

統計的機械学習（応用計量分析2）第6回

条件付き平均因果効果の推定法（参考pdf 6章）

振り返り

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

本日の内容

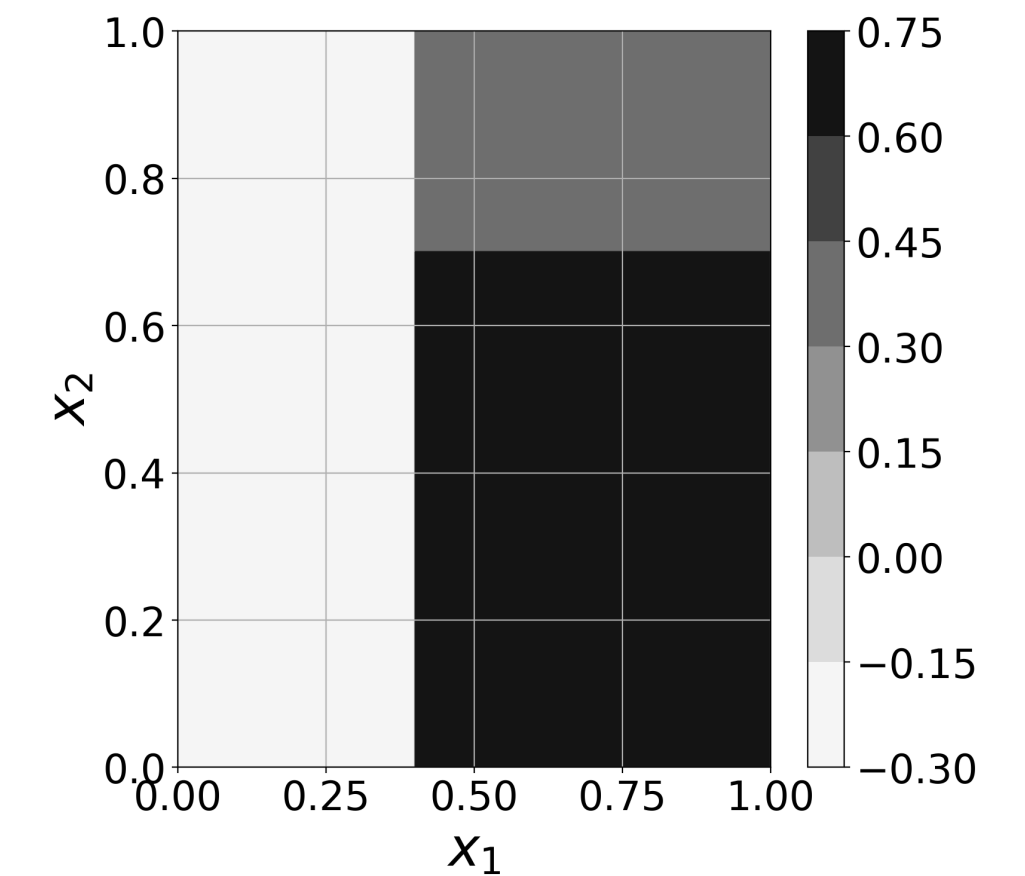
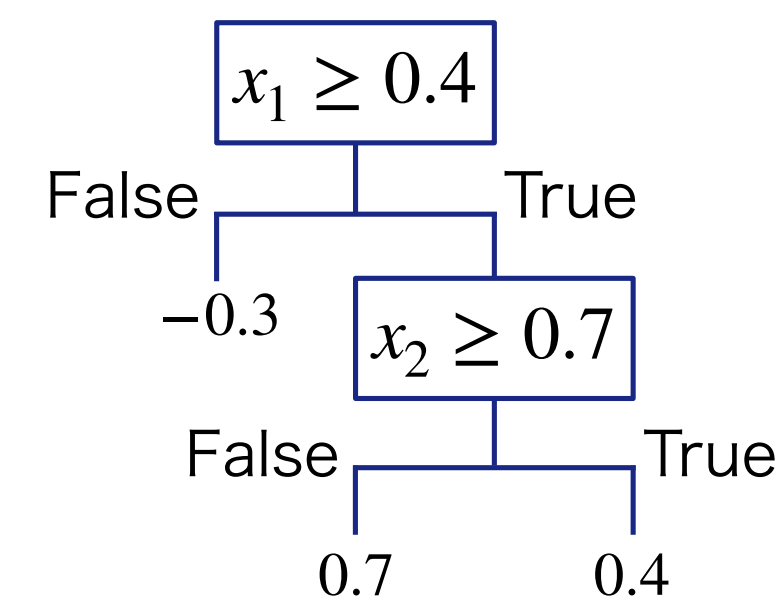
CATE推定法の中でも決定木系の手法について学ぶ ベイス推論についてもとりあげる

