

文科省・情報科学技術分野における
戦略的重要研究開発領域に関する検討会

資料9
科学技術・学術審議会 情報委員会
情報科学技術分野における戦略的重要研究
開発領域に関する検討会(第1回)
令和6年4月24日

2024年4月24日

AIのサイエンス

杉山 将



理化学研究所 革新知能統合研究センター
／東京大学



<http://www.ms.k.u-tokyo.ac.jp/sugi/>



東京大学
THE UNIVERSITY OF TOKYO



自己紹介

2

■ 現職:

- 理化学研究所・センター長: 研究者とともに
- 東京大学・教授: 学生とともに
- 企業・技術顧問: 経営者, エンジニアとともに

■ 専門分野:

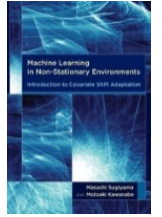
- 機械学習の理論とアルゴリズム開発
- 英語著書→

■ 学会活動:

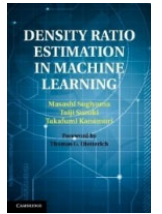
- NeurIPS2015, AISTATS2019, ACML2010/2020 プログラム委員長
- ICLR2023キーノート講演
- 通信学会IBISML研究会委員長(2022-24)



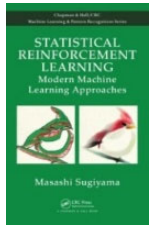
Sugiyama & Kawanabe,
Machine Learning in Non-Stationary Environments,
MIT Press, 2012



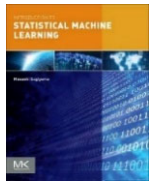
Sugiyama, Suzuki & Kanamori,
Density Ratio Estimation in Machine Learning,
Cambridge University Press, 2012



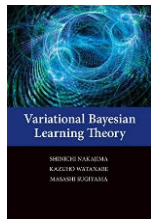
Sugiyama, **Statistical Reinforcement Learning**,
Chapman and Hall/CRC, 2015



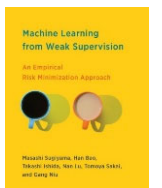
Sugiyama, **Introduction to Statistical Machine Learning**,
Morgan Kaufmann, 2015



Nakajima, Watanabe & Sugiyama, **Variational Bayesian Learning Theory**,
Cambridge University Press, 2019



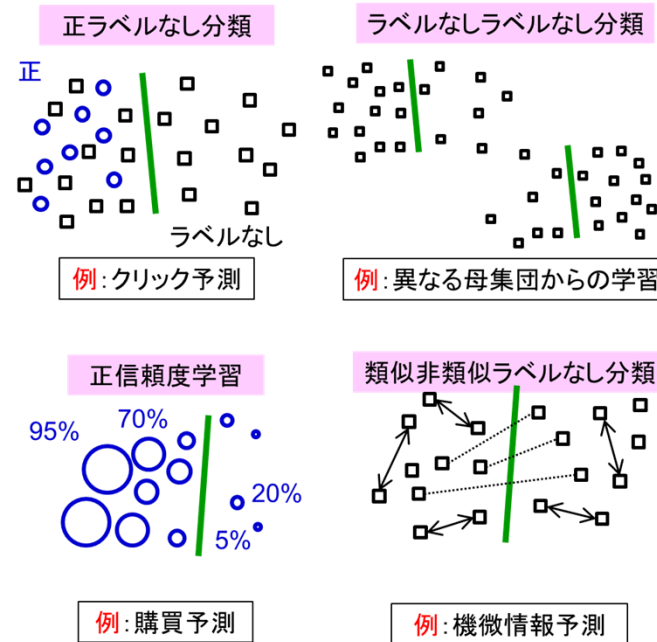
Sugiyama, Bao, Ishida, Lu, Sakai & Niu.
Machine Learning from Weak Supervision,
MIT Press, 2022.



