

cate - 条件付き平均処置効果 【 評価版 】

cate はいくつかの変数について条件付けが行われた形の平均処置効果 (ATE: average treatment effects) である条件付き ATE (CATE: conditional average treatment effects) の推定を行います。CATE の推定により処置効果の異質性 (heterogeneity) を評価することが可能になります。CATE はパラメトリックな回帰手法を用いて推定することもできますが、lasso やランダムフォレスト (random forest) といった新たな手法を用いた推定も可能です。

1. 機能概要	
1.1 CATE とは	
1.2 CATE の種別	
1.3 cate コマンド群	
2. ワークフロー	Workflow 1
	Workflow 2
	Workflow 3
	Workflow 4
	Workflow 5
3. 用例集	Example 1
	Example 2
	Example 3
	Example 4
	Example 5
	Example 6
	Example 7
	Example 8
補足 1	

1. 機能概要

処置効果 (treatment effects) というのはアウトカムに対する処置の因果的効果 (causal effect) を推定するものです。この効果は母集団中で一定のこともあれば、部分集団ごとに異なったものであることもあります。医療の分野であれば、個人のストレスレベルに対する喫煙の効果を年齢群別に評価するといった例が考えられます。

ATE というのは処置効果を母集団上で平均化したものです。これは処置効果が母集団上で均質な場合には有効な指標となります。しかし処置効果が均質ではなく、ユニットが同一の処置に対して異なる反応を示す場合には、平均値の推定のみでは不十分です。処置効果の異質性 (heterogeneity) にまで踏み込んだ分析が必要となります。

CATE を用いることによって処置効果の異質性に関する理解を深めることができます。ATE と同様、CATE も処置効果を平均化したものではあるのですが、その対象が母集団中のサブグループである点が異なります。ATE に比べてより高い解像度を持った顕微鏡が手に入ったことに相当します。また処置効果の異質性についてのより正確な情報は処置割当てプロセスの改善をもたらすものでもあります。

CATE を分析することによるメリットをまとめると次のようになります。

1. 処置効果の異質性を理解する上で助けとなる。
2. 処置割当てを最適化する上での基盤となる情報をもたらす。

1.1 CATE とは

潜在的アウトカムフレームワーク (potential-outcomes framework) のもとで、ユニット i が処置を受けたときの潜在的アウトカムを $y_i(1)$ 、ユニット i が処置を受けなかったときの潜在的アウトカムを $y_i(0)$ と書くことにします。またユニット i の特徴を表現するベクトルを \mathbf{x}_i で表すことにします。このとき CATE は次のように定義されることになります。

$$\text{CATE} \equiv \tau(\mathbf{x}) = E\{y_i(1) - y_i(0) \mid \mathbf{x}_i = \mathbf{x}\}$$

すなわち CATE は特徴 \mathbf{x}_i が \mathbf{x} であるという条件のもとでの潜在的アウトカム間の差の期待値ということができます。

$\tau(\mathbf{x})$ は部分的線形モデル (partial linear model)、あるいは完全交互作用モデル (fully interactive model) を介して特定することができます。ここでは部分的線形モデルを前提に話を進めることにします^{*1}。なお、以下においては添え字 i を省略し、観測されたアウトカムを単に y 、処置の状態を表す 2 値変数を単に d と表記することにします。

^{*1} 完全交互作用モデル (fully interactive model) については Example 5 を参照ください。

処置効果が一定である場合、あるいは関心の対象が ATE の推定にある場合の部分線形モデルは次のようになります。

$$\begin{aligned}y &= d * \tau + g(\mathbf{x}, \mathbf{w}) + \epsilon \\d &= f(\mathbf{x}, \mathbf{w}) + u\end{aligned}$$

ここでは変数を \mathbf{x} と \mathbf{w} の 2 つのグループに分けてある点に注意してください。区別のしかたについては追って説明します。アウトカムモデルについて見ると、観測されたアウトカムが処置効果 $d * \tau$ 、ノンパラメトリックな関数 $g(\mathbf{x}, \mathbf{w})$ 、誤差項 ϵ の和として表現されているため、部分的線形ということになります。一方、処置割当ては関数 $f(\mathbf{x}, \mathbf{w})$ に誤差項 u を加えたものとなっています。このとき、潜在的アウトカムのモデルは次のように表現できることとなります。

$$\begin{aligned}y(1) &= \tau + g(\mathbf{x}, \mathbf{w}) + \epsilon \\y(0) &= g(\mathbf{x}, \mathbf{w}) + \epsilon\end{aligned}$$

従って τ は ATE を規定するものとなります。

$$\text{ATE} \equiv E\{y(1) - y(0)\} = \tau$$

今度は処置効果が一定ではなく \mathbf{x} に依存するケースについて考えてみましょう。この場合、アウトカムモデルは次のようになります。

$$y = d * \tau(\mathbf{x}) + g(\mathbf{x}, \mathbf{w}) + \epsilon$$

ただし \mathbf{x} は処置効果を条件付ける変数ベクトルを意味します。 $\tau(\mathbf{x})$ は \mathbf{x} の関数で処置 d と交互作用を持つ形となっています。なお、 $\tau(\mathbf{x})$ は \mathbf{x} の関数ですが、 \mathbf{w} の関数ではない点に注意してください。このモデルは $\tau(\mathbf{x})$ や $g(\mathbf{x}, \mathbf{w})$ に対しパラメトリックな仮定を全く課していないので、柔軟で汎用性が高いモデル式であると言えます。このとき潜在的アウトカムは次のようになります。

$$\begin{aligned}y(1) &= \tau(\mathbf{x}) + g(\mathbf{x}, \mathbf{w}) + \epsilon \\y(0) &= g(\mathbf{x}, \mathbf{w}) + \epsilon\end{aligned}$$

従って $\tau(\mathbf{x})$ は CATE を意味するわけです。

$$E\{y(1) - y(0) | \mathbf{x}\} = \tau(\mathbf{x})$$

仮に $\tau(\mathbf{x}) = \mathbf{x}'\beta$ とか $g(\mathbf{x}, \mathbf{w}) = \mathbf{x}'\gamma_1 + \mathbf{w}'\gamma_2$ といった具合にパラメトリックな仮定を置くのであれば、回帰による推定が可能になります。しかしパラメトリックな仮定は強すぎてデータへの適合が困難となるようなケースも想定されます。従って cate ではパラメトリックな仮定を置かない、より一般的なケースに焦点を当てた実装としています。関心の対象である $\tau(\mathbf{x})$ の推定は Athey, Tibshirani, and Wager (2019) で提案されている一般化ランダムフォレスト (generalized random forest) を介してノンパラメトリックな形で実行されます。一方、 $g(\mathbf{x}, \mathbf{w})$ と $f(\mathbf{x}, \mathbf{w})$ については lasso、あるいはランダムフォレストによって推定が行われます。しかし $\tau(\mathbf{x})$ や $g(\mathbf{x}, \mathbf{w})$, $f(\mathbf{x}, \mathbf{w})$ に対してパラメトリックな前提を設けることも可能で、cate はそのようなモデルのフィットにも対応しています。

1.2 CATE の種別

評価版では割愛しています。

1.3 cate コマンド群

評価版では割愛しています。

2. ワークフロー

評価版では割愛しています。

3. 用例集

以下の用例では、処置効果の異質性を分析したり処置割当てポリシーを評価したりするのに cate をどのように使うかを紹介します。Example 1-7 においては 401(k) プログラム資格の金融資産形成に対する効果を cate を使って評価します。具体的に関心のある事項は次の通りです。

- (1) 金融資産形成に対する 401(k) プログラム資格の効果は均質と言えるのだろうか？言い換えるなら処置効果は個人やグループによって変動するのだろうか？
- (2) 処置効果が不均質であるとしたとき、それは想定されるグループ — 収入のカテゴリ、住宅の所有、教育レベル、等 — のレベルによって変動するのだろうか？
- (3) データによって処置効果の特に大きな、あるいは小さなグループを検出できるのだろうか？

Example 8 においては 2 種類の肺移植が患者の健康に及ぼす影響について評価します。具体的には処置割当てに関する推奨ルールが与えられたとして、そのルールを実際に試行したときの全体的な効果について評価したいものとします。

▷ Example 1: 処置効果の異質性

ここでは Example データセット `assets3.dta` を用いて 401(k) プログラム資格の金融資産形成に対する効果を推定してみます。

```
. use https://www.stata-press.com/data/r19/assets3.dta *2
(Excerpt from Chernozhukov and Hansen (2004))
```

*2 メニュー操作 : File ▷ Example Datasets ▷ Stata 19 manual datasets と操作、Causal Inference and Treatment-Effects Estimation Reference Manual [CAUSAL] の cate の項よりダウンロードする。

データセット中に含まれている変数の一覧を確認しておきます。

```
. describe
```

```
. describe
```

Contains data from <https://www.stata-press.com/data/r19/assets3.dta>
 Observations: 9,913 Excerpt from Chernozhukov and Hansen (2004)
 Variables: 11 27 Feb 2025 19:19
 (_dta has notes)

Variable name	Storage type	Display format	Value label	Variable label
assets	float	%9.0g		Net total financial assets
age	byte	%9.0g		Age
income	float	%9.0g		Household income
educ	byte	%9.0g		Years of education
pension	byte	%16.0g	lbpen	Pension benefits
married	byte	%11.0g	lbmar	Marital status
twoearn	byte	%9.0g	lbytes	Two-earner household
e401k	byte	%12.0g	lbe401	401(k) eligibility
ira	byte	%9.0g	lbytes	IRA participation
ownhome	byte	%9.0g	lbytes	Homeowner
incomecat	byte	%9.0g		Income category

Sorted by: e401k

データセット上には 9,913 件の世帯主に関するデータが記録されています。その中でキーとなる変数は 401(k) プログラム資格の有無を表す指標変数 (0/1 変数) e401k と純金融資産 (net financial assets) の額を表す assets です。この他、関係する変数をリストアップしておくと次のようになります。

- incomecat - 収入レベルを表すカテゴリ変数。0/1/2/3/4 というコードで表される。
- age - 年齢
- educ - 教育歴を意味する年数 (1-18)
- pension - 年金受給者か否かを表す指標変数
- married - 婚姻状況を示す指標変数
- ira - IRA (Individual Retirement Accounts) と呼ばれる個人年金への参画を表す指標変数
- ownhome - 持ち家の有無を示す指標変数
- twoearn - 共働きか否かを示す指標変数

ここで調べようとしているのは e401k の assets に対する処置効果が incomecat, age, educ, pension, married, ira, ownhome, twoearn の組合せ (x と表記する) によって変動するか否かという点です。すなわち e401k の assets に対する効果を変数群 x に条件付けた形で — 別の言い方をすると x の関数として — 推定したいわけです。asset(1) を 401(k) プログラム資格ありとしたときの潜在的アウトカム、asset(0) を該当資格なしとしたときの潜在的アウトカムとするなら、推定したいものは

$$E\{\text{asset}(1) - \text{asset}(0) \mid \mathbf{x}\}$$

と表現することができます。

この場合、x は個人の特徴を表しているため、推定対象の CATE は IATE (individualized average treatment effects) ということになります。なお、cate の構文上、x は *catevarlist* と呼ばれます。

この用例では部分的線形モデルを対象に PO 推定法を用いて IATE 関数の推定を行うことにします。制御変数を何も含まないとしたときの部分的線形モデルは次のように表現できます。

$$\text{assets} = \text{e401k} * \tau(\mathbf{x}) + g(\mathbf{x}) + \epsilon$$

$\tau(\mathbf{x})$ は処置 e401k との交互作用を有する x の関数、 $g(\mathbf{x})$ はノンパラメトリックな nuisance 関数、 ϵ はアウトカムに対する誤差項です。一方、処置 e401k に関する処置割当てモデルは

$$\text{e401k} = f(\mathbf{x}) + u$$

のように表現されます。 $f(\mathbf{x})$ はノンパラメトリックな nuisance 関数、 u は処置に関する誤差項です。

このとき潜在的アウトカムは

$$\begin{aligned} \text{asset}(1) &= \tau(\mathbf{x}) + g(\mathbf{x}) + \epsilon \\ \text{asset}(0) &= g(\mathbf{x}) + \epsilon \end{aligned}$$

のように記述できるため、関数 $\tau(\mathbf{x})$ が IATE 関数を表すこととなります。

$$\tau(\mathbf{x}) = E\{\text{asset}(1) - \text{asset}(0) \mid \mathbf{x}\}$$

なお、 $\tau(\mathbf{x})$ の関数形について何の仮定も置いていない点に注意してください。ここでは $\tau(\mathbf{x})$ について特定の関数形を仮定するというアプローチではなく、データ自体にそれを語らせるというアプローチを取ります。すなわち Athey, Tibshirani, and Wager (2019) によって提案されている一般化ランダムフォレスト (generalized random forest) を介して、関数 $\tau(\mathbf{x})$ をノンパラメトリックな形で推定しようというわけです。この手法は因果フォレスト (causal forest) という名でも知られており、cate ではデフォルトの推定法という形で位置付けられています。

cate の実行に先立ち、条件付けのための変数群 x をグローバルマクロ catecovars として定義しておきます。グローバルマクロの用法については [P] macro, もしくは BR01 のプログラミング機能の項を参照ください。

```
. global catecovars age educ i.(incomecat pension married twoearn ira ownhome)
```

cate の推定法としてはデフォルトの po を選択します。po というのは partialing-out の略語ですが、これは他の変数の影響を排除するという意味を持ちます。アウトカム変数としては assets を、catevarlist としてはマクロ \$catecovars を、処置割当て変数としては e401k を指定します。

- Statistics ▷ Causal inference/treatment effects ▷ Continuous outcomes
 - ▷ Conditional average treatment effects と操作
- Model タブ: Estimator: Partialing out (デフォルト)
 - Dependent variable: assets
 - Covariates for the CATE model: \$catecovars
 - Treatment variable: e401k
 - Random-number seed: 12345671

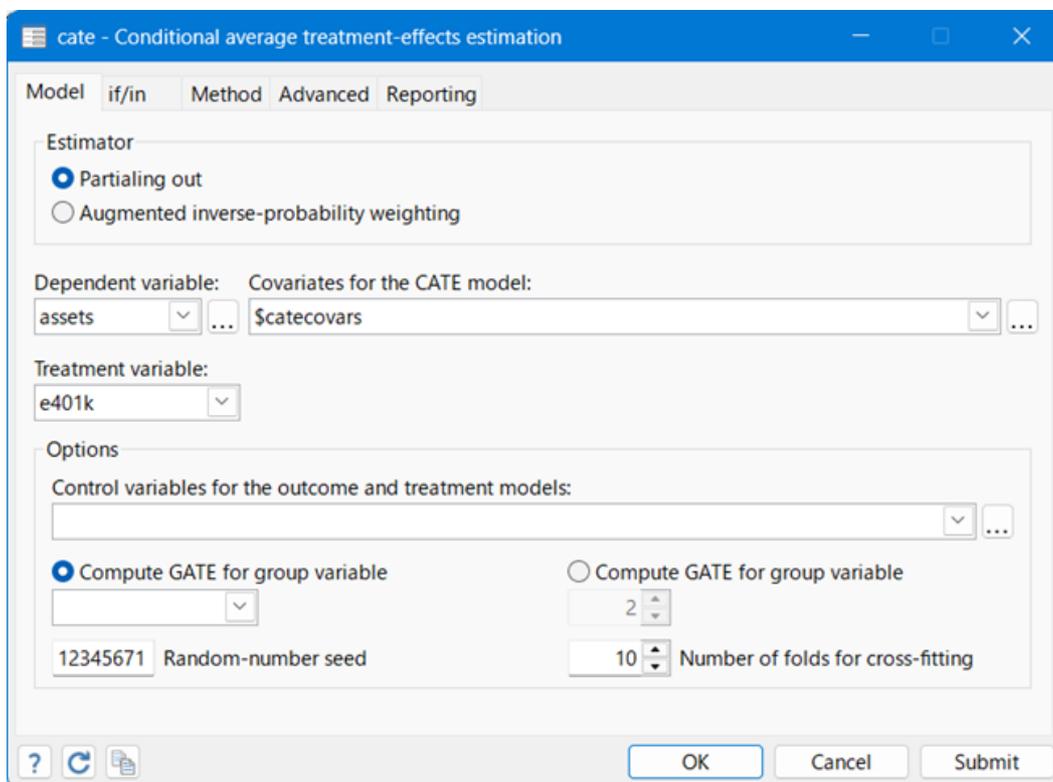


図1 cate ダイアログ - Model タブ

```

. cate po (assets $catecovars) (e401k), rseed(12345671)

Cross-fit fold 1 of 10 ...
Performing lasso for outcome assets ...
Performing lasso for treatment e401k ...

( output omitted )

Cross-fit fold 10 of 10 ...
Performing lasso for outcome assets ...
Performing lasso for treatment e401k ...

Performing random forest for IATE ...
Estimating AIPW scores ...
Estimating ATE ...

Conditional average treatment effects      Number of observations      = 9,913
Estimator:      Partialing out              Number of folds in cross-fit = 10
Outcome model:  Linear lasso                 Number of outcome controls  = 17
Treatment model: Logit lasso                 Number of treatment controls = 17
CATE model:     Random forest                Number of CATE variables    = 17

```

	assets	Robust Coefficient	std. err.	z	P> z	[95% conf. interval]	
ATE							
	e401k (Eligible vs Not eligible)	7937.182	1153.017	6.88	0.000	5677.309	10197.05
POmean							
	e401k Not eligible	14016.38	833.4423	16.82	0.000	12382.87	15649.9

反復ログは `assets` についてのアウトカムモデルと `e401k` についての処置モデルを対象とした cross-fitting プロセスの進行を示しています。IATE 関数 $\tau(x)$ を推定する場合、PO 推定法は nuisance (邪魔な) 関数である $g(x)$ と $f(x)$ の影響を排除する必要があります。そのためには x に条件付けられた形のアウトカム変数、処置変数の期待値の推定が必要となります。デフォルトの場合、アウトカム `assets` の推定には線形モデル用の lasso が、処置 `e401k` の推定にはロジットモデル用の lasso が使用されます。Nuisance 関数のフィットに伴う誤差 — lasso 使用時の変数選択誤差、及びランダムフォレスト使用時の予測誤差 — の影響を防ぐべく、`cate` は cross-fitting のプロセスを使用しているわけです。デフォルトの場合には上記反復ログに示されているように 10 重の cross-fitting が行われます。

アウトカムモデルと処置モデルの推定にはランダムフォレストやパラメトリックモデルを使用するという選択肢もありますが、デフォルトでは共に lasso が選択されることとなります (Method タブ参照)。

Cross-fitting に伴う反復ログに続けてランダムフォレストによる推定結果が出力されています*3。最初に IATE 関数が、次に AIPW (augmented inverse-probability weighting) スコアが推定されます。AIPW スコアというのは 2 重にロバストな個体レベルの処置効果推定値を意味します。これら AIPW スコアの平均値が ATE となります。推定された ATE の値からすると、母集団の全員が 401(k) の資格を有するとした場合、誰もが 401(k) の資格を持たないとしたときに比べて純金融資産が平均で \$7,937 大きくなることがわかります。また潜在的アウトカム平均の推定値からすると、誰もが 401(k) の資格を持たないとしたときの純金融資産は \$14,016 であることが示されています。

出力中に表示されている ATE に加えて cate は IATE 関数 $\tau(x)$ の推定も行っています。従ってそれを用いると各観測データごとの処置効果予測値を得ることができます。ここでは categraph histogram を用いて予測される $\tau(x)$ 関数のヒストグラムを作成してみます。

- Statistics ▷ Postestimation ▷ Diagnostic and analytic plots
 - ▷ Histogram of the IATE predictions ▷ Launch と操作
- Main タブ: デフォルト設定のまま実行

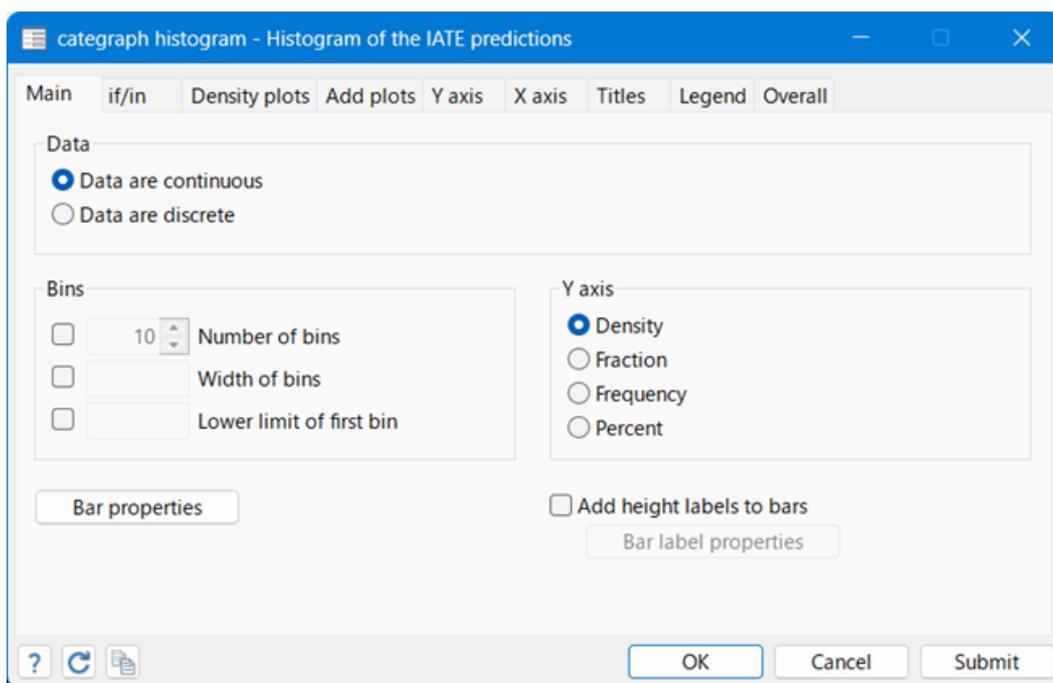