- 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
 - CATEの推定法2：二重機械学習
- **6. CATEの推定法3：決定木と決定森**
 - 深層学習に基づく方法
- 7. 構造方程式モデルとバックドア基準
- 8. 因果探索
- 9. 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
- 10. 続き
- 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

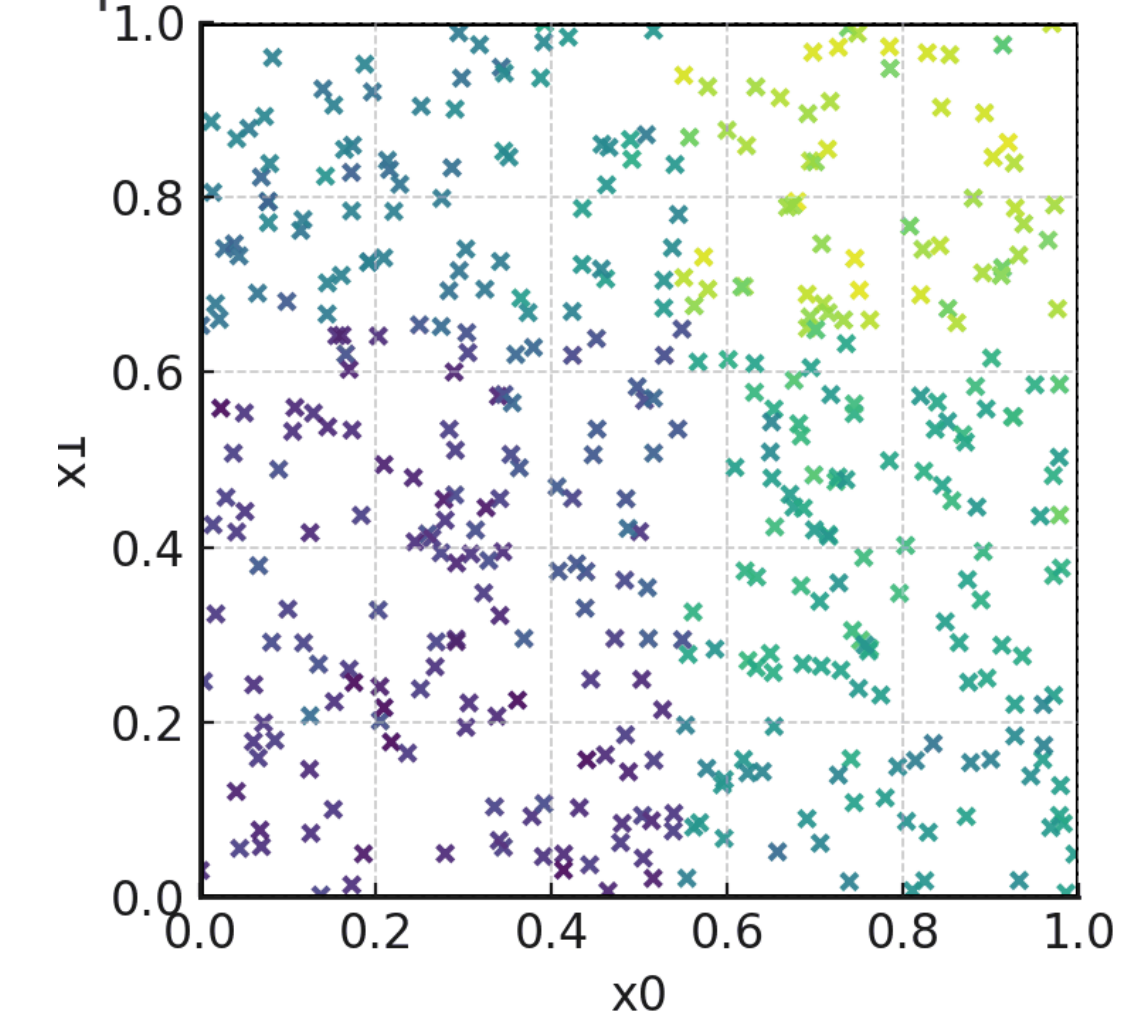
決定木の基本

決定木は各次元で領域分割した上で定数で予測するモデル

- 決定木は説明変数の各次元の値の大小によって場合分けした領域ごとに定数を予測値として出力
 - 分割ごとの変数としきい値、領域ごとの予測値が学習可能なパラメタ
- 決定木の基本的な学習法
 1. 現状の領域分割に追加して最も目的変数を改善する分割を逐次的に探す
 - 領域の予測値は平均値とする
 2. 停止条件に合致したら停止
 - 木の深さが一定になるか、領域内のサンプルサイズが一定以下になったらその領域は分割終了
- 各分割は説明変数軸に沿って行われる
 - 「斜め切り」はしない。ただ、表データは各軸が意味的に独立していることが多く、現実には有益な制約であることが多いとされる



epth 0 — Start node with 420 sampl



因果推論向けの決定木

因果木：データ分割の工夫と因果効果の直接推定

- 因果効果の直接推定を行う

- 決定木＝領域内は同一の値で近似

- x の属する領域を $\ell(x)$ として、（※葉ノード＝リーフのL。損失関数の ℓ とは全く別）

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

- 領域内の曝露群と統制群でそれぞれ期待値を平均で置き換える

- 領域分割の最適化に用いる訓練データ D^{tr} 、予測値の計算に用いる推定用データを分ける：**Honestアプローチ**