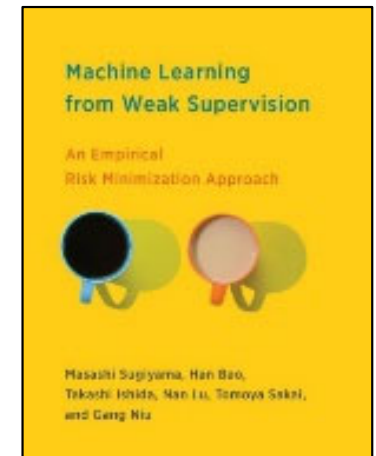
最近の研究：弱教師付き学習

■ 収集コストの低い「弱い」教師データだけから学習

Wang, Xiao, Li, Feng, Niu, Chen & Zhao
(ICLR2022 Outstanding Paper Honorable Mention),
Sugiyama (ICLR2023 Plenary Talk)



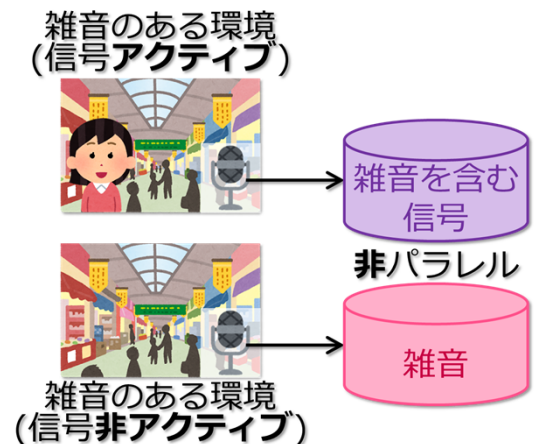
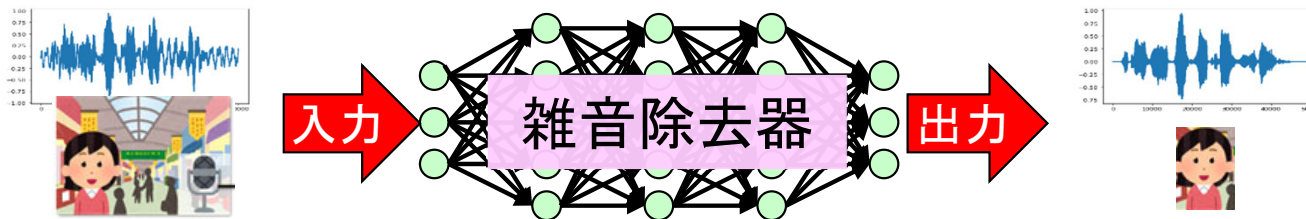
Sugiyama, Bao, Ishida, Lu, Sakai & Niu.
(MIT Press 2022)



■ 例：音声の雑音除去

Ito & Sugiyama
(ICASSP2023 Best Paper Award)

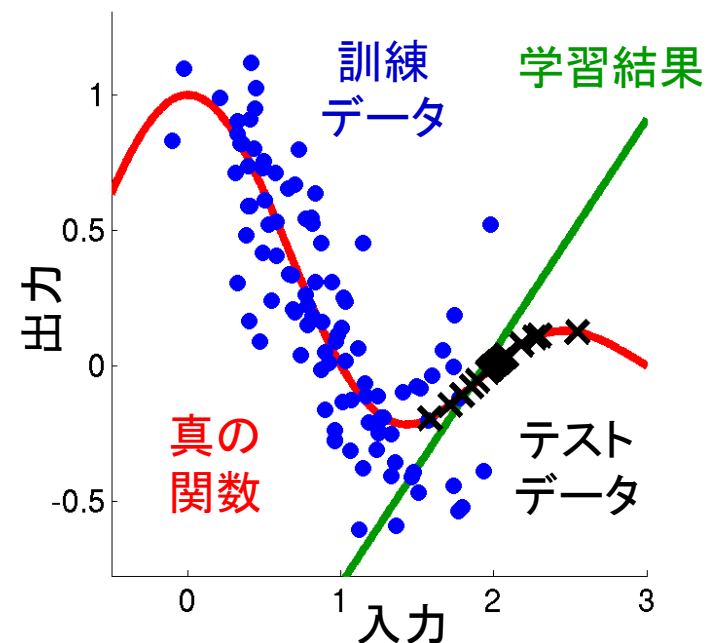
- 雑音を含む音声と雑音だけから学習できる！（収集コストのかかる雑音を含まない音声は不要）



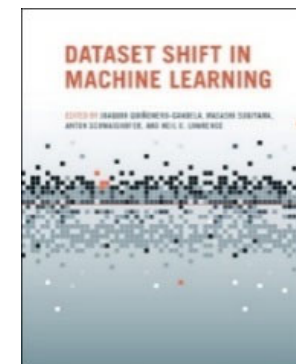
最近の研究：転移学習

■ 訓練データとテストデータの分布の違いを吸収

- 古典的な設定：共変量シフト
(入力分布だけが変化) Shimodaira (JSPI2000)



Quiñonero-Candela,
Sugiyama, Schwaighofer &
Lawrence (MIT Press 2009)

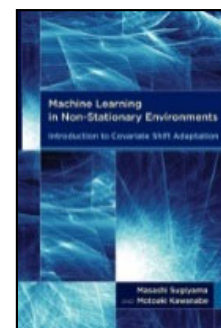


- 重要度重み付き学習:

$$\operatorname{argmin}_{f \in \mathcal{F}} \left[\sum_{i=1}^{n_{\text{tr}}} \frac{p_{\text{te}}(\mathbf{x}_i^{\text{tr}})}{p_{\text{tr}}(\mathbf{x}_i^{\text{tr}})} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right]$$

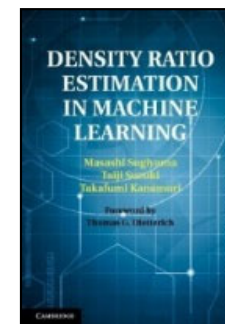
重要度 損失

重要度(密度比)の推定が重要!



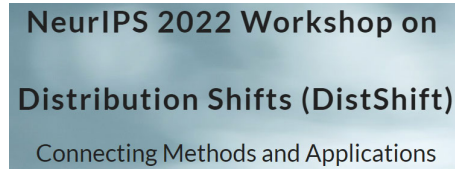
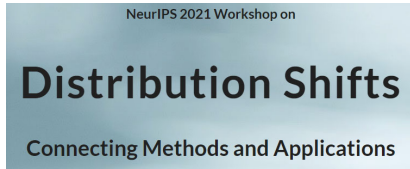
Sugiyama &
Kawanabe
(MIT Press 2012)

Sugiyama, Suzuki &
Kanamori (Cambridge
University Press 2012)



最近の研究: 転移学習

■ 近年再注目: NeurIPS2021, 2022, 2023ワークショップ他



■ 最近の発展:

● 重みと予測器の同時学習

Zhang, Yaname, Lu & Sugiyama
(ACML2020 Best Paper Award, SNCS2021)

● 連続分布シフトへの拡張

Bai, Zhang, Zhao, Sugiyama & Zhou (NeurIPS2022)
Zhang, Zhang, Zhao & Sugiyama (NeurIPS2023)

● 同時分布シフトへの拡張

Fang, Lu, Niu & Sugiyama
(NeurIPS2020 Spotlight)

● 分布外適応への拡張

Fang, Lu, Niu & Sugiyama
(NeurIPS2023 Spotlight)

■ 予測誤差の上界の同時最小化:

$$\min_{r, f} J_{\ell'}(r, f) \quad \begin{aligned} J_{\ell'}(r, f) &\geq \frac{1}{2} R_{\ell}(f)^2 \\ R_{\ell}(f) &= \mathbb{E}_{p_{\text{tr}}(\mathbf{x}, y)}[\ell(f(\mathbf{x}), y)] \\ \ell' &\leq 1, \ell' \geq \ell, r \geq 0 \end{aligned}$$

$J_{\ell'}(r, f) = \mathbb{E}_{p_{\text{tr}}(\mathbf{x})}[(r(\mathbf{x}) - r^*(\mathbf{x}))^2] \leftarrow$ 最小二乗重要度推定
 $+ (\mathbb{E}_{p_{\text{tr}}(\mathbf{x}, y)}[r(\mathbf{x})\ell'(f(\mathbf{x}), y)])^2 \leftarrow$ 重要度重み付き学習

- 従来法は上界の二段階最小化に相当
- 収束性を理論保証:
 $R_{\ell}(\hat{f}) \leq \sqrt{2} \min_{f \in \mathcal{F}} R_{\ell}(f)$

訓練フェーズ テストフェーズ

$t=1$ $t=2$... $t=T$

■ 連続クラス事前分布シフト:

- クラスの比率 $p_t(y)$ だけが変化

■ 連続共変量シフト:

■ 与えられるデータ: 訓練とテストの入出力標本

$$\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y) \quad \{(\mathbf{x}_j^{\text{te}}, y_j^{\text{te}})\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}, y)$$

■ 各ミニバッチ $\{(\bar{\mathbf{x}}_i^{\text{tr}}, \bar{y}_i^{\text{tr}})\}_{i=1}^{\bar{n}_{\text{tr}}}, \{(\bar{\mathbf{x}}_j^{\text{te}}, \bar{y}_j^{\text{te}})\}_{j=1}^{\bar{n}_{\text{te}}}$ に対して、重要度をカーネル平均適合で推定: Huang, et al. (NeurIPS2007)

■ ドメイン外への拡張:

- 訓練ドメインの外では重要度が発散
- 外れ値検知を用いて、テストデータを訓練ドメイン内外に分割:
 $\{(\mathbf{x}_j^{\text{te, in}}, y_j^{\text{te, in}})\}_{j=1}^{n_{\text{te, in}}}, \{(\mathbf{x}_j^{\text{te, out}}, y_j^{\text{te, out}})\}_{j=1}^{n_{\text{te, out}}}$
- 損失を個別に計算:
$$\frac{n_{\text{te, in}}}{n_{\text{tr}} n_{\text{te}}} \sum_{i=1}^{n_{\text{tr}}} \frac{p_{\text{te}}(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})}{p_{\text{tr}}(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) + \frac{1}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te, out}}} \ell(f(\mathbf{x}_j^{\text{te, out}}), y_j^{\text{te, out}})$$

今後の課題：無仮定機械学習

6

■ これまで：

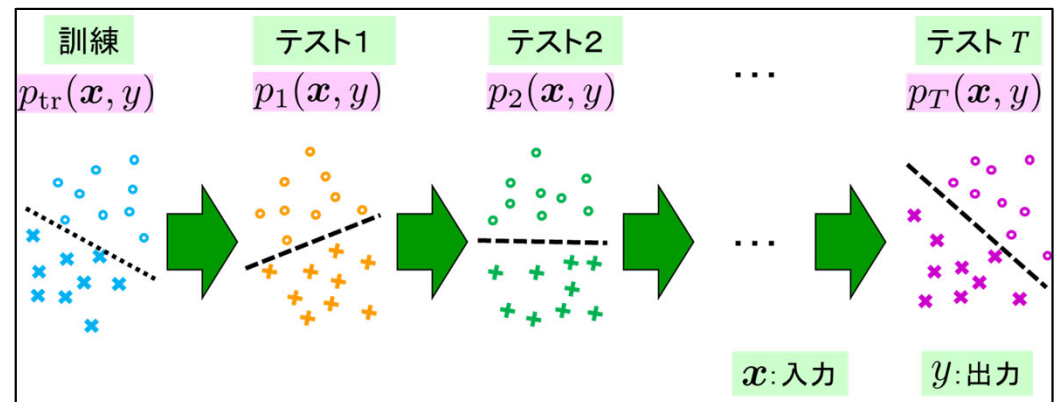
- (共変量シフトなど) **適切な仮定**のもとでの理論とアルゴリズム
- 仮定が正しければうまくいく. そうでなければ保証なし

■ 現実：

- 仮定が成り立っているか分からない
- やってみてうまくいったらラッキー

■ 目標：

- (ほぼ) **無仮定機械学習**の理論とアルゴリズム開発
- (ほぼ) 任意の場面で最低限の性能が保証される
- 現場で最初に試すべき
極限的な汎用手法！
- 例：**弱教師付き**
連続同時分布シフト適応



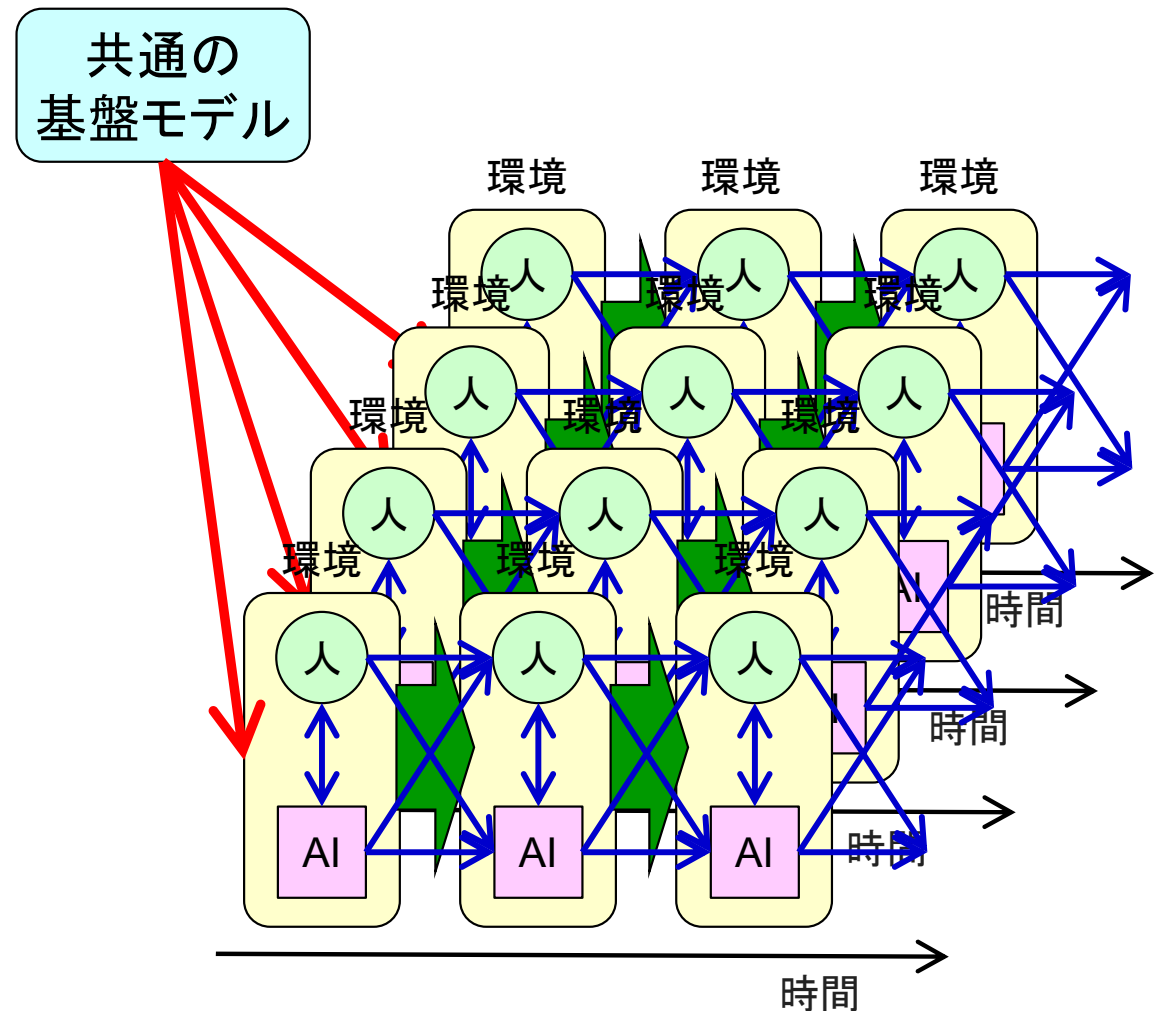
今後の課題：個別化AIの全体最適化⁷

- 中央にある巨大な共通の基盤モデルを、個人ごと、タスクごとに**個別化**：

- 各AIがユーザと相互作用しながら**時間発展**
- AI同士，ユーザ同士も**相互作用**

- 社会全体で最適化：

- 効率性
- 信頼性
- 説明性
- 公平性
- 安全性



今後の課題：個別化AIの全体最適化⁸

■ 学習の信頼性向上：

- 不完全情報からの弱教師付き学習
- 低信頼性データからのロバスト学習
- 予測の不確実性に基づく棄却
- データの異常・変化検知
- データ改変に対する敵対的学習
- 国、文化、人種等の多様性のドメイン適応
- 非定常環境、タスク変化に対する継続学習
- 言語、画像、音響、科学データ等のマルチモーダル学習

■ 理論解析：

- 深層学習の汎化解析
- 非凸・確率的・分散最適化の収束解析
- 学習理論と最適化理論の融合

■ 大規模モデル・データ対応：

- 個別化モデルのマルチタスク学習
- モデルを逐次更新するオンライン学習
- 不確実性を考慮した近似ベイズ推論
- モデルを小型化するモデル圧縮
- 圧縮データからの直接学習
- 分散データからの連合学習
- 巨大モデルの分散最適化

■ 新しい推論機構：

- データからの因果推論
- 基盤モデルを用いた文脈内学習
- 長期的な予測に基づく逐次的意思決定
- 動的システムに基づく時系列予測
- データ駆動型仮説検定

■ 社会的観点：

- 予測・モデルの説明性
- プライバシ保護
- 意思決定・資源配分の公平性
- 多様なAIと人間の相互作用
- モデル・ユーザ間通信の安全性

画像生成やテキスト対話型の人工知能（AI）システムが注目を集めている。専門的な知識は必要とせず、簡単な言葉を入力するだけで美的な出力が得られるため、AIの活用は新たな局面に突入したといえる。本稿では、AIシステムの基本となる機械学習の技術的な概念と、その本質について述べていく。機械学習とは、機械（コンピュータ）に学習能力を持たせるための情報技術の総称である。音声や画像の認識、言語の翻訳、顧客情報の分析、株価の予測、サイバー攻撃の検知、自動運転車の制御など、様々な場面で用いられている。



すぎやま・まさし
74年生まれ。東京工業大学博士（工業）。専門は機械学習。理化学研究所センター長

AI開発の現在地①

杉山将 東京大学教授

原理解明さらなる研究必要

日経新聞 経済教室(2023年3月30日)

AI研究人材

- 日本の国際的な研究競争力は限定的：
 - 日本人の博士進学者数は増加せず
 - 博士支援を改善しても、タイパの悪い研究はしたくない
 - 博士学生不足が研究者・教員不足に直結
 - 少子高齢化の加速で、今後も抜本的な改善は見込めない
 - 頼みの留学生も、少子化・経済安保の影響を受ける見込み
 - 社会人の再教育は、必ずしも高度研究人材育成につながらない
 - 女子学生の情報系進学への支援は、効果は低いが必要
 - (親御さん・小中高の進路担当教員・予備校への働きかけ)
- 世界の優秀な若手人材はスター研究者の周りに結集：
 - 安定してボトムアップに研究ができる環境を整え、日本発の国際的スター研究者を育成すべき
 - 海外人材のクローポ雇用できる仕組みの導入