

## 統計的機械学習（応用計量分析2）第10回

---

代理変数法（参考pdf 12.3節）

回帰不連続デザイン（参考pdf 13章）

- CausalML2025\_CATE\_solution.ipynb
- 問1 : R-Learner (DML) の実装間違い
  - $Y - \mathbb{E}[Y | X, W] = \theta(X) \cdot (a - \mathbb{E}[a | X, W]) + \varepsilon$
  - $\mathbb{E}[Y | X, W]$ の推定は $a$ を用いない
- 問2 : X-Learnerの有利条件
  - 曝露・統制間のサンプル不均衡など
- 問3 : 任意の手法の有利・不利

```
def fit(self, X, T, Y):
    if self.cv > 1:
        m_hat, e_hat = self._oof(X, T, Y)
    else:
        XT = np.column_stack([X, T.reshape(-1, 1)])
        # self.m.fit(XT, Y) # 間違い
        self.m.fit(X, Y) # 正解
        self.g.fit(X, T)
        # m_hat = self.m.predict(XT) # 間違い
        m_hat = self.m.predict(X) # 正解
        e_hat = clip01(self.g.predict_proba(X)[ :, 1])
    y_tilde = Y - m_hat
    t_tilde = T - e_hat
    eps = 1e-3
    z = y_tilde / (t_tilde + np.sign(t_tilde) * eps)
    w = t_tilde ** 2 + 1e-6
    self.tau.fit(X, z, sample_weight=w)
    return self
```

# 振り返り

未観測交絡因子によりバックドア基準を満たせない場合の手法を2つ紹介  
フロントドア調整は媒介変数で分解して総合、操作変数法は外生変数をランダムマイザとみなす

- 未観測交絡因子があってもバックドア基準を満たせる場合もある
- どうしても満たせない場合、または未観測交絡の疑いがいくらかでも出てくる場合は別のアプローチも考えられる
- 2つのアプローチを紹介
  - フロントドア調整は未観測交絡因子から影響を受けていない、かつ $a \rightarrow y$ の因果効果を完全に媒介する変数を用いて分解して総合
  - 操作変数法は未観測交絡因子と独立かつ $a$ を通してしか $y$ に影響しない（かつ $a$ には影響する）変数をクジ（ランダムマイザ）として、その影響による $a$ の変動の $y$ への伝播を分析