- 説明変数の空間 \mathcal{X} の領域 $\ell_k \subset \mathcal{X}$ によるある分割 $\Pi = \{\ell_k\}_k$ ($\cup_k \ell_k = \mathcal{X}$) に対して推定用データ \mathcal{D}^{est} を用いて

- $\hat{\tau}(x; \mathcal{D}^{\text{est}}, \Pi) = \frac{1}{|\mathcal{D}_1^{\ell(x)}|} \sum_{i \in \mathcal{D}_1^{\ell(x)}} y^i - \frac{1}{|\mathcal{D}_0^{\ell(x)}|} \sum_{i \in \mathcal{D}_0^{\ell(x)}} y^i$ と推定する。ここで $\mathcal{D}_a^{\ell(x)}$ は \mathcal{D}^{est} のうち $\ell(x)$ に属し行動が a のサンプル

- 目的関数も推定分散を考慮して修正

- 以下の目的関数を**最大化**するように領域 Π を学習

- $-\text{EMSE}_{\tau}(\Pi) = \mathbb{E}_{x^i} \left[\tau^2(x^i; \Pi) \right] - \mathbb{V}_{x^i, \mathcal{D}^{\text{est}}} \left[\hat{\tau}(x^i; \mathcal{D}^{\text{est}}, \Pi) \right]$ の推定として以下の目的関数を得る

- $-\widehat{\text{EMSE}}_{\tau}(\mathcal{D}^{\text{tr}}, \Pi) := \frac{1}{N^{\text{tr}}} \sum_{i \in \mathcal{D}^{\text{tr}}} \hat{\tau}^2(X_i; \mathcal{D}^{\text{tr}}, \Pi) - \left(\frac{1}{N^{\text{tr}}} + \frac{1}{N^{\text{est}}} \right) \cdot \sum_{\ell \in \Pi} \left(\frac{(S_1^{\text{tr}}(\ell))^2}{p} + \frac{(S_0^{\text{tr}}(\ell))^2}{1-p} \right)$ ただし $(S_a^{\text{tr}}(\ell))^2$ は領域内の各群の結果の経験分散

- $\hat{\tau}$ は平均が入るので、それが領域ごとに異なっているほど x ごとの異質性を捉えており近似精度が高い（第1項）

かつ領域内の分散は小さいほど推定分散が小さい（第2項）

(参考) 因果木の目的関数の導出

推定用データと訓練データの独立性を利用して 訓練の楽観バイアスと推定分散を推定できる

- 分割 Π に関して、 \mathcal{D}^{est} を用いて前ページのように $\hat{\tau}$ を定義し、これを \mathcal{D}^{te} で評価した場合のMSEを考える

$$\bullet \text{MSE}_{\tau}(\mathcal{D}^{\text{te}}, \mathcal{D}^{\text{est}}, \Pi) := \frac{1}{N^{\text{te}}} \sum_{i \in \mathcal{D}^{\text{te}}} \left\{ (\tau^i - \hat{\tau}(x^i; \mathcal{D}^{\text{est}}, \Pi))^2 - (\tau^i)^2 \right\}$$

- ここで最後の $-(\tau^i)^2$ は解析の都合で付与した、モデルによらない定数

- \mathcal{D}^{est} 、 \mathcal{D}^{te} に関して期待値をとる

$$\bullet \text{EMSE}_{\tau}(\Pi) := \mathbb{E}_{\mathcal{D}^{\text{te}}, \mathcal{D}^{\text{est}}} [\text{MSE}_{\tau}(\mathcal{D}^{\text{te}}, \mathcal{D}^{\text{est}}, \Pi)]$$

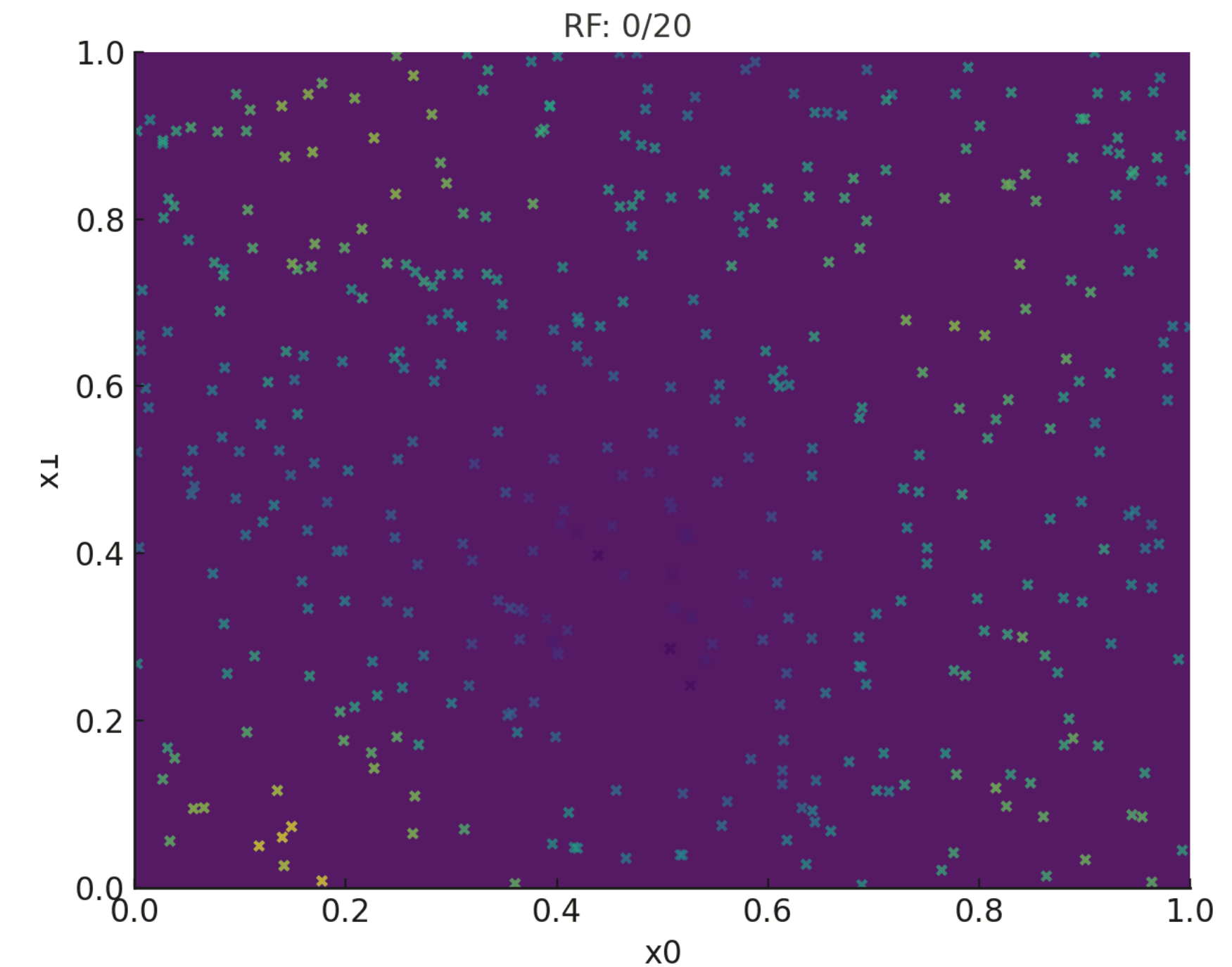
- 負号をつけて式展開する。領域 $\ell(x)$ における真の平均因果効果を $\tau(x^i; \Pi)$ として分解して

$$\begin{aligned} -\text{EMSE}_{\tau}(\Pi) &= -\mathbb{E}_{(x^i, \tau^i), \mathcal{D}^{\text{est}}} [(\tau^i - \hat{\tau}(x^i; \mathcal{D}^{\text{est}}, \Pi))^2 - (\tau^i)^2] \\ &= -\mathbb{E}_{(x^i, \tau^i)} [(\tau^i - \tau(x^i; \Pi))^2 - (\tau^i)^2] \\ &\quad - \mathbb{E}_{x^i, \mathcal{D}^{\text{est}}} \left[(\hat{\tau}(x^i; \mathcal{D}^{\text{est}}, \Pi) - \tau(x^i; \Pi))^2 \right] \\ &= -\mathbb{E}_{x^i} \left[-2\mathbb{E}[\tau^i | x^i] \tau(x^i; \Pi) + (\tau(x^i; \Pi))^2 \right] \quad \leftarrow \mathbb{E}[\tau^i | \ell(x^i)] = \tau(x^i; \Pi) \\ &\quad - \mathbb{V}_{x^i, \mathcal{D}^{\text{est}}} [\hat{\tau}(x^i; \mathcal{D}^{\text{est}}, \Pi)] \\ &= \mathbb{E}_{x^i} [\tau^2(x^i; \Pi)] - \mathbb{V}_{x^i, \mathcal{D}^{\text{est}}} [\hat{\tau}(x^i; \mathcal{D}^{\text{est}}, \Pi)] \end{aligned}$$

因果森：因果木のアンサンブル

- ランダム森（教師あり学習法）
 - 決定木をベース学習器とするアンサンブル法
 1. データからより小さいサンプルを重複ありで再抽出し決定木を学習する
 - 木の多様性を高める（木どうしの相関を下げる）ため
 2. これを繰り返し、出力は各木の出力の平均とする
- 因果森（Causal Forest）
 - 因果木をベース学習器 $\tau_j(x)$ とするランダム森

- $$\hat{\tau}(x) = \frac{1}{m} \sum_{j=1}^m \hat{\tau}_j(x)$$



ベイズ的決定森

BARTはベイズ推論 (MCMC) による決定森 もともとは因果推論用ではないが、因果推論にもよく使われる

● ベイズ推論

- ベイズの定理を用いてデータからパラメタの分布を推論する

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta')p(\theta')d\theta'} p(\theta) \quad \cdots(\star)$$

- 予測時はパラメタに関して期待値をとる： $\hat{y} = \mathbb{E}_{\theta \sim p(\theta|D)} \mathbb{E}[y|x, \theta]$

- 実際には $p(\theta|D)$ からサンプリングしたパラメタ集合 $\{\theta_1, \dots, \theta_K\}$ の平均をとるなど

- 実際には(★)式の分母の積分は実質的に計算不能

→ マルコフ連鎖モンテカルロ (MCMC) 法

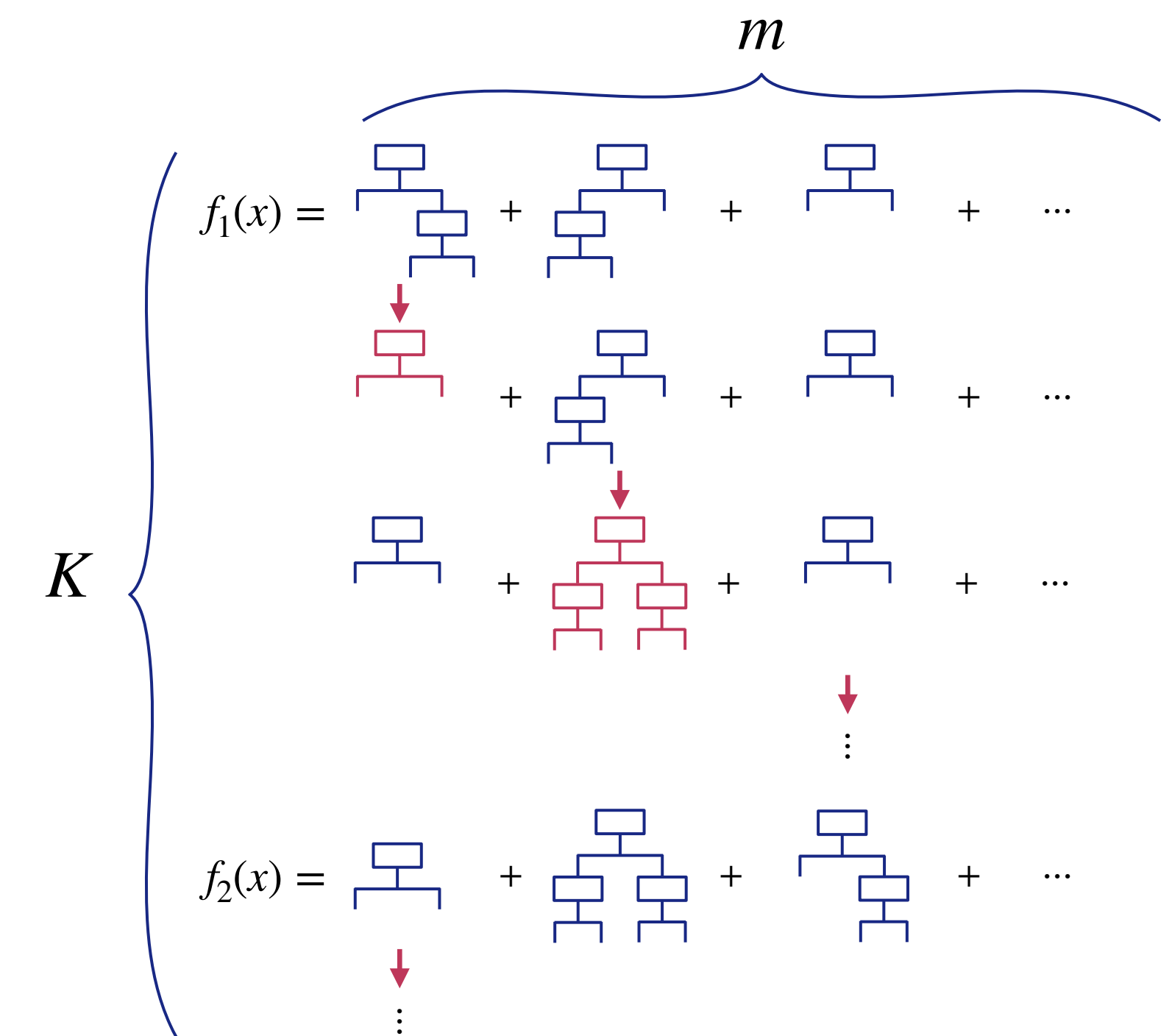
- あるパラメタ θ_{t-1} からランダムに少し動かした次のパラメタ θ^* に遷移するか止まるかを

事後確率の比 $p(\theta_t|D)/p(\theta_{t-1}|D)$ (計算可能) に基づいて確率的に決める (1以下ならその確率で遷移: メトロポリス法) ことを繰り返すと長期的には $\{\theta_t\}_t$ は $p(\theta|D)$ からサンプリングしたものと一致

- [MCMC Interactive Gallery](#)

● BART (Bayesian Additive Regression Trees)

- 決定森を1つのパラメタ θ としてMCMCを行う回帰モデル
- これをS-Learnerとして因果推論に用いることが提案されている [Hill 2011]



木を1つずつ更新する
(backfitting MCMC)

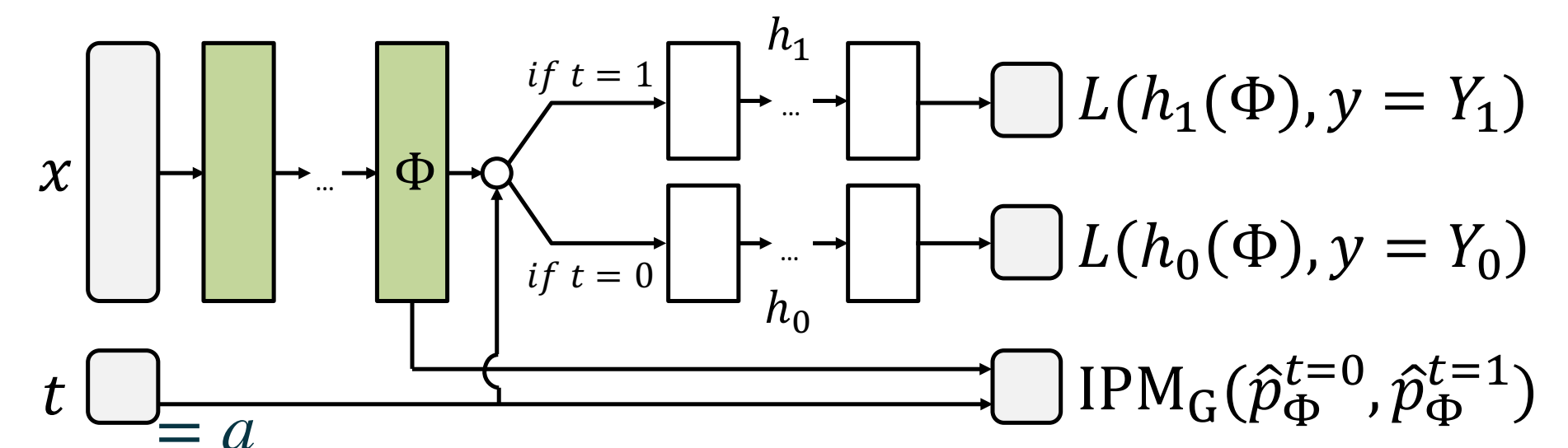
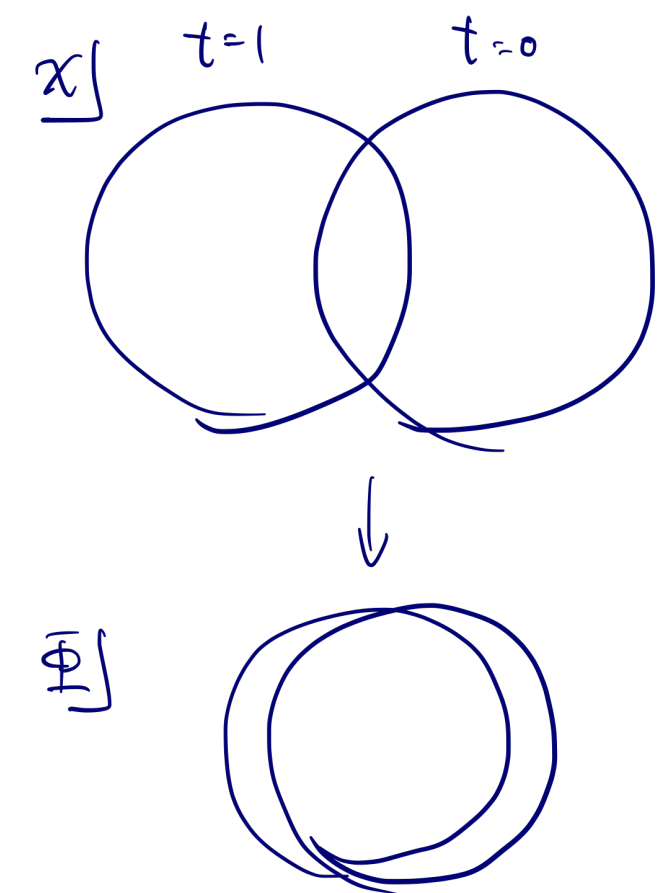
深層学習（表現学習）に基づく手法

深層学習による敵対的ドメイン適応として解く

- 実績分布が傾向スコア $\mu(a|x)$ によって偏ることが問題
- CounterFactual Regression (CFR)
 - 行動 a によって分布変化しない表現 $\phi_x = \phi(x)$ を抽出し、
 $p(\phi_x|a=0) \simeq p(\phi_x|a=1)$ になるようにすれば、
 ϕ_x 空間上で構築したモデルはバイアスがなくなるはず
- CFRは積分確率計量 Integral Probability Metric (IPM) によって
 分布間不一致を測り、これを同時最小化（損失に加える）

$$\text{IPM}_G(p_1, p_2) := \sup_{g \in G} \left| \int_{\Phi} g(\phi)(p_1(\phi) - p_2(\phi)) d\phi \right|.$$

- 関数クラス G の設定は1-リップシッツ関数などが用いられる
 - この場合のIPMはワッサーズタイン計量とも呼ばれる
- 分類器 g を同時に学習する（敵対的学習）
 - （上記supの双対問題と呼ばれる別問題を毎回解く方法もある）



(参考) 敵対的ドメイン適応の理論

(観測できない) x による損失の差を上限 \sup で置き換える

● 理論

- 表現抽出器 ϕ が逆関数を持つなら、一様分布上の損失 MSE^u が、データ分布上の損失 MSE と分布間IPMを用いて上から抑えられる
- $\rightarrow \text{MSE} + \text{IPM}$ を最小化すれば MSE^u が抑えられる

反事実 (観測されなかった
潜在結果) の誤差

$$\text{MSE}_{CF}(h, \Phi) \leq \underbrace{(1 - u)\text{MSE}_F^{t=1}(h, \Phi) + u\text{MSE}_F^{t=0}(h, \Phi)}_{\text{事実に (観測された結果) 誤差}} + B_\Phi \cdot \text{IPM}_G(p_\Phi^{t=1}, p_\Phi^{t=0}),$$

$\downarrow p(a=1)$

事実に (観測された結果) 誤差

- Φ 上の真のモデル h と損失 L (MSE) の合成 $\ell_{h,\Phi}(x, t) = \int_{\mathcal{Y}} L(Y_t, h(\Phi(x), t))p(Y_t|x)dY_t$ が関数クラス G に入るように定数(パラメタ) B_Φ を決めれば、その未知の関数 $\ell_{h,\Phi}$ に関して最悪 \sup をとれば上界になる

まとめ

決定森と深層学習を用いた手法

- 因果木は出力を τ とし分割学習用と予測値推定用にデータ分割した決定木
- 因果森は因果木をベース学習器とするランダム森
 - ランダム森はデータをサブサンプリングして学習した決定木を平均したもの
- BARTはベイズ推論を決定森モデルに適用したもの
- 深層学習を用いるとIPMによる正則化を敵対的学習で実現可
 - 重み付けとは異なるアプローチ