

統計的機械学習（応用計量分析2）第9回

未観測交絡因子への対処法（参考pdf 12章）

言語モデルの知識を用いる手法

データだけから因果構造全体を復元するのは難しい問題

→ 知識を言語モデルから抽出する

- データのみから因果構造全体を復元することはかなり難しく、実用的には人間がある程度の構造を事前知識として与えることが多い
- 「(変数名A)を変更すると(変数名B)に変化が生じますか？」と大規模言語モデルに質問
 - 変数名が存在すること、その因果関係が一般的知識として知られていることが前提
 - 大規模言語モデルにその一般的知識が学習され埋め込まれている前提
 - 変数ペアの因果方向判定タスクで既存のデータに基づく手法を凌駕
- ただし新たな知識構築（科学の発展）にはデータも必要

Model	Acc.	Wt. Acc.
Slope (Marx & Vreeken, 2017)	0.75	0.83
bQCD (Tagasovska et al., 2020)	0.68	0.75
PNL-MLP (Zhang & Hyvarinen, 2012)	0.75	0.73
Mosaic (Wu & Fukumizu, 2020)	0.83	0.82
ada	0.50	0.50
text-ada-001	0.49	0.50
babbage	0.51	0.50
text-babbage-001	0.50	0.50
curie	0.51	0.52
text-curie-001	0.50	0.50
davinci	0.48	0.47
text-davinci-001	0.50	0.50
text-davinci-002	0.79	0.79
text-davinci-003	0.82	0.83
LMPrior (Choi et al., 2022)	0.83	-
gpt-3.5-turbo	0.81	0.83
gpt-3.5-turbo (causal agent)	0.86	0.87
gpt-3.5-turbo (single prompt)	0.89	0.92
gpt-4 (single prompt)	0.96	0.97

Kiciman, Emre, et al. "Causal reasoning and large language models: Opening a new frontier for causality." Transactions on Machine Learning Research (2023).

振り返り

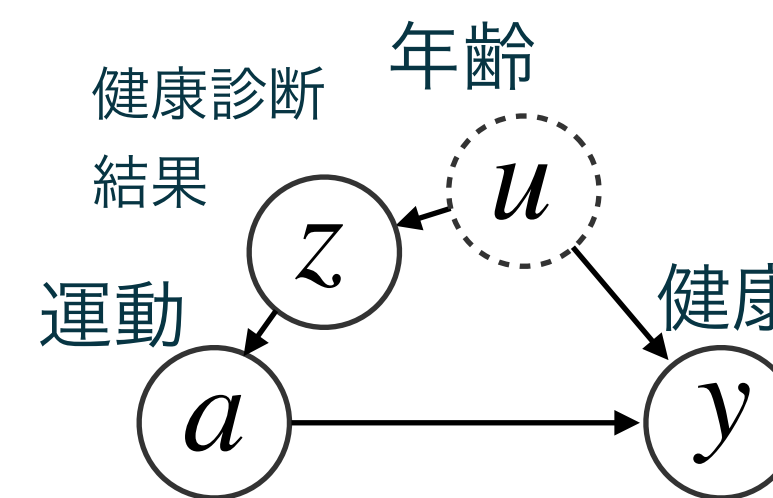
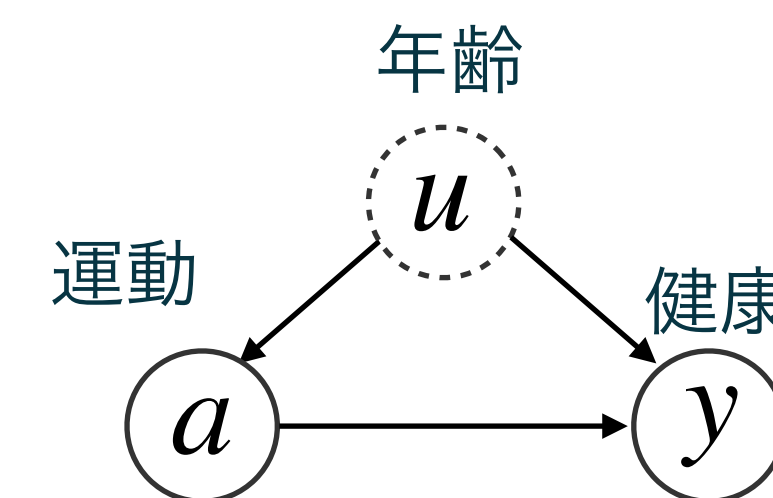
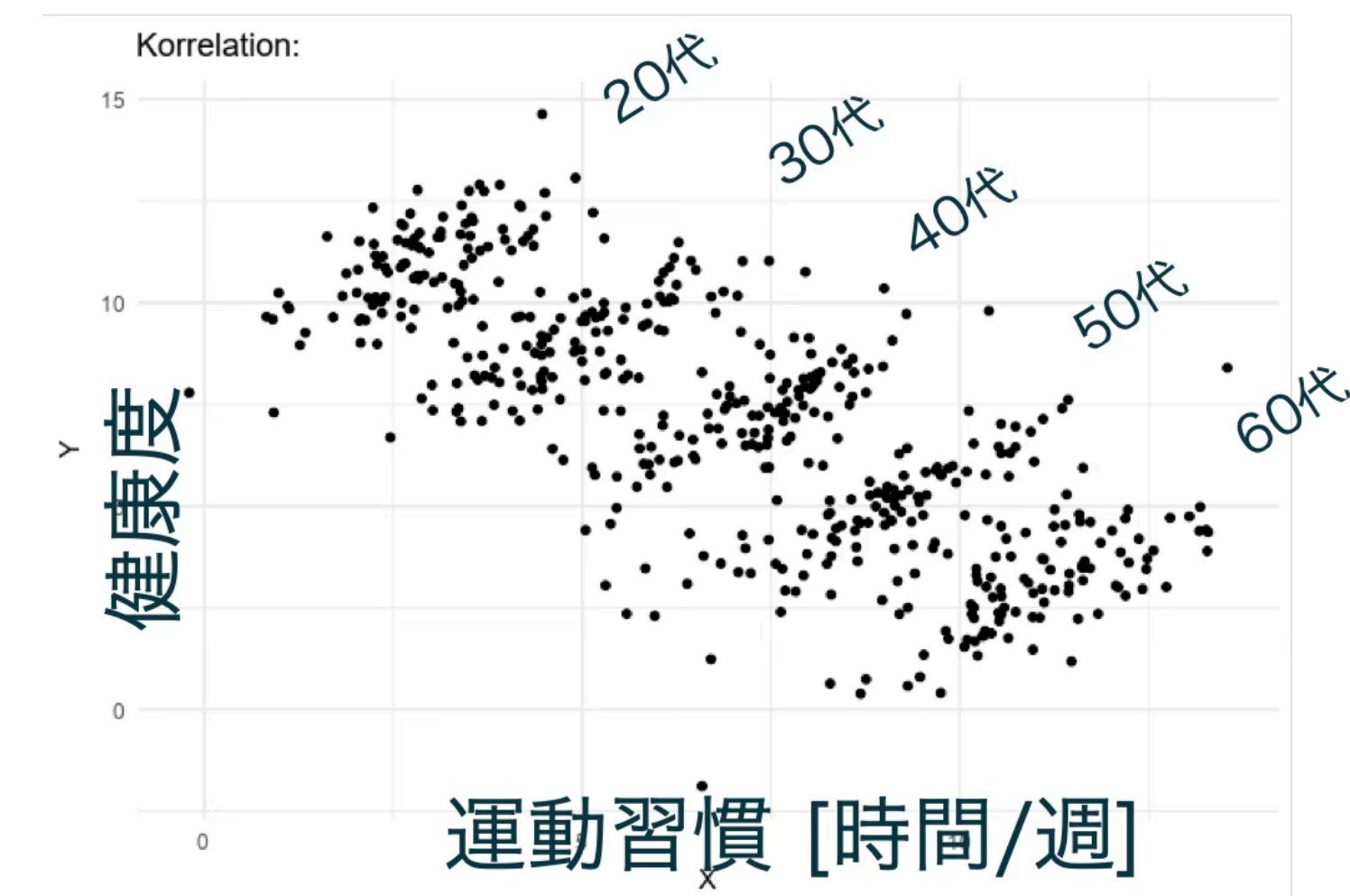
因果構造を推定する問題＝因果探索

仮定の置き方と得られる結論の強さにより幾つかのアプローチ

- 観察データのみから因果構造を完全に特定することは不可能
 - 一部の辺の向きは諦めるか、追加でモデル等に仮定を置くことで識別する
- IC/PC法
 - (条件付き) 独立性のみに基づく辺の枝刈りとV構造を用いた向き付け
 - 仮定が少ないが、向きが判定できない辺が残る可能性がある (CPDAGまで識別)
- LiNGAM
 - 線形性を仮定、ノイズ分布はガウス分布でないことを仮定
- NOTEARS
 - 構造方程式 f のモデルをノイズ含めて仮定
 - トレース指数関数により非巡回制約を連続空間で表現
- 言語モデルからの知識抽出
 - 変数ペアの因果関係を問い合わせることで知識を抽出

未観測交絡因子の問題は重大

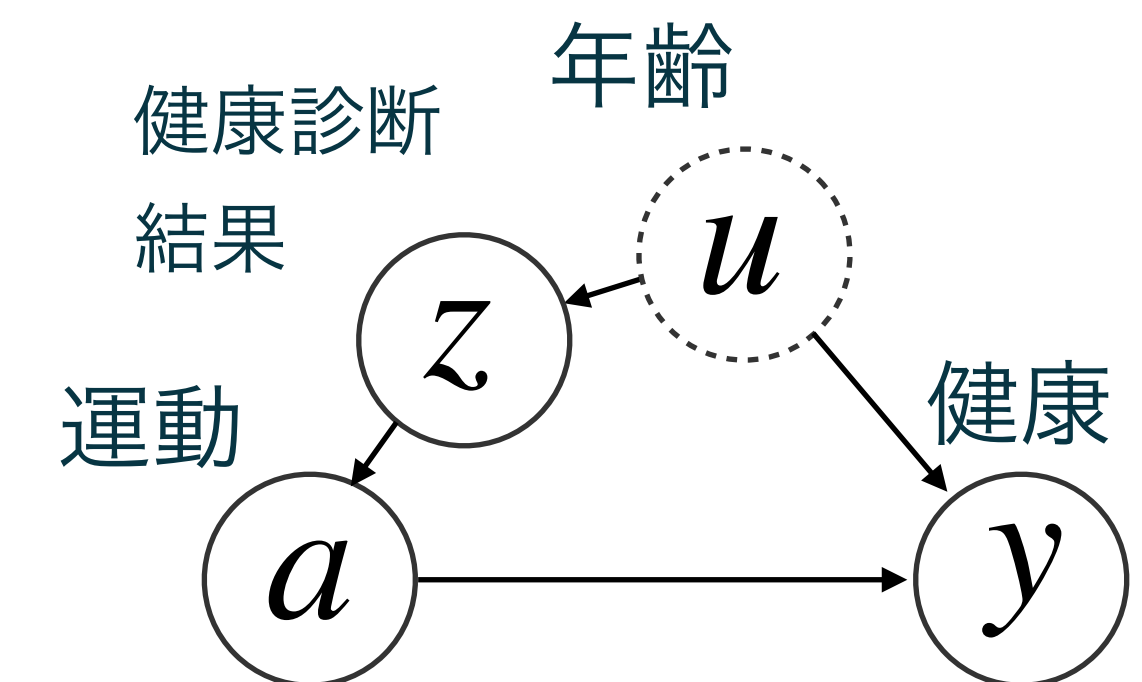
- 交絡因子 (u) が観測できていなかったら
一般に調整できない
 - 無視可能性を満たすには交絡因子の観測が必要
 - $(y_{a'})_{a' \in \mathcal{A}} \perp\!\!\!\perp a \mid u$
- ただし、バックドア基準を満たせる場合もある
 - 右下図の健康診断結果 z を調整することでも代替可
- 現実の問題では交絡因子の候補も不明なことも
 - データだけで未観測交絡因子の不存在を確認する術はない
 - ※ LiNGAMを仮定すると未観測交絡因子の”検出”は可能な場合も
- そのような場合にも識別可能な方法は？



バックドア基準を用いて確認

因果構造の情報を用いて調整する変数集合を選ぶ基準

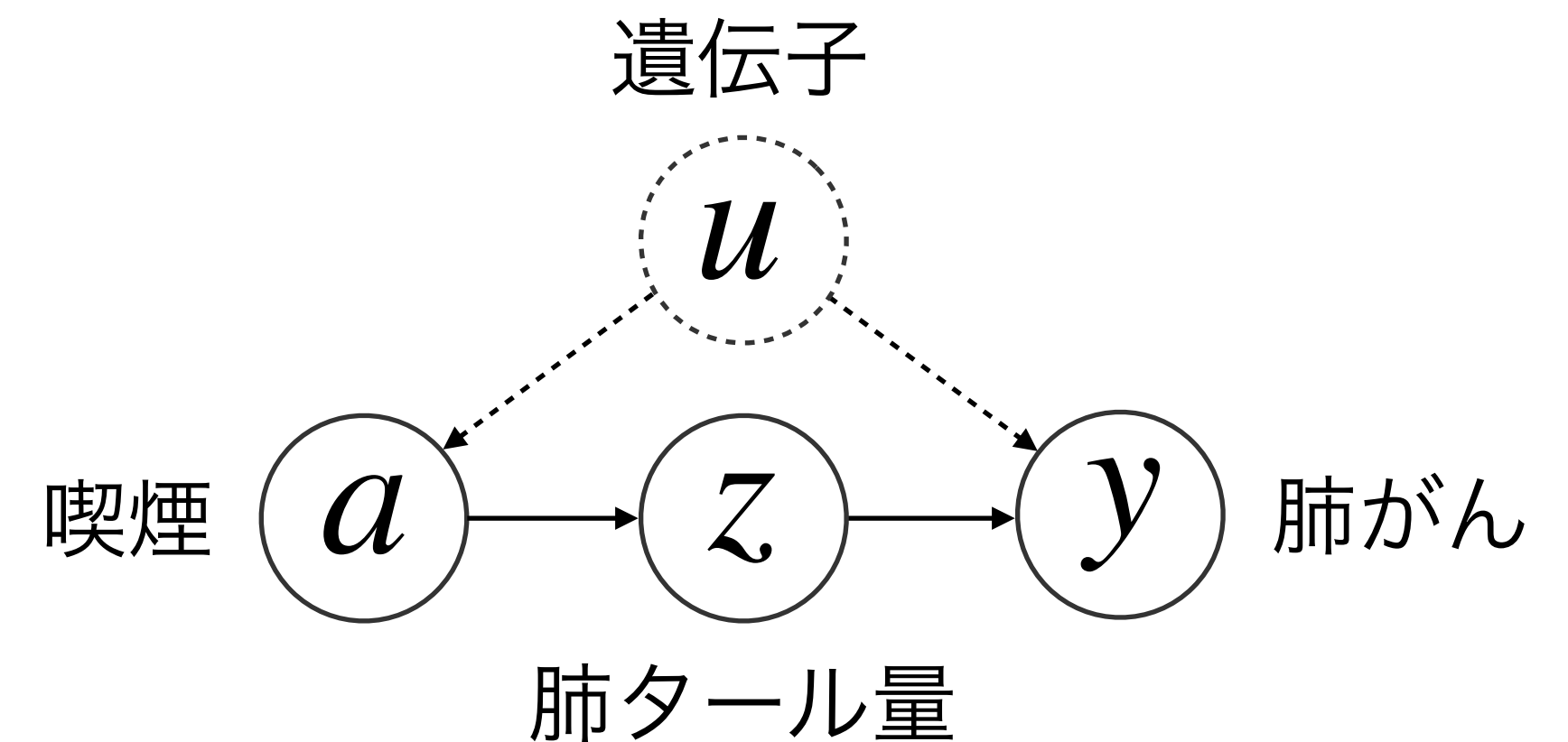
- 健康診断結果 z はバックドア基準を満たす
 - バックドアパスの鎖上の変数を用いてブロック
- バックドア基準（再掲）
 - 因果ダイアグラム G において、 a から y への有向パスがあるとする。次の2条件を満たすとき、**変数集合 Z は順序対 (a, y) についてバックドア基準を満たす**という
 - (B1) a から Z への有向パスがない（行動より下流の変数を含まない）
 - (B2) a に入るパスを含む、 a と y を結ぶパス（**バックドアパス**）において、 Z が a と y を **有向分離**（ブロック）する
 - ただし、 a - y 間の全てのパス p に対して Z が以下の条件のいずれかを満たすとき、 Z は a と y を **有向分離する**という
 - 鎖 $i \rightarrow m \rightarrow j$ またはフォーク $i \leftarrow m \rightarrow j$ を含み、 m は Z に含まれる
 - 合流点 $i \rightarrow m \leftarrow j$ を含み、 m 及びその子孫は Z に含まれない



フロントドア調整とは

フロントドア調整は因果パスを組み合わせる バックドア基準の応用

- 例：喫煙→肺がんの因果機序の論争（1950—'60年代）
 - 「喫煙を欲させる遺伝子 u がある」主張
 - ※実際には、リスク増加率が大きく遺伝子だけでは説明つかないことなどで結論
- 因果機序を分解する
 - 喫煙 a → 肺のタール沈着量 z → 肺がん y
 - 遺伝子→肺のタール沈着量とは考えづらい
- $a \rightarrow z$ と $z \rightarrow y$ はそれぞれ識別可能
 - $a \rightarrow z$ はバックドアパス無し
 - $z \rightarrow y$ はバックドアパスがあるが、 a が有向分離する
- 媒介変数 z を介した因果機序を推定して組み合わせれば $a \rightarrow y$ の因果効果ATEは識別可能



フロントドア調整の定理

フロントドア基準が満たされればフロントドア調整によってATEが識別される

- フロントドア基準
 - 変数の集合 z が順序対 (a, y) に関して**フロントドア基準**を満たすとは、以下の3つの条件を満たすこと
 1. z は a から y へのすべての有向パスを遮断する
 2. a から z へのバックドアパスが存在しない
 3. z から y へのすべてのバックドアパスが a によってブロックされる
- z がフロントドア基準を満たすとき、 $a \rightarrow z$ と $z \rightarrow y$ は以下

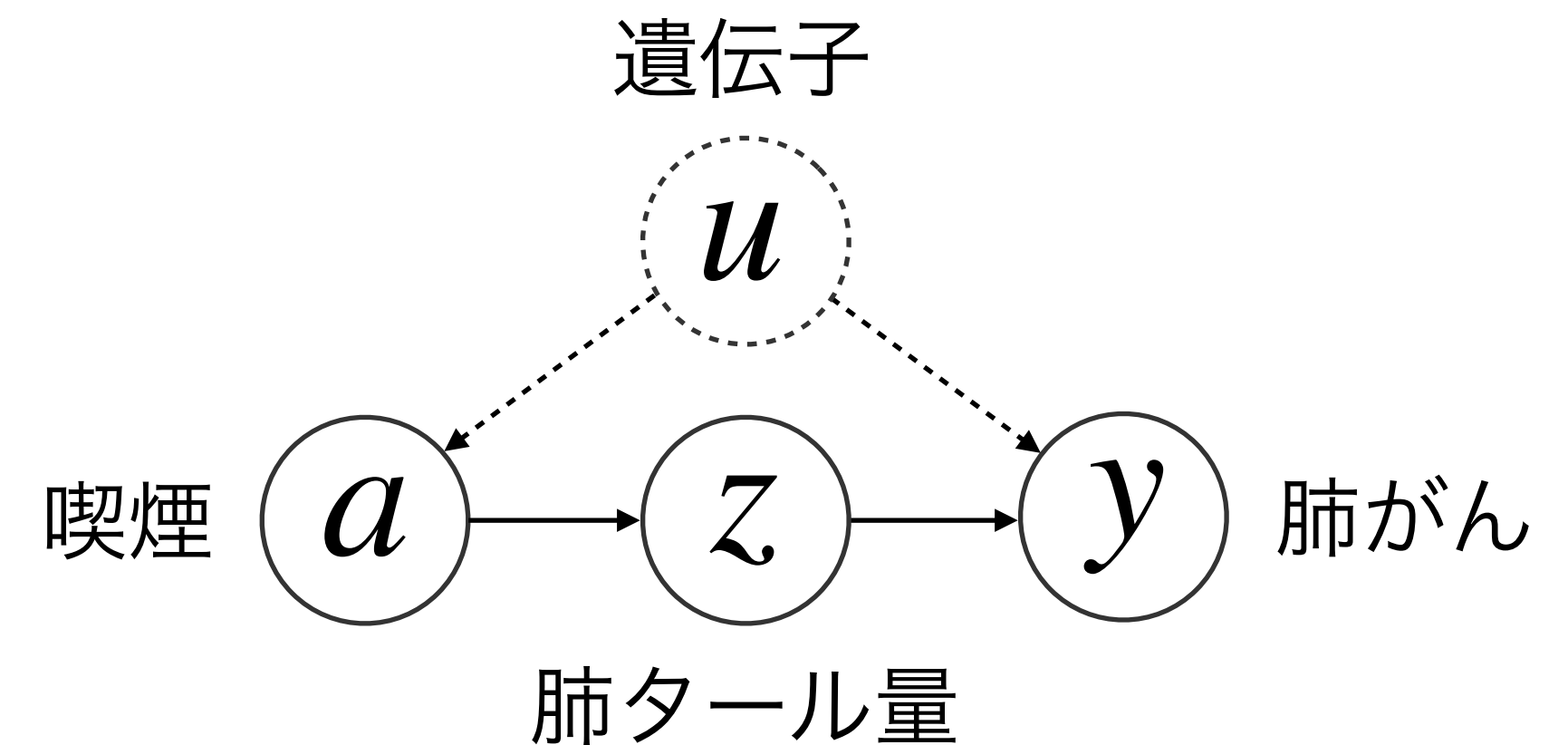
$$p(z | \text{do}(a)) = p(z | a)$$

$$p(y | \text{do}(z)) = \sum_a p(y | a, z)p(a)$$

- これらを組み合わせると以下が成立

$$p(y | \text{do}(a)) = \sum_z p(z | \text{do}(a))p(y | \text{do}(z))$$

$$= \sum_z p(z | a) \sum_{a'} p(y | a', z)p(a')$$



フロントドア調整の具体的方法（線形モデル）

$a \rightarrow z$ と $z \rightarrow y$ をそれぞれ推定して総合する

- 線形モデルを仮定

- $\mathbb{E}[y|a, u] = \beta'_{y0} + \beta_{ay}a + \beta_{uy}u$ — (★)

- $\beta_{uy}u$ の項が邪魔

- 分解する

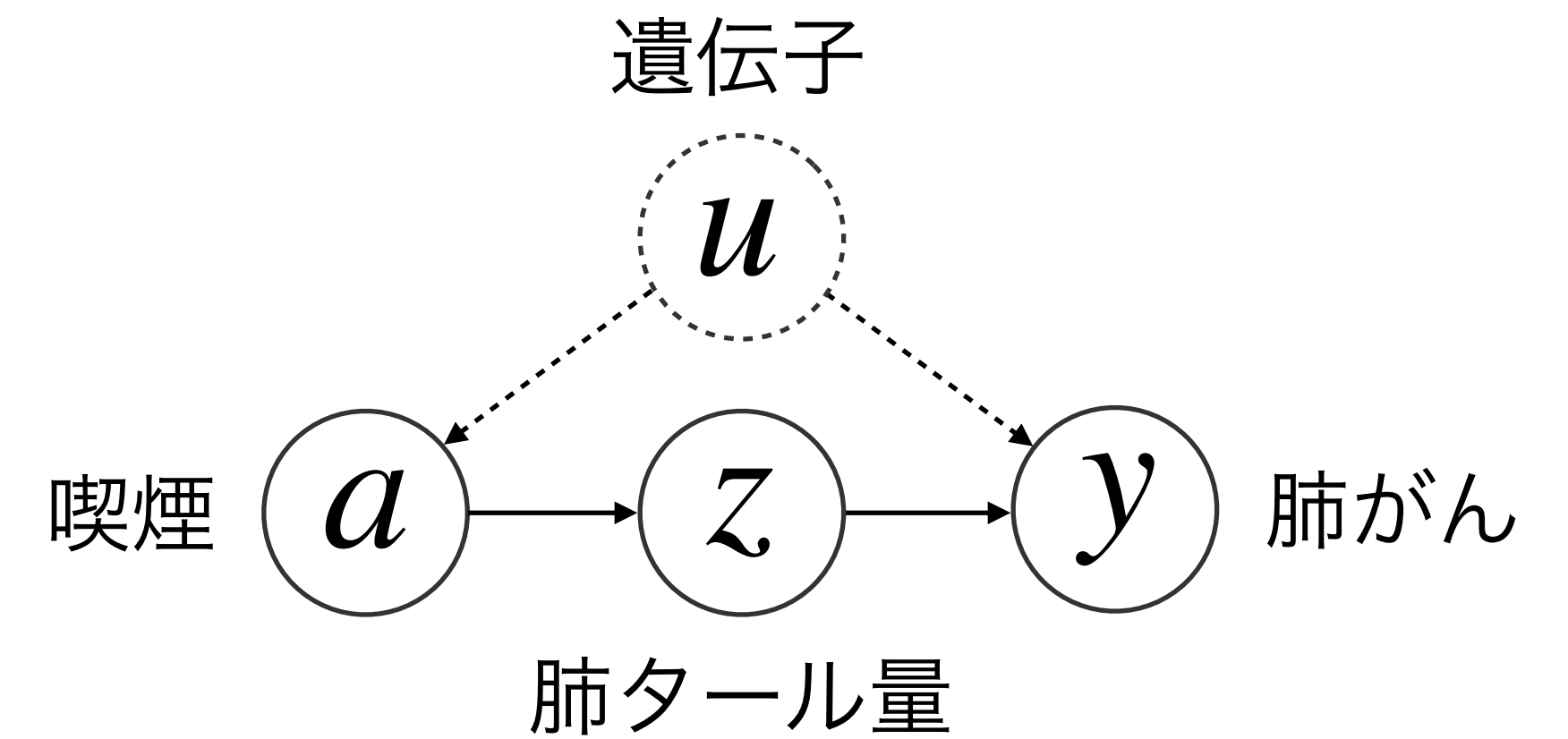
- $\mathbb{E}[z|a] = \beta_{z0} + \beta_{az}a$

- $\mathbb{E}[y|z, a] = \beta_{y0} + \beta_{zy}z + \beta'_{ay}a$...推定可能

- 式展開

$$\begin{aligned} \mathbb{E}[y|a, u] &= \mathbb{E}_z[\underbrace{\mathbb{E}[y|z, a, u]}_{=\mathbb{E}[y|z, u]} | a, u] \\ &= \mathbb{E}_z[\beta''_{y0} + \beta_{zy}z + \beta_{uy}u | a, u] \\ &= \beta''_{y0} + \beta_{zy} \underbrace{\mathbb{E}[z|a, u]}_{=\mathbb{E}[z|a]} + \beta_{uy}u \\ &= \beta''_{y0} + \beta_{zy}(\beta_{z0} + \beta_{az}a) + \beta_{uy}u \end{aligned}$$

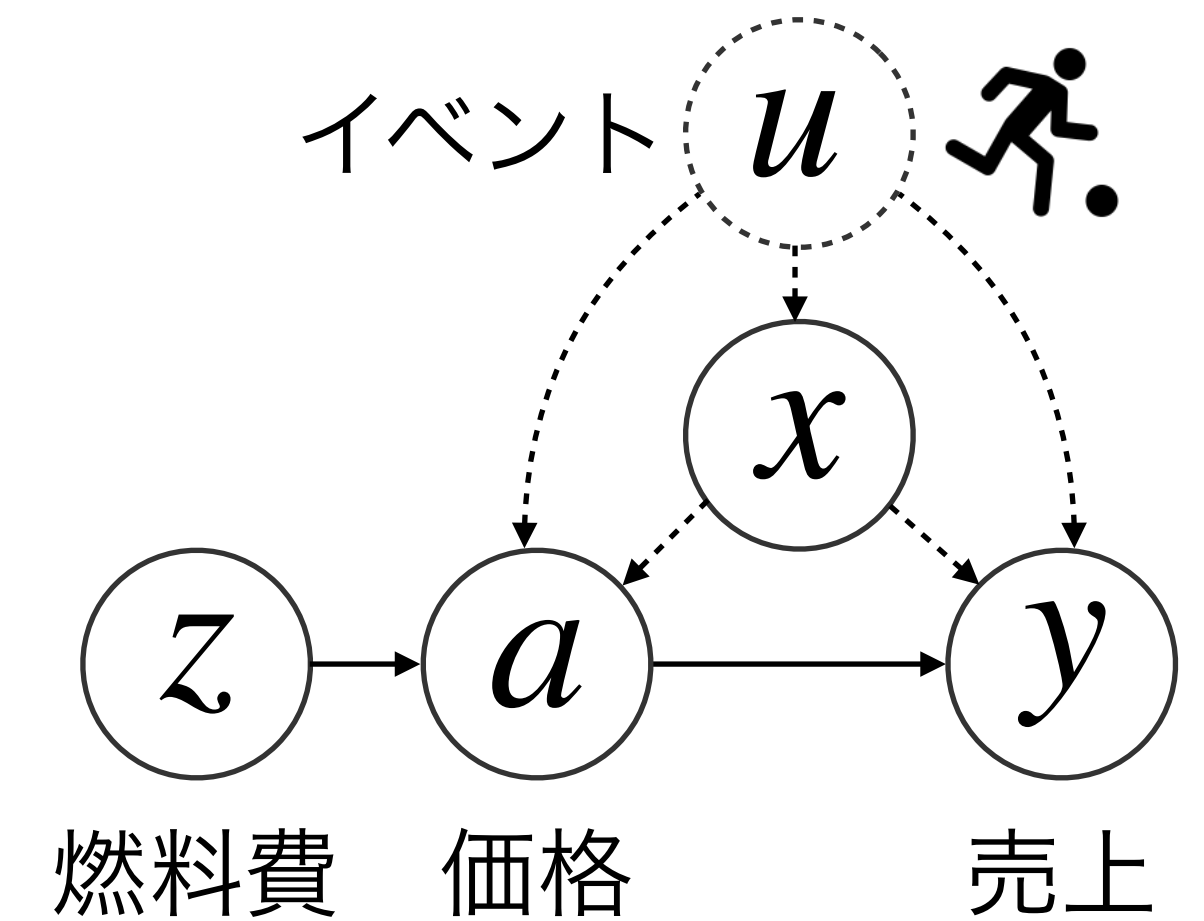
- (★) と比較して、 $\beta_{ay} = \beta_{zy}\beta_{az} \rightarrow$ 係数をそれぞれ推定して積をとればよい



操作変数 (Instrumental Variable ; IV) 法とは

独立した変数 (外生変数) を弱いクジとしてその影響を分析

- 例：航空機チケット価格弾力性
 - 近隣でイベント u があると売上が上昇
 - イベント情報は得られないとする
 - 早期から売れ行きが良いのでアルゴリズムにより価格設定が上昇
- 交絡因子と独立した行動決定要因 z を仮定
 - 燃料費 (原油価格) はイベントと独立 $z \perp u$ かつ
チケット価格を通してしか売上に影響しない (除外制約)
 - 燃料費 $z \rightarrow$ 売上 y の因果効果は価格 $a \rightarrow$ 売上 y の因果を通ったもの
 - 燃料費 $z \rightarrow$ 価格 a の因果効果と併せて考えれば価格 $a \rightarrow$ 売上 y がある程度わかる
- 2021年ノーベル経済学賞はIV法関連
 - IV法を用いて教育効果を分析したカード
 - IV法の理論的基盤を整備したアングリスト・インベンス



III. Niklas Elmehed © Nobel Prize Outreach.

デイビッド・カード教授



III. Niklas Elmehed © Nobel Prize Outreach.

ヨシュア・アングリスト教授



III. Niklas Elmehed © Nobel Prize Outreach.

グイド・インベンス教授

操作変数法（線形モデル）

$z \rightarrow a$ と $z \rightarrow y$ の分解を用いて式展開

- 線形モデルを仮定

- $\mathbb{E}[y|a, u] = \beta_{y0} + \beta_{ay}a + \beta_{uy}u$

- u を条件に入ればバックドア基準を満たし
 β_{ay} が求めたい因果効果。ただし $\beta_{uy}u$ の項が邪魔

- 簡単のため x は無いとする

- 操作変数 z から a, y への因果効果

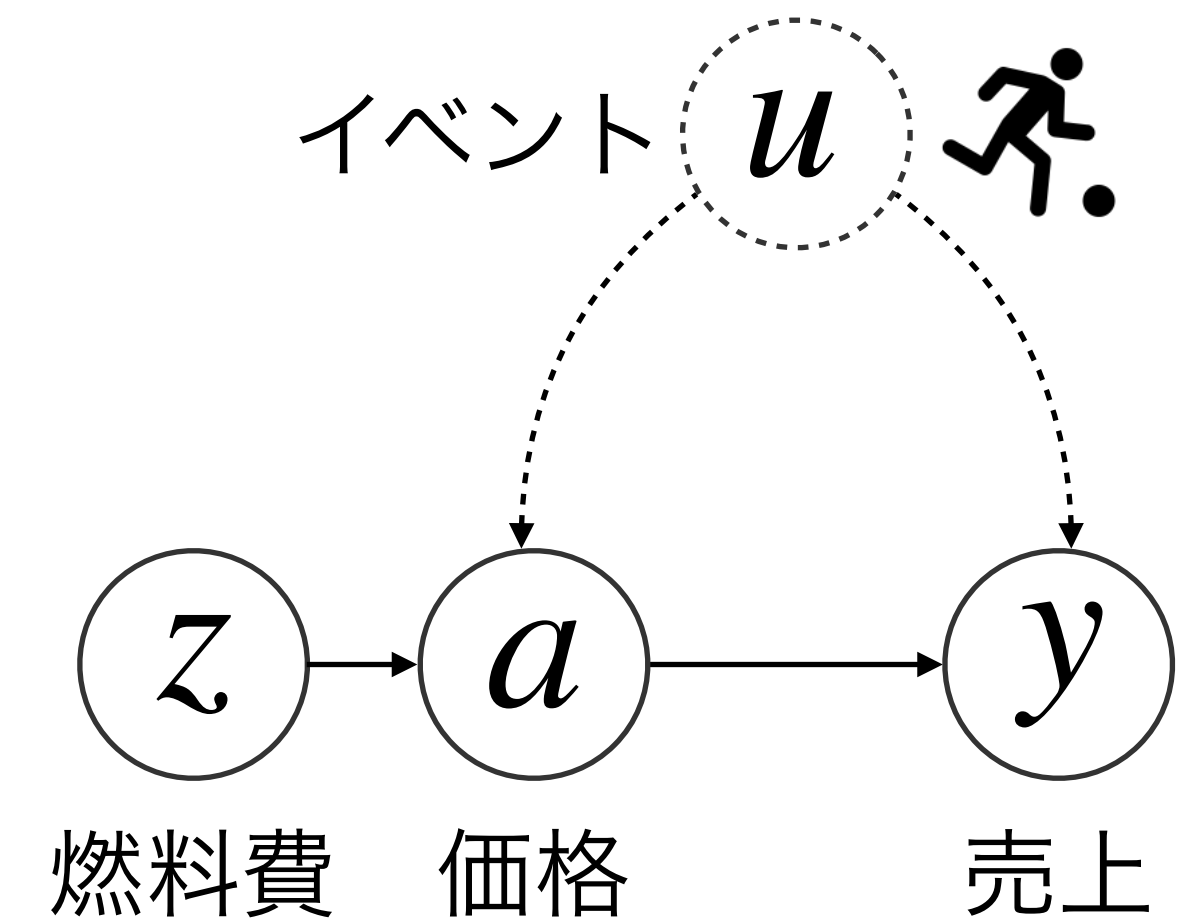
- $\mathbb{E}[a|z] = \beta_{a0} + \beta_{za}z$

- $\mathbb{E}[y|z] = \beta'_{y0} + \beta_{zy}z$

- ここで、 $z \rightarrow y$ を $z \rightarrow a$ と $a \rightarrow y$ に分解して表式

- $\mathbb{E}[y|z] = \mathbb{E}_{a,u}[\mathbb{E}[y|a, z, u]|z] = \mathbb{E}_{a,u}[\mathbb{E}[y|a, u]|z]$

- 2つ目の等式は z と u の独立性から従う



操作変数法（線形モデル）の理論と手順

係数の比（変数数が同じ）か2段階最小二乗法

- 式展開

$$\begin{aligned}
 \mathbb{E}[y|z] &= \mathbb{E}_{a,u}[\mathbb{E}[y|a,u]|z] \\
 &= \beta_{y0} + \beta_{ay}\mathbb{E}[a|z] + \beta_{uy}\underbrace{\mathbb{E}[u|z]}_{=\mathbb{E}[u]} && \because \mathbb{E}[y|a,u] = \beta_{y0} + \beta_{ay}a + \beta_{uy}u \\
 &= \beta''_{y0} + \beta_{ay}\mathbb{E}[a|z] && \leftarrow (\star) \beta''_{y0} := \beta_{y0} + \beta_{uy}\mathbb{E}[u] \\
 &= \beta''_{y0} + \beta_{ay}(\beta_{a0} + \beta_{za}z) && \because \mathbb{E}[a|z] = \beta_{a0} + \beta_{za}z \\
 &= \beta'''_{y0} + \beta_{ay}\beta_{za}z
 \end{aligned}$$

- モデル $\mathbb{E}[y|z] = \beta'_{y0} + \beta_{zy}z$ と見比べると

- $\beta_{zy} = \beta_{ay}\beta_{za} \Rightarrow \beta_{ay} = \beta_{zy}/\beta_{za}$
 - a, z 共に1変数の場合は推定した係数の比を取るだけで推定可
- より一般には、 $\mathbb{E}[a|z]$ を回帰で予測した上で(★)に代入する**二段階最小二乗法**

操作変数法の推定対象

操作変数法はATEを識別しない 単調性を追加で仮定すれば局所因果効果LATEを識別する

- 線形モデルが正しければ β_{ay} はATEを識別する
 - $\mathbb{E}[y|a, u] = \beta_{y0} + \beta_{ay}a + \beta_{uy}u$
- しかしモデルが不正確な場合はそうはいえない
- 追加で**単調性**を仮定する
 - 各ユニット*i*において、 z^i の増加に対して a^i は非減少
 - 行動を z によって「促された」ら、（反応しないことはあっても）
「逆の行動をとる」あまのじゃく（非遵守者；Defiers）は居ない
- このとき、IV法は遵守者（Compliers）に限定した平均因果効果を識別
 - $c^i = 1 \Leftrightarrow a_{z=1}^i > a_{z=0}^i$ と定義して
 - $LATE = \mathbb{E}[y_{a=1} - y_{a=0} | c^i = 1]$
- 行政施策等（ z ）を打った際に、それに反応した人だけに限定した効果と解釈できる

まとめ

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