# 本日の内容

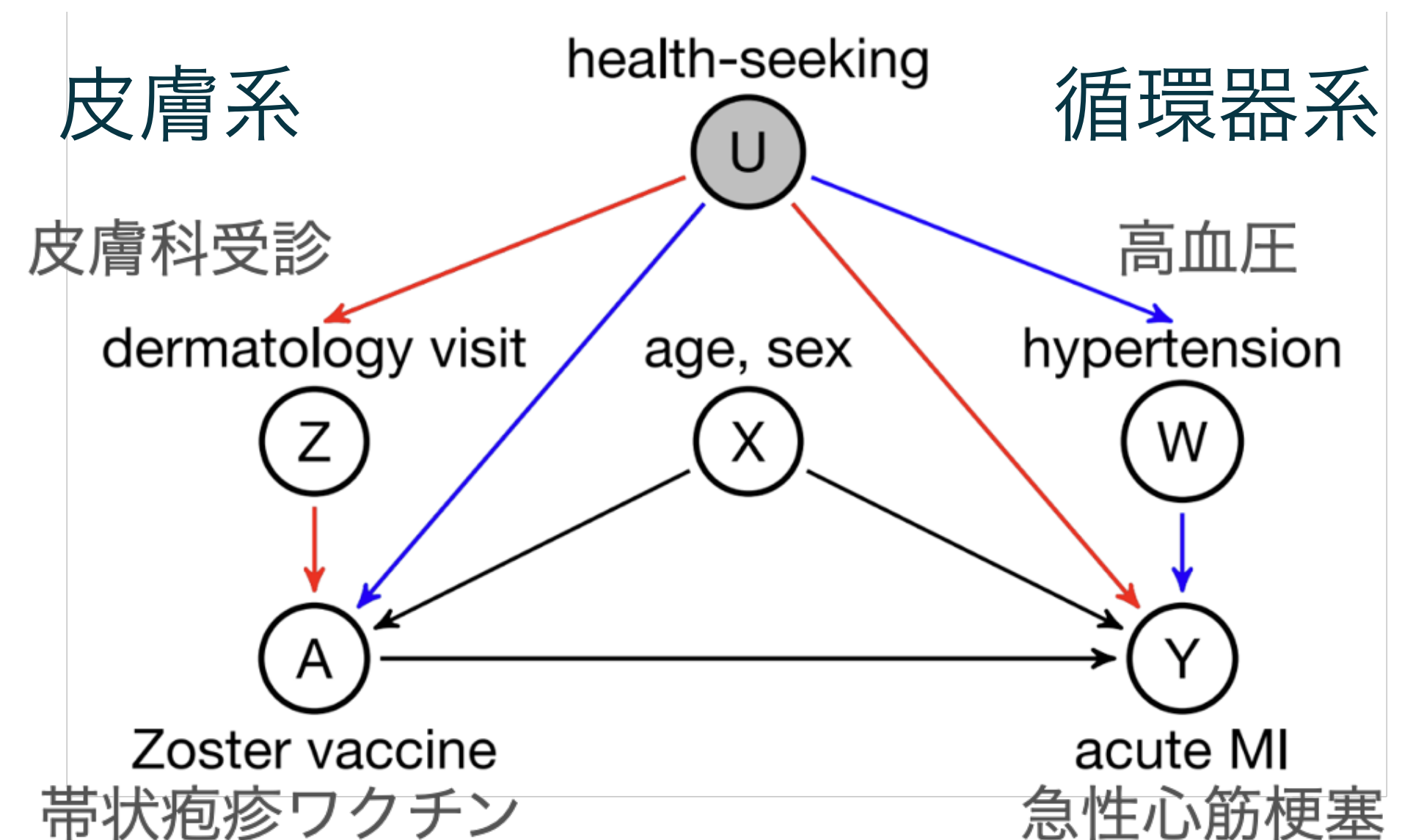
## 代理変数法・回帰不連続デザイン

- 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

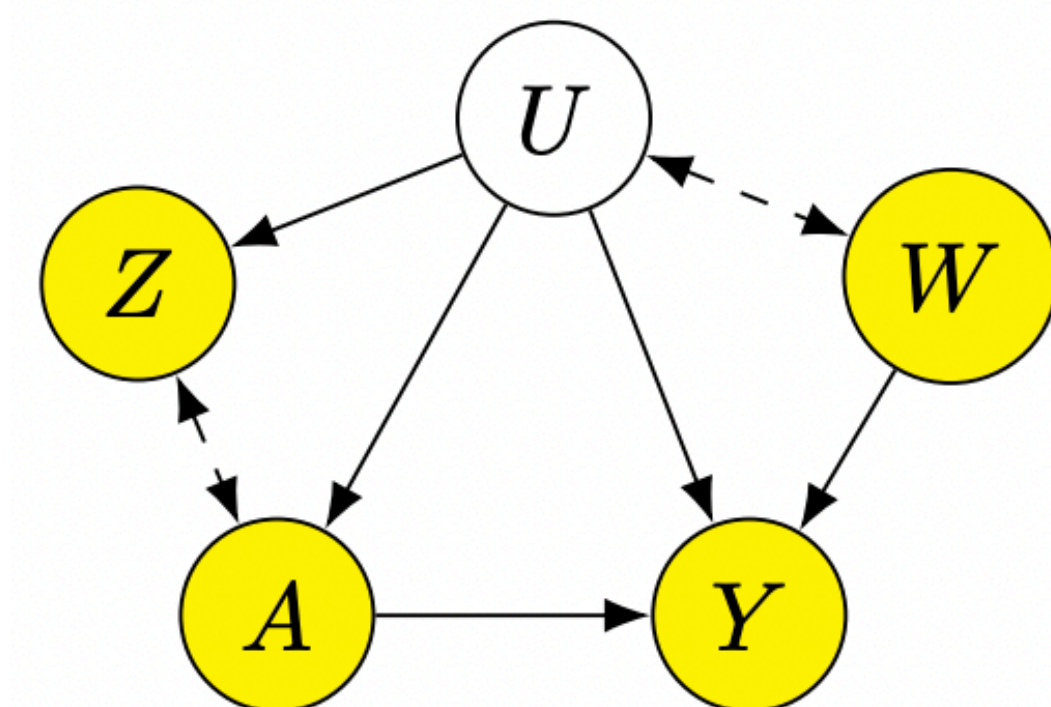
# 代理変数法とは

## 未観測交絡因子によるバックドアパスの効果を代理変数を通して推定して調整

- 例：帯状疱疹ワクチンと心筋梗塞
  - 元々高血圧の人が健康に対する意識が高く、ワクチンを打ちやすく、ワクチンとは無関係に高血圧のため心筋梗塞になりやすいかもしれない
- 未観測交絡因子 $u$ から $(a, y)$ と同様に影響を受けていて、行動 $a$ と結果 $y$ にそれぞれ関連する2つの因子を未観測交絡因子 $u$ の代理として用いる
  - 行動側の $z$ （陰性対照曝露）は $y$ と独立  $y \perp\!\!\!\perp z \mid a, u$  と仮定
  - 結果側の $w$ （陰性対照結果）は $a, z$ と独立  $w \perp\!\!\!\perp (a, z) \mid u$  と仮定
  - → 帯状疱疹は皮膚の病気なので皮膚科受診 $z$ と関係しそう、心筋梗塞と高血圧は循環器系なので相互には $u$ を通してしか関連しない
- 陰性対照 $z, w$ は $u$ の情報を十分に受け継いでいると仮定
- このときATE ( $x$ があればCATE) が識別可能



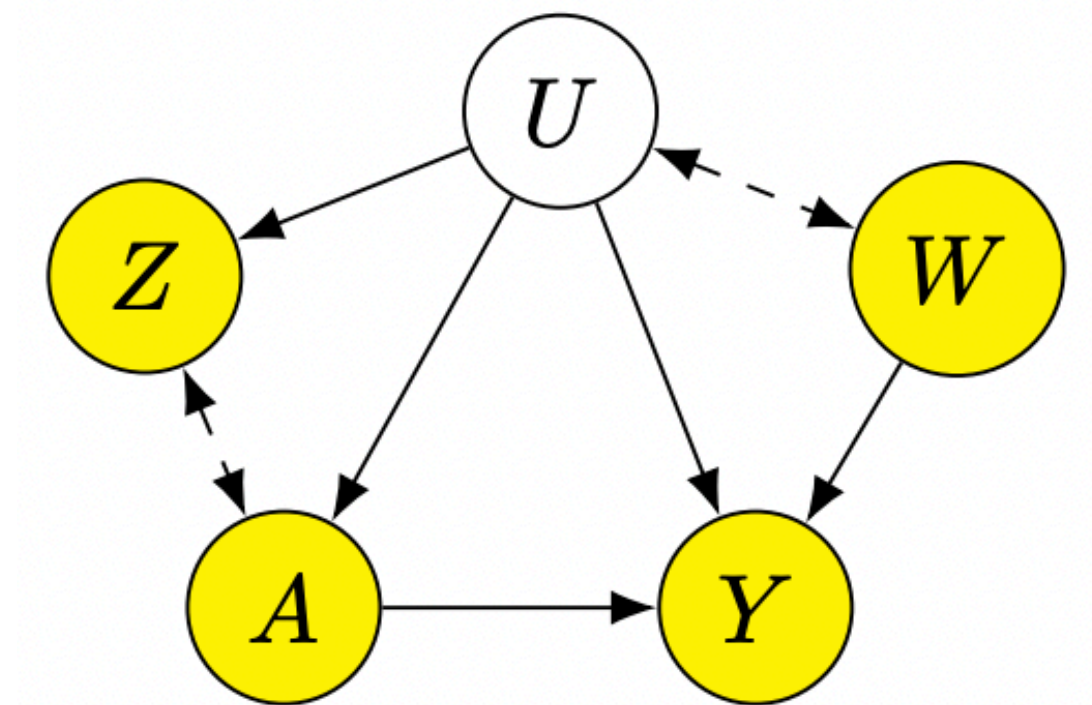
因果関係は一部逆でもよい



# 代理変数法（線形モデルの場合）

## $(z, a)$ から予測した $w$ を $u$ の代わりに説明変数として用いる

- 推定対象：下記の $\beta_{ay}$ 
  - $\mathbb{E}[y | a, z, u] = \beta_{y0} + \beta_{ay}a + \beta_{uy}u$ （独立性 $y \perp\!\!\!\perp z | a, u$ から $z$ の項は無し）
  - $u$ があればバックドア基準を満たすため、 $\beta_{ay}$ は因果効果
- 同様に $w$ については、以下で書ける
  - $\mathbb{E}[w | a, z, u] = \beta_{w0} + \beta_{uw}u$ （独立性 $w \perp\!\!\!\perp (a, z) | u$ から $a, z$ の項は無し）
- 上2式を  $p(u | a, z)$  について期待値を取り  $\mathbb{E}[u | a, z]$  を消去すると
  - $\mathbb{E}[y | a, z] = \beta_{y0} + \beta_{ay}a + \frac{\beta_{uy}}{\beta_{uw}} (\mathbb{E}[w | a, z] - \beta_{w0})$
  - 係数を整理して  $\mathbb{E}[y | a, z] = \beta'_{y0} + \beta_{ay}a + \beta_{wy}\mathbb{E}[w | a, z]$
- 結局、 $w$ を $a, z$ に回帰して得た予測値 $\hat{w}$ を $u$ の代わりに用いて回帰する



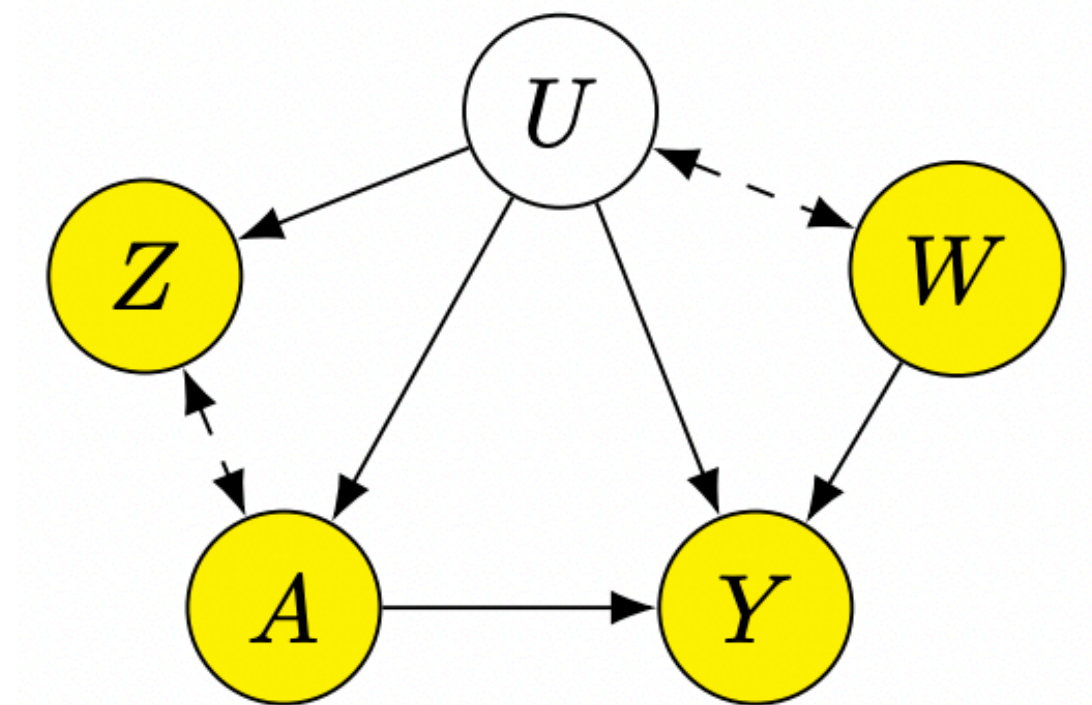
# 代理変数法の直観

## 余計な相関を代理変数で除去

- 用いている式

- $$\mathbb{E}[y | a, z] = \beta_{y0} + \beta_{ay}a + \frac{\beta_{uy}}{\beta_{uw}} (\mathbb{E}[w | a, z] - \beta_{w0})$$

- 知りたい因果効果 $\beta_{ay}$ とバックドアの相関 $AU \times UY$ を分離したい
- $AU \times UW$ は $\mathbb{E}[w | a, z]$ の $a$ の係数として推定可能
  - しかし $UY/UW$ だけスケールが間違っている
- $UY, UW$ それぞれは推定不可能だが、その比は推定可能
  - $UY/UW = (ZU+AU) \times UY / (ZU+AU) \times UW$
  - $= \mathbb{E}[y | a, z] / \mathbb{E}[w | a, z]$



# 回帰不連続デザイン

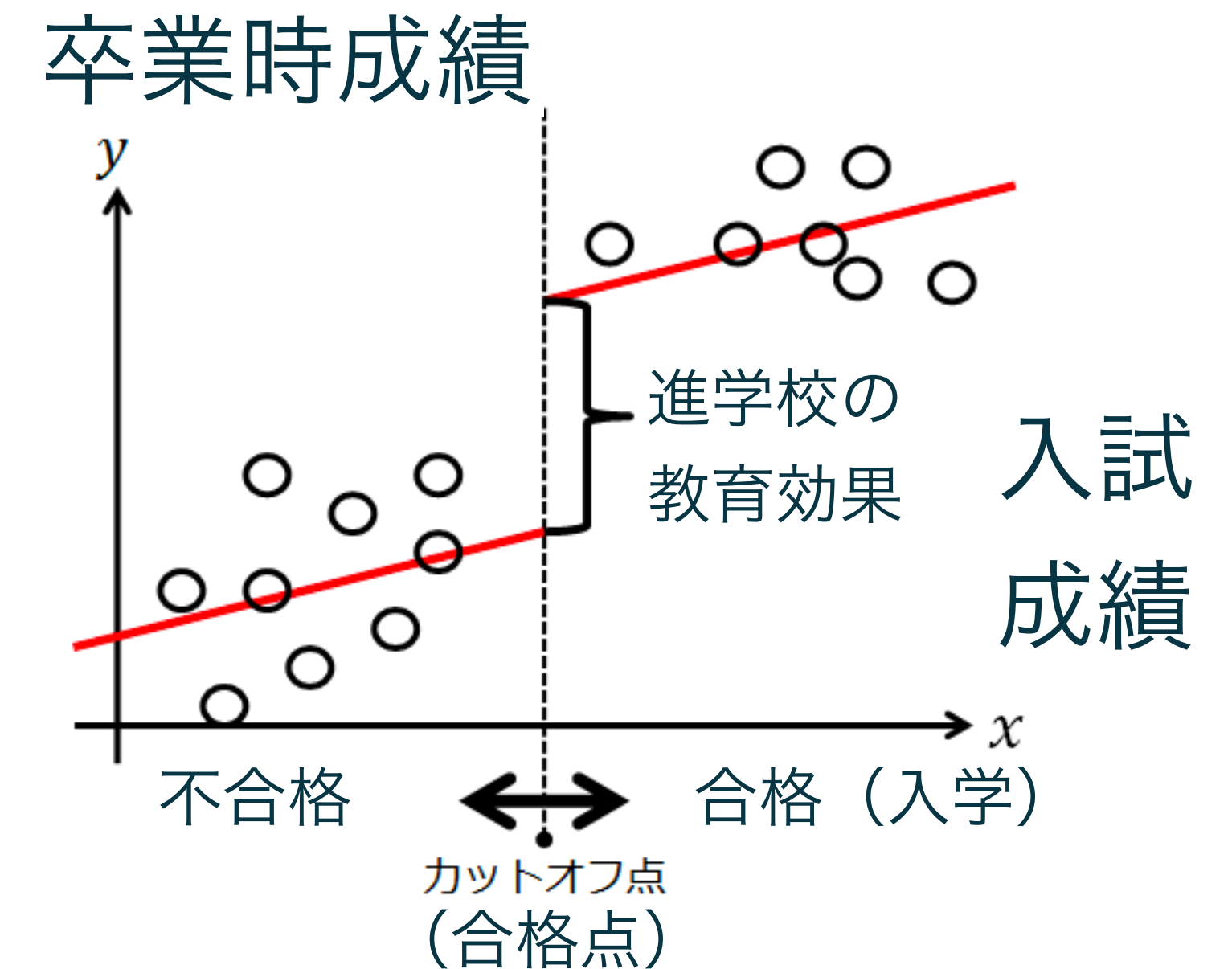
---

正值性  $\mu(a|x) > 0$  が満たされない場合

決定論的に  $a$  が変化する点の周辺を比較する

## 割り当てカットオフ点周辺はランダムとみなせることを利用

- 正值性  $\mu(a|x) > 0 \quad \forall(x, a)$  が満たされない場合  
→  $x$  全体では推定のしようがない
  - 中でも決定論的に割り当てが決まる場合を考える
  - $a = \begin{cases} 1 & g(x) \geq c \\ 0 & g(x) < c \end{cases}$
- カットオフ点  $z := g(x) = c$  の周りではほとんどランダムに  $a$  が決まる  
→ それらの差は ( $z = c$ での) 因果効果
  - カットオフ点における平均因果効果 (ATE at the cutoff)
$$\tau_c = \mathbb{E}[y_1 | z = c] - \mathbb{E}[y_0 | z = c] = \lim_{z \downarrow c} \mathbb{E}[y | z] - \lim_{z \uparrow c} \mathbb{E}[y | z]$$
  - 例：進学校の教育効果を、入試成績がギリギリ合格・不合格の群を比較



<https://ja.wikipedia.org/wiki/回帰不連続デザイン>

# RDDの詳細

## 真の関数が連続的+カットオフ点周りの正值性を仮定 関数をフィットするパラメトリック法とカットオフ付近を平均するノンパラメトリック法

### ● 仮定

- 合格点付近が「割り当て以外同様の人」と見なせる必要あり

$$\lim_{z \downarrow c} E[y_a | z] = \lim_{z \uparrow c} E[y_a | z] \quad \forall a \in \{0,1\}$$

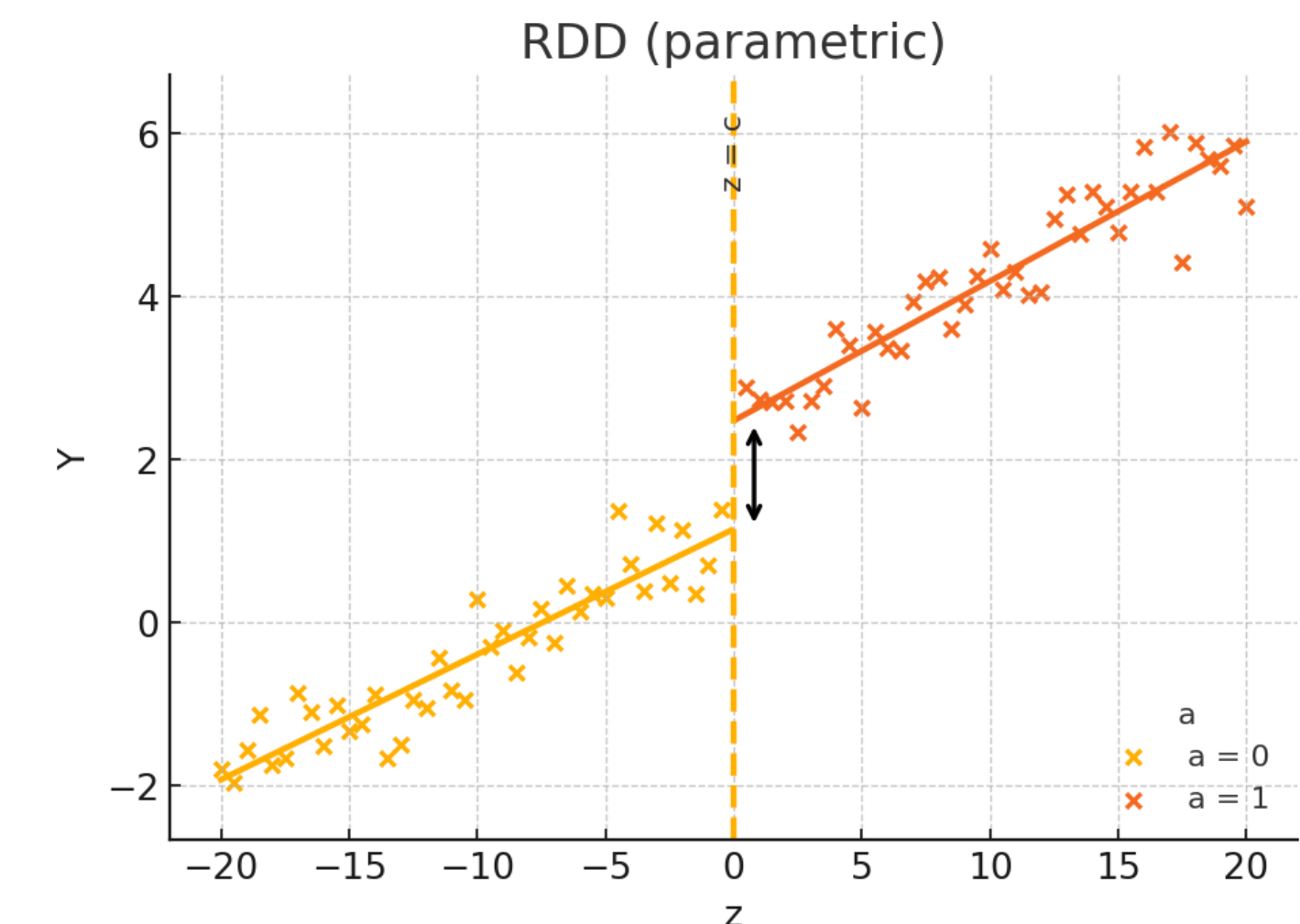
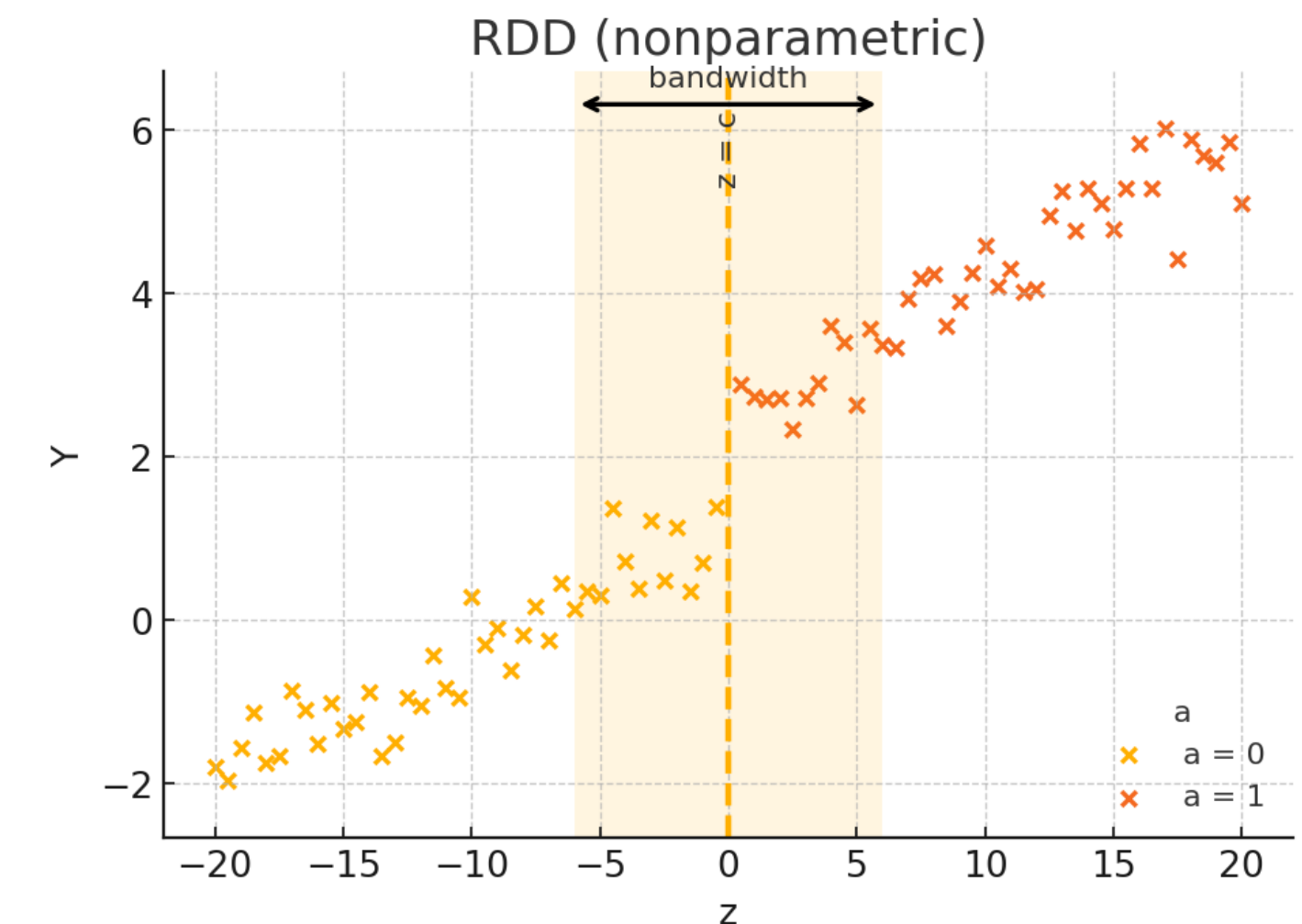
- 18歳から大学に入るしパチンコもできる場合、  
大学入学の効果かパチンコの効果かわからなくなる
- カットオフ点周りにサンプルが存在：  $p(z = c) > 0$

### ● 手法2：ノンパラメトリックな方法

- 一定のバンド幅  $\delta$  をとり、その領域の平均間を比較

### ● 手法1：パラメトリックな方法

- カットオフ点以後かどうか ( $a = I(z \geq c)$ ) で分割
- $a$  ごとに重回帰等をフィット



## まとめ

# 未観測交絡因子への対処法（その3）代理変数法 正值性を満たさない決定論的行動割り当て状況の手法RDD

- 代理変数法は未観測交絡を通じたバックドアパスの相関を代理変数を用いて推定して差し引く
- RDDはカットオフ点周りで決定論的に行動の割り当てが決まる場合に、その周辺を比較することでカットオフ点周りの因果効果を推定する
  - パラメトリックな手法とノンパラメトリックな手法