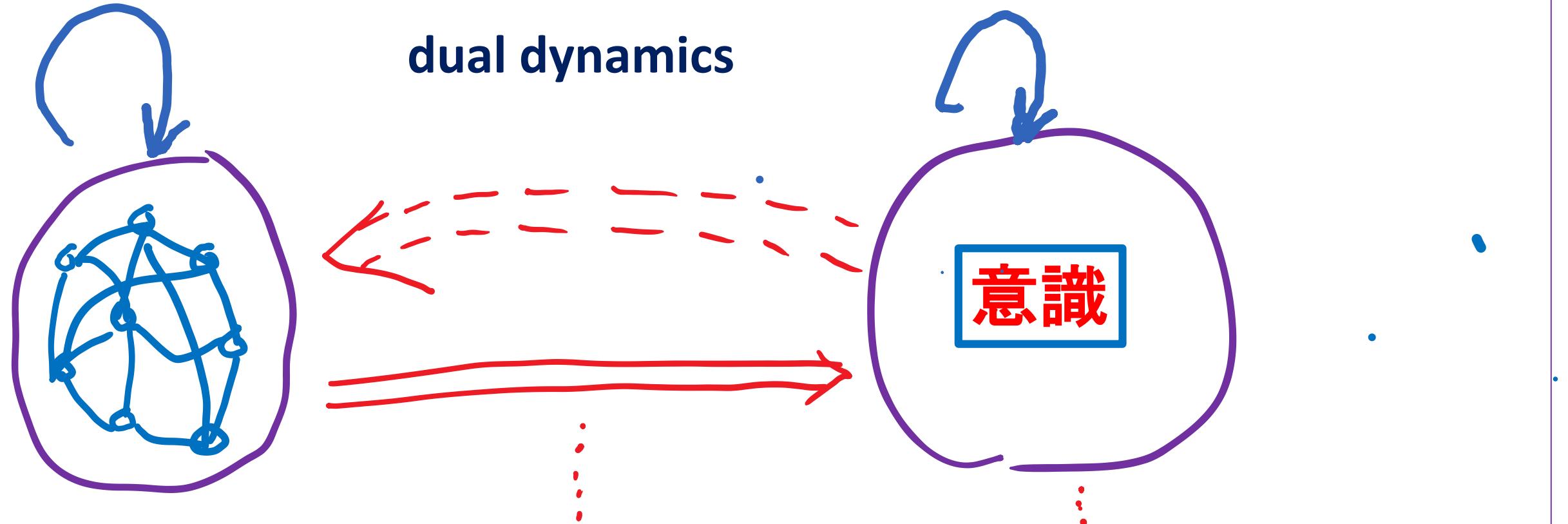
共同作業、自分の意図を自分で知る

言語： 論理的思考、数学

Libet の実験：自由意志



予測(先付け)と後付け Prediction and Postdiction



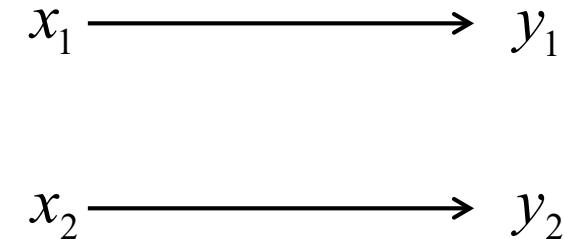
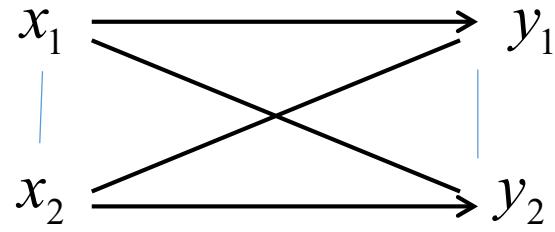
ダイナミックス

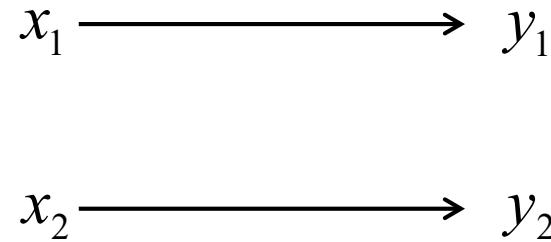
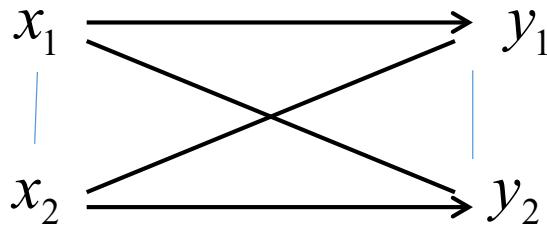
意思決定と行動

反省、正当化、論理

意識の情報統合理論 Tononi IIT

システムの情報理論——情報幾何





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

Disconnected model:

$$S_{dis} = \{q(\mathbf{x}, \mathbf{y})\} \quad q(\mathbf{y} \mid \mathbf{x}) = \prod q(y_i \mid x_i)$$

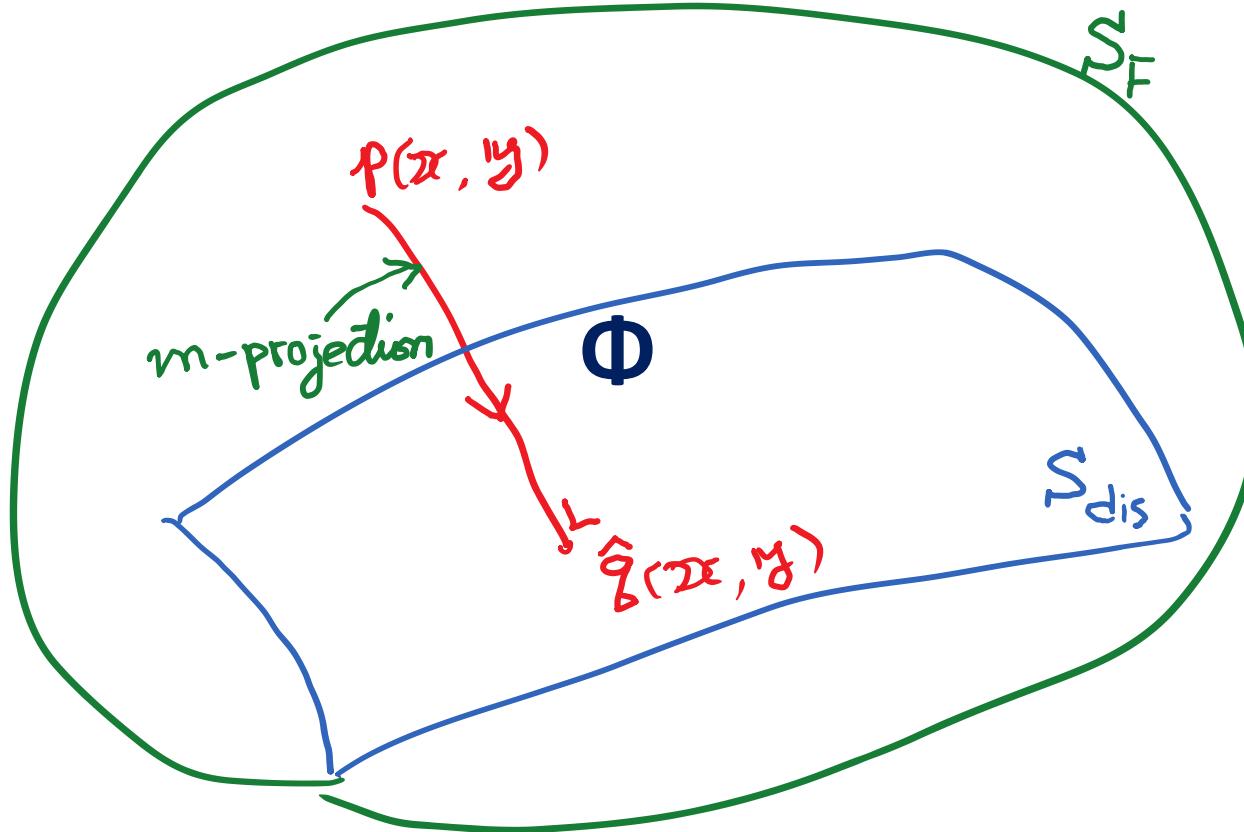
measure of interaction : N. Ay

information integration : Tononi

Barrett and Seth

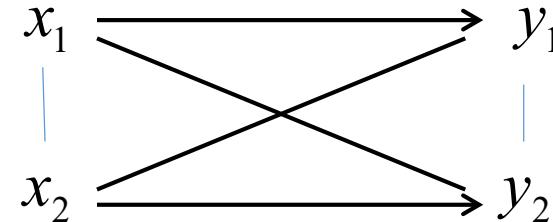
Many other Φ

Integrated Information Theory G. Tononi



Necessary condition; sufficient?

Markov Condition



(1→2) branch deleted: **Markov condition:** $x_1 \rightarrow x_2 \rightarrow y_2$

$$p(x_1, y_2 \mid x_2) = p(x_1 \mid x_2)p(y_2 \mid x_2)$$

$$X_1 - X_2 - Y_2$$

$$X_2 - X_1 - Y_1$$

$S_{dis} :$ all $x_i \rightarrow y_j$ ($i \neq j$) deleted

人工知能が脳に学ぶべきこと：数理的理解

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

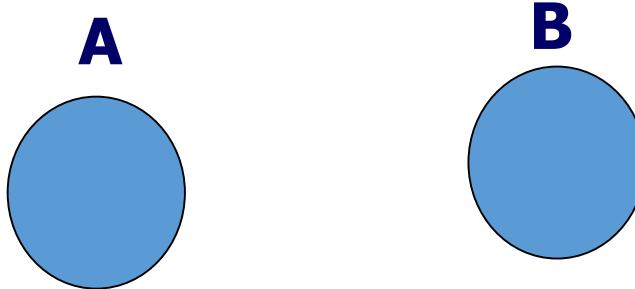
心の理論



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

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：使命感

人間は不合理； 芸術、喜び、愛、苦悩

ただ一度の、かけがいのない人生

ロボット（金融システム）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: 技術は止まらない: 制御できるか

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

我々はは何をなすべきか？

日本のAIの進むべき道: 政府の戦略 ブームは終わる



超大国 文化国家

物量作戦はだめ
理論とアイデア
中小企業を含む現場との交流; 産業の情報化