

統計的機械学習 (応用計量分析2) 第5回

条件付き平均因果効果の推定法 (参考pdf 6章)

メタ学習器

二重機械学習

振り返り

ATE/ATT推定法を学んだ

- マッチング法
 - 距離の定義によってマッチの取り方がいくつかある
 - 高次元では距離は難しいので1次元の傾向スコアを使う方法も
 - 傾向スコアはバランシングスコアなのでそれが類似ならペアとしてよい
 - 傾向スコアの推定にはプロパーな損失を用いる
- IPW法（ホーヴィッツ＝トンプソン推定量）
 - 行動と結果の積は観測可能、行動をかけて歪んだ分布を傾向スコアで重み付けて戻す
- G計算（CATE推定量を用いる）
- 二重に頑健な推定量（AIPW）
 - $\hat{\mu}$ と \hat{f} のどちらかが正しければ正しい推定

本日の内容

CATE推定法を進める（シラバスより巻きを進める想定）

- 1. ガイダンス・因果推論と機械学習の概論
- 2. 意思決定理論（期待効用理論）の復習、因果推論との関係
- 3. 潜在結果モデルに基づく因果推論の枠組み
- 4. 平均因果効果の推定法
- **5. 条件付き平均因果効果（CATE）の推定法1：メタ学習器**
- **6. CATEの推定法2：二重機械学習**
- 7. CATEの推定法3：決定木と決定森
- 8. 構造方程式モデルとバックドア基準
- 9. 因果探索
- 10. 発展的な因果推論手法：フロントドア調整、操作変数法、回帰不連続デザイン、代理変数法
- 11. 発展的な意思決定理論
- 12. 強化学習
- 13. オフライン強化学習
- 14. バンディット
- 15. まとめ

条件付き平均因果効果と評価指標

因果推論の典型的な精度は潜在結果や効果の精度

- 推定対象

- 条件付き平均因果効果 (Conditional Average Treatment Effect; **CATE**)

$$\tau(x) = \mathbb{E}[y_1 - y_0 | x]$$

- ある程度個別化された効果。因果機械学習 (機械学習分野における因果推論研究) でよく推定対象とされる

- 条件付き平均潜在結果 $\mathbb{E}[y_a | x]$

- 評価指標

- 因果効果の精度 $\text{PEHE}(\hat{\tau}) = \mathbb{E}_x \left[(\tau(x) - \hat{\tau}(x))^2 \right]$

- 潜在結果の精度 $\text{MSE}^u(\hat{f}) := \mathbb{E}_x \left[\frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} \left(\mathbb{E}[y_a | x] - \hat{f}(x, a) \right)^2 \right]$

- 意思決定価値 $V(\pi) := \mathbb{E}_x \mathbb{E}_{a \sim \pi(a|x)} \mathbb{E}[y_a]$

CATE推定のためのメタ学習器 (1/4)

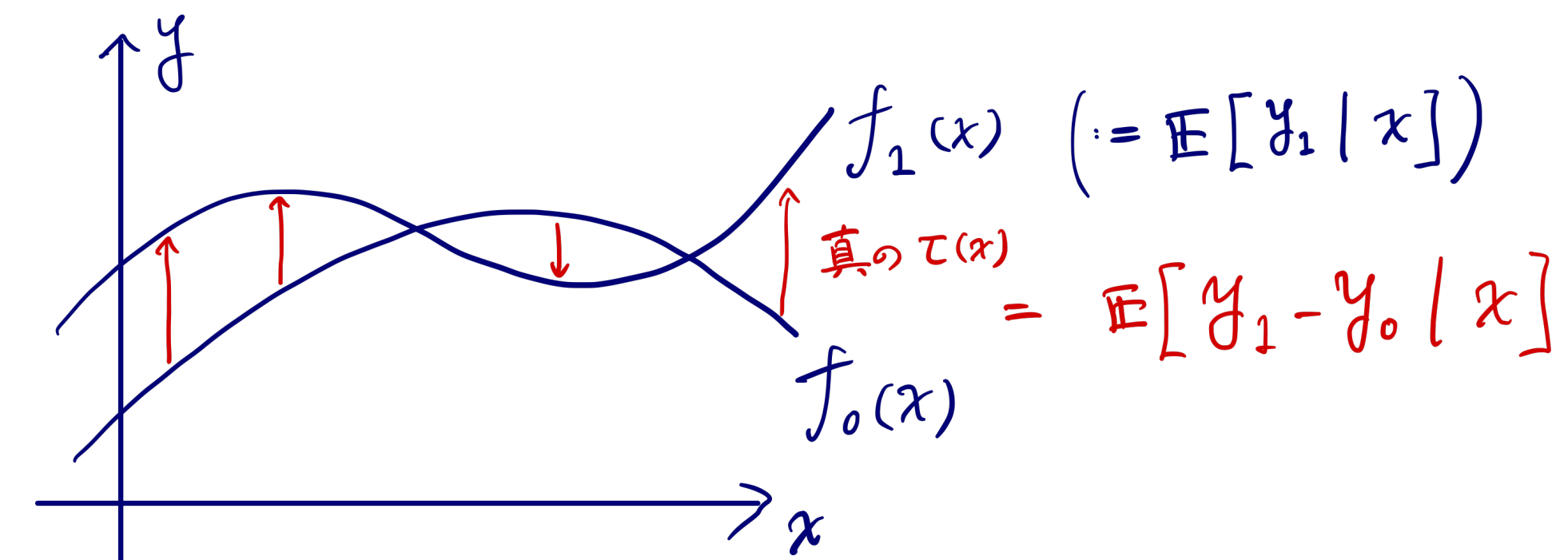
メタ学習器は既存の機械学習器をブラックボックス的に使う推定量 a ごとにモデルを分けるT-Learner、 a も入力とするS-Learner

- メタ学習器 (Meta learners)
 - 教師あり学習のためのベース学習器 (Base learners) を用いて因果推論用の学習器を構成する手法
 - 参考) 機械学習のメタ学習: 「学習する方法を学習する」ことで新たなデータでの学習を効率化する手法群 とは別物
- S(ingle)-Learner
 1. 行動 a と説明変数 x を区別せず一緒に入力変数として用いて学習: $f(x, a)$
 2. $a = 1, 0$ をそれぞれ入力した予測値の差分を効果の推定とする
$$\hat{\tau}(x) = \hat{f}(x, a = 1) - \hat{f}(x, a = 0)$$
- T(wo)-Learner
 1. 行動ごとにデータを分け、それぞれ別々の潜在アウトカム予測モデルを学習
 2. 学習した \hat{f}_1, \hat{f}_0 にそれぞれ x を入力した予測値の差を効果の推定量とする:
$$\hat{\tau}(x) = \hat{f}_1(x) - \hat{f}_0(x)$$
 - もっとも基本的な選択肢

T-Learnerの課題 (X-Learnerへの伏線)

結果モデル f_a は複雑、曝露群が少ないとき推定精度悪化

- T-Learnerは曝露群・統制群両方の結果モデル推定精度に依存
 - $\hat{\tau}(x) = \hat{f}_1(x) - \hat{f}_0(x)$
- 一方で、真の因果効果 $\tau(x)$ は単純かもしれない
 - x に依存しない一定値 $\tau(x) = \tau$ かもしれない
 - x の一部の変数にしか依存しないかもしれない
 - にもかかわらず、結果モデル全体 f_a の推定精度に依存してしまう
 - とくに、暴露群は小さい ($\mu(a=1|x)$ が小さい) ことが多く、 \hat{f}_1 は信頼できない
- 因果効果自体のデータはないので $\tau(x)$ を直接学習することはできない
 - しかし、 \hat{f}_0 は比較的信頼できる。ならば、 $\hat{\tau}^i = y_1 - \hat{f}_0(x)$ と推定し、これを学習すればよいのではないか？
 - → これを、 x ごとにサンプルの多い方を”信頼”して最終的な推定を行うように一般化したのが **X-Learner**



CATE推定のためのメタ学習器 (2/4)

X-Learner : 群ごとの効果モデルの重み付き平均

- X-Learnerの手順

1. T-learner と同様に、行動ごとの潜在結果予測モデルを学習

$$\hat{f}_1(x), \hat{f}_0(x)$$

2. 行動ごとのデータに対し上記モデルをプラグインして個別の効果推定値を計算

$$\hat{\tau}_0^i = \hat{f}_1(x^i) - y_0^i \quad (\text{for } i : a^i = 0)$$

$$\hat{\tau}_1^i = y_1^i - \hat{f}_0(x^i) \quad (\text{for } i : a^i = 1)$$

- 欠損している方を潜在結果モデルで補完

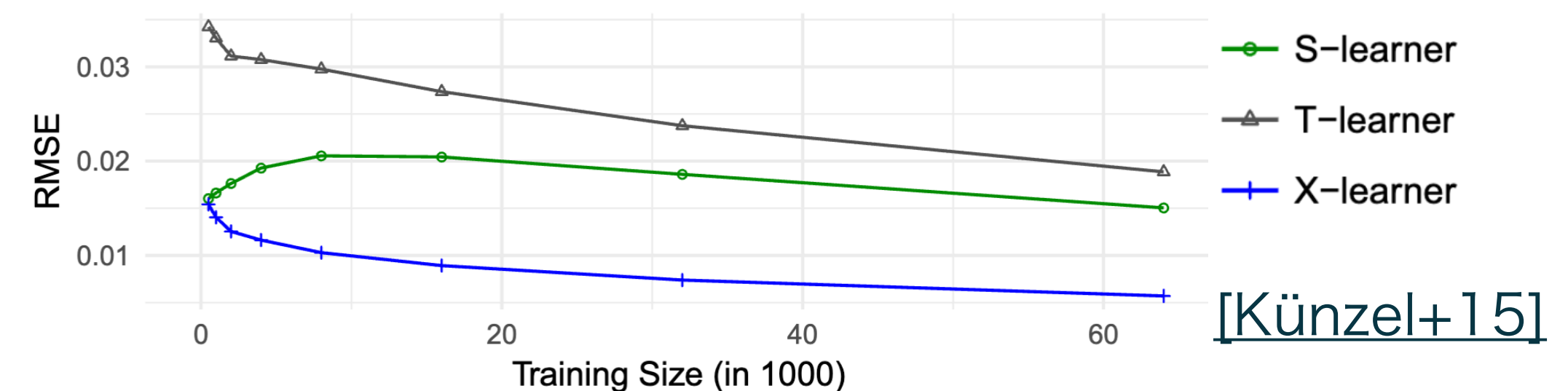
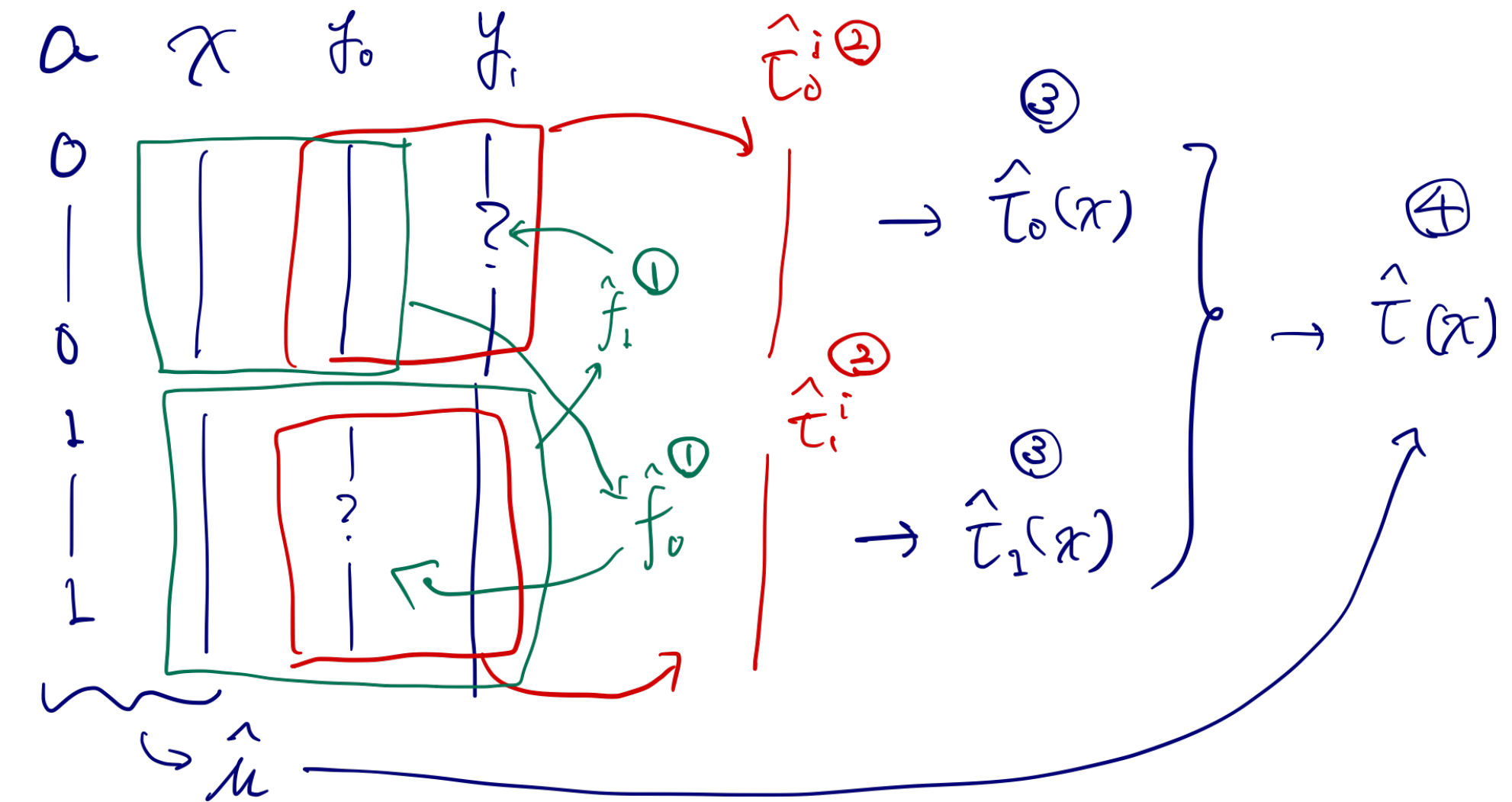
3. 行動ごとの効果モデルを学習

$$\hat{\tau}_0 = \mathcal{A}(\{\hat{\tau}_0^i, x^i\}_i), \quad \hat{\tau}_1 = \mathcal{A}(\{\hat{\tau}_1^i, x^i\}_i)$$

4. 最終的に傾向スコアで重み付きアンサンブル (足し合わせ)

$$\hat{\tau}(x) = \hat{\mu}(x)\hat{\tau}_0(x) + (1 - \hat{\mu}(x))\hat{\tau}_1(x)$$

- $\mu(x)$ が小さい = 曝露群が小さい x では \hat{f}_0 を途中で用いた $\hat{\tau}_1(x)$ を信頼



CATE推定のためのメタ学習器 (3/4)

目的変数変換法：ホーヴィッツ・トンプソン推定量の応用 傾向スコアによるノイズ付き個別因果効果を用いた直接回帰

- 目的変数変換法：ATE推定のためのホーヴィッツ＝トンプソン推定の応用

1. 傾向スコア推定モデル $\hat{\mu}(x)$ を学習
2. 各サンプル $(x^i, a^i, y_{a^i}^i)$ に対し、目的変数を以下により変換

$$z^i = y_1^i \frac{a^i}{\hat{\mu}(x^i)} - y_0^i \frac{(1-a^i)}{1-\hat{\mu}(x^i)}$$

- $\hat{\mu}$ が正しければ、 z^i の条件付き期待値はCATEになる $\mathbb{E}[z^i | x^i] = \tau(x^i)$

3. z を x に回帰

$$\hat{\tau} = \mathcal{A}(\{x^i, z^i\}_i)$$

- 利点：結果モデル $f(x, a)$ の推定を経ることなく直接的に $\hat{\tau}$ のモデルを学習できる
- 欠点：傾向スコアの推定 $\hat{\mu}$ の値が0や1に近い場合、 z^i に加わるノイズが大きく推定分散が高い

CATE推定のためのメタ学習器 (4/4)

DR-Learner：二重頑健推定量の応用

結果予測モデルか傾向スコアモデルのどちらかが良ければ良い

- DR-learner

1. T-learnerと同様、行動ごとの潜在結果モデルを学習

$$\hat{f}_1(x), \hat{f}_0(x)$$

2. DR推定量の結果変換を適用

$$z^i = \hat{f}_1(x^i) - \hat{f}_0(x^i) + \frac{y_1^i - \hat{f}_1(x^i)}{\hat{\mu}(x^i)} a^i - \frac{y_0^i - \hat{f}_0(x^i)}{1 - \hat{\mu}(x^i)} (1 - a^i)$$

- 直接法による予測 + 残差を目的変数変換法と同様逆傾向スコア重み付け

3. その差を x に回帰して効果モデルを学習

$$\hat{\tau} = \mathcal{A}(\{x^i, z^i\}_i)$$

- 利点

- 結果モデルと傾向スコアモデルの**いずれかが**正しければ後段の $\hat{\tau}$ の推定への悪影響が少ない
- 理論的には、全体の収束レートが2つのモデルの収束レートの**積**になる
 - 片方の仮説集合が間違っていて $O(1)$ レートでも、もう片方が正しくて $O(N^{-1/2})$ レートであればよい
 - 両方のモデルがそれぞれ $O(N^{-1/4})$ 以上であれば効率的なレート $O(N^{-1/2})$ が達成される

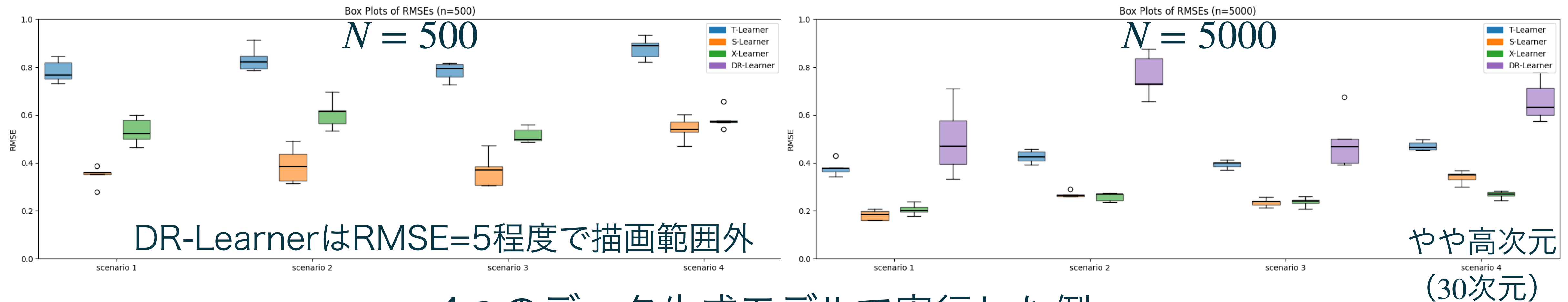
- 欠点

- $\hat{\mu}$ の値が0や1に近い場合、残差なのでマシとはいえやはりノイズが大きく推定が不安定

性能比較例

DR-Learnerは逆数をとるので不安定になりがち

- 実際の精度は諸条件に依存するので参考まで
 - サンプルサイズ N 、次元 p 、 μ の偏り、真のモデル $f(x, a)$ の複雑さ、ベース学習器、等
- 理論では $N \rightarrow \infty$ の漸近的な性質を議論しがち（理論的に言えることが他にあまりない）だが、DR-Learnerなど推定値の逆数をとる推定法はサンプルサイズが大きいとき不安定がち



4つのデータ生成モデルで実行した例

(ベース学習器：勾配ブースティング決定木 GBDT)

二重機械学習

二値行動の一般化：行動線形モデル

二値/連続値行動向けセミパラメトリック法：DML (1/3)

二重機械学習：連続的行動 $a \in \mathbb{R}$ 等にも適用可能

● 二重機械学習 (Double/Debiased Machine Learning; DML, R-Learner)

- $y = \theta(x) \cdot a + g(x, w) + \varepsilon \quad (\mathbb{E}[\varepsilon | x, w] = 0) \cdots (1)$

- 説明変数のうち、 a の係数に関与する部分 (効果修飾因子) を x 、関与しない部分を w
- 行動 a に関して線形性を仮定、その他の部分 g は非線形を許容 (セミパラメトリック)

- $a = \mu(x, w) + \eta \quad (\mathbb{E}[\eta | x, w] = 0)$

- 無視可能性から $\mathbb{E}[\eta \cdot \varepsilon | X, W] = 0$

● 残差を線形モデルで学習

- (1) の期待値をとると $\mathbb{E}[y | x, w] = \theta(x) \cdot \mathbb{E}[a | x, w] + g(x, w) \cdots (2)$

- (1) の両辺から (2) の両辺を引くと、

- $Y - \mathbb{E}[Y | X, W] = \theta(X) \cdot (a - \mathbb{E}[a | X, W]) + \varepsilon$
“普段”の売上との差 “普段”の価格との差

● 結果を (x, w) に回帰した残差を行動の残差に回帰する

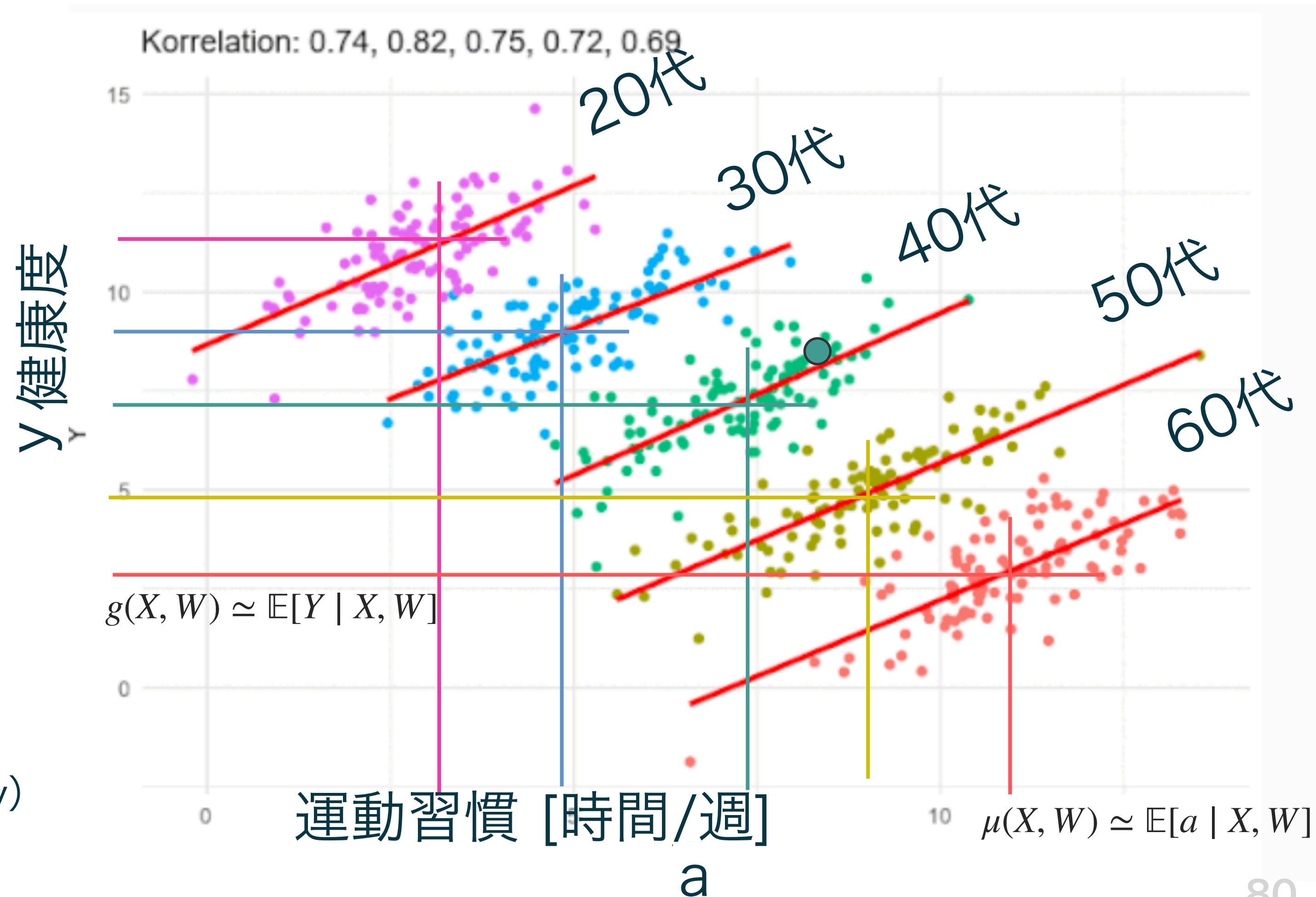
← 推定の邪魔だった $g(x, w)$ が消えている!

Robinsonの分解とも呼ばれることから
DMLによるCATE推定法のことをR-learnerとも呼ぶ

DML (2/3) イメージ

行動以外の背景因子の効果を差っ引くことで 背景因子ごとに原点をずらす

- $y - \mathbb{E}[y | x, w] = \theta(x) \cdot (a - \mathbb{E}[a | x, w]) + \varepsilon$
- 属性 (x, w) ごとに「原点」を設定
そこからの差分にフォーカス
- g と μ の推定誤差が θ の推定誤差
に影響する度合いが小さい
 - セミパラメトリック推論の理論
 - 誤差 Δ_θ を (Δ_μ, Δ_g) についてテイラー
展開した1次の係数がゼロ
 - ネイマン直交性 (Neyman orthogonality)
という



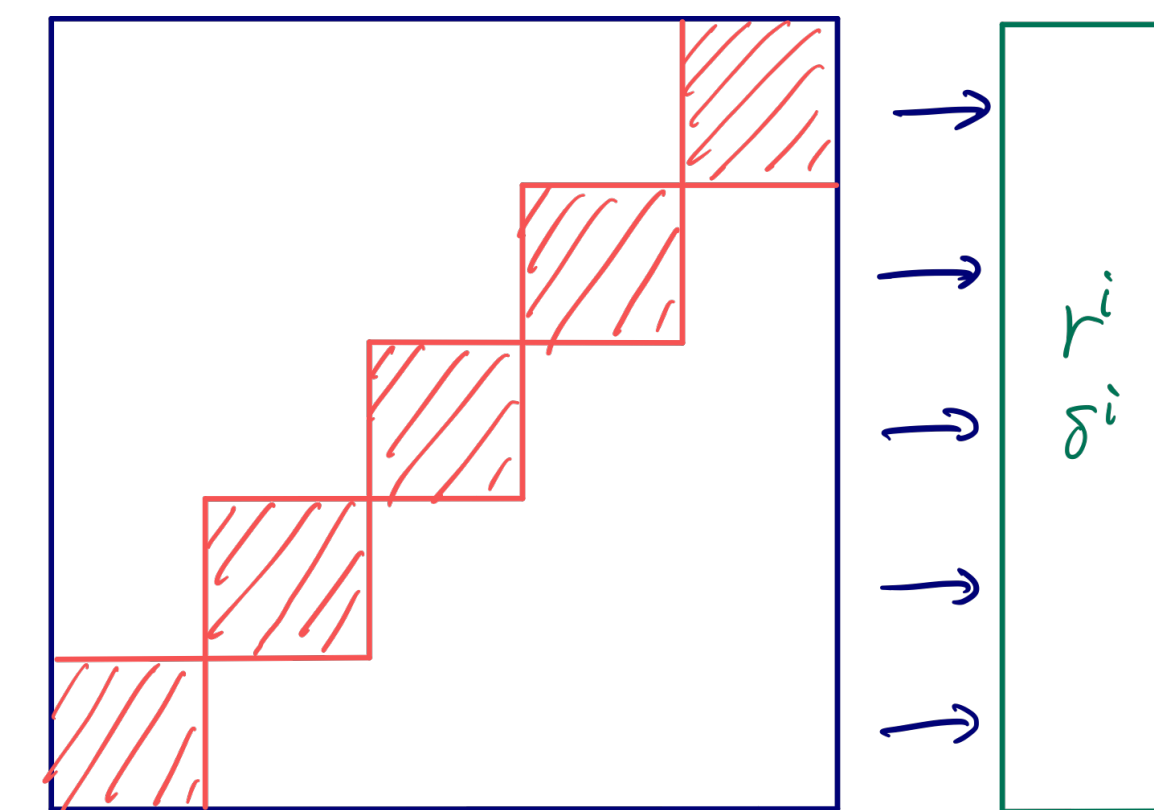
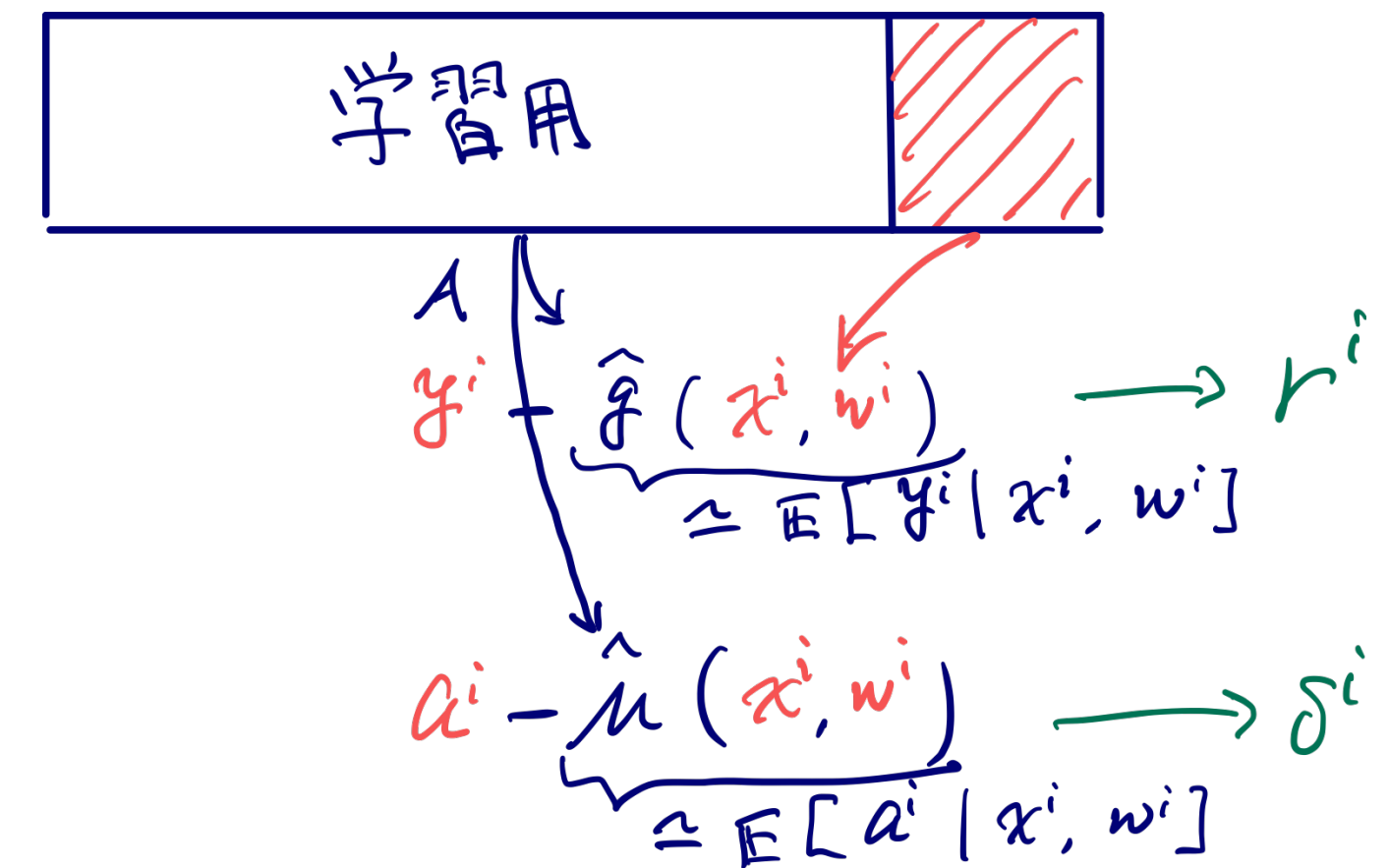
DML (3/3) データ分割法 (クロスフィット)

残差を偏りなく推定するため、K-foldに分割して推定

- 学習用データと同じデータで残差を評価すると過小評価になる
 - 訓練誤差はテスト誤差より通常小さい
- データをK分割し、未学習データで残差 r^i, δ^i を算出
- これらを集め、 (x^i, r^i, δ^i) から $\theta(x)$ を推定

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{N} \sum_i^N (r^i - \theta(x^i) \delta^i)^2$$

- (または各foldで学習したモデルを平均する)



$$\rightarrow \hat{\theta} = \arg \min_{\theta} \frac{1}{N} \sum_i^N (r^i - \theta(x^i) \delta^i)^2$$

(参考) Neyman直行性

モーメント条件が迷惑パラメタの推定誤差に関して直行

- 真のパラメタ θ_0 , 真の迷惑モデル $\eta_0 = (g, \mu)$ に対して、
損失関数がモーメント条件 $E_P [\psi(W; \theta_0, \eta_0)] = 0$ を満たすと仮定する
 - ψ はスコア関数。損失関数の θ 微分だと思えばOK
- Neyman直行性の定義
 - $\partial_\eta E_P [\psi(W; \theta_0, \eta_0)] [\eta - \eta_0] = 0$, for all $\eta \in \mathcal{T}_N$
 - 微少なズレ $\Delta_\eta \ll 1$ では期待モーメントは0から変化しない
→真のパラメタへの影響も小さい
- なお、Doubly RobustならNeyman直行だが、その逆は必ずしも真ではない
 - R-Learnerは二重頑健ではない

(参考) DMLの構造的行動への拡張

構造のある行動空間の特徴付けに対して線形なモデル

- Causal Effect Inference for Structured Treatments (NeurIPS '21)

- 積効果モデル: $y = g(\mathbf{x})^\top h(\mathbf{a}) + \varepsilon$

- cf) DML: $y = \theta(x) \cdot a + g(x, w) + \varepsilon \rightarrow a$ は二値か1次元か1-hotベクトルくらい

- 行動を特徴付ける空間 h を同時学習する

- $e^h(x) := \mathbb{E}[h(a)|x]$

- 一般化Robinson分解

- $y - m(x) = g(\mathbf{x})^\top (h(\mathbf{a}) - e^h(x)) + \varepsilon$

- 学習ステップ

- $m(x) \approx \mathbb{E}[y|x]$ を学習

- while not converged

- g, h をK回更新

- e^h を更新

Causal graph

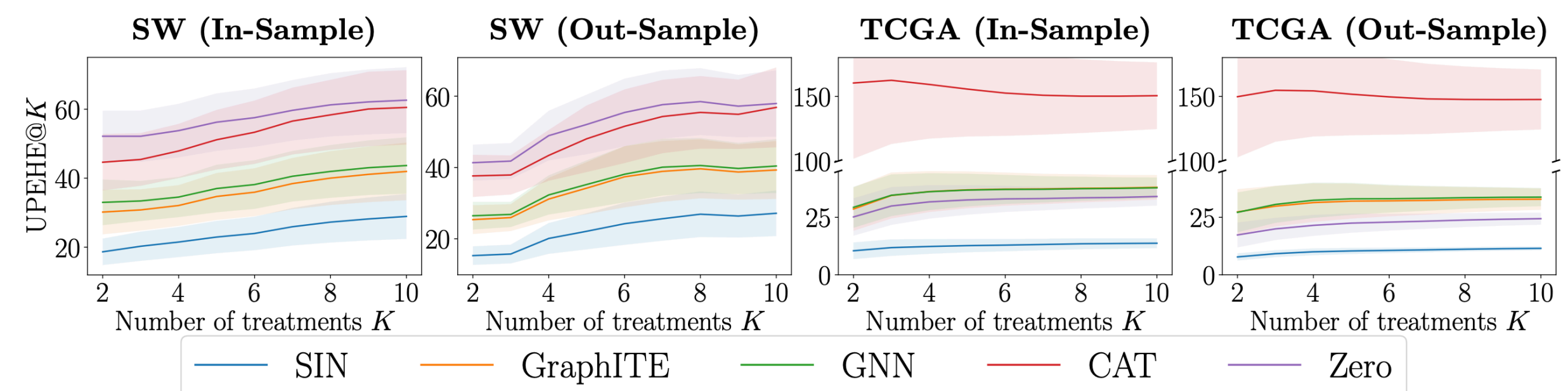
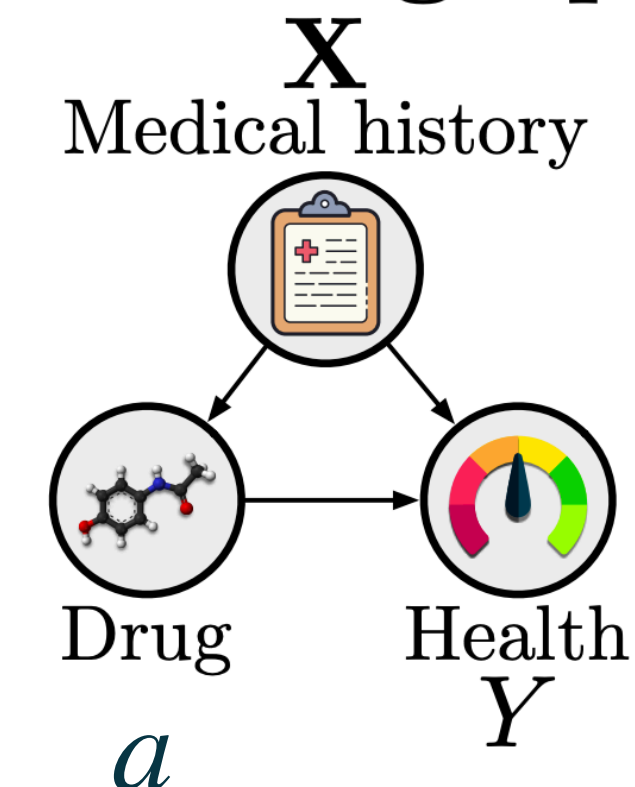


Figure 3: UPEHE@K for $K \in \{2, \dots, 10\}$.

まとめ

CATE推定のためのメタ学習とDML

- メタ学習器
 - S-Learnerは行動を入力変数に含めてモデル化
 - T-Learnerは行動ごとに分けてモデル化
 - X-LearnerはT-Learnerをベースに曝露/統制群間の偏りによって結果モデルの推定精度がネックにならないよう群ごとに $\hat{\tau}$ を学習して重みつき平均
 - DR-LearnerはAIPW法を応用した損失関数で二重の頑健性を保証
- 二重機械学習 (DML)
 - 行動が二値 $a \in \{0,1\}$ の場合に加えて連続値の場合に対して線形モデルとしたセミパラメトリックモデル (a 以外に関しては非線形を許容)
 - 結果 y を説明変数 (a 以外) に回帰したモデルの残差を、行動を説明変数に回帰したモデルの残差に回帰
 - 残差が過小推定とならないようにデータ分割、クロスフィットを行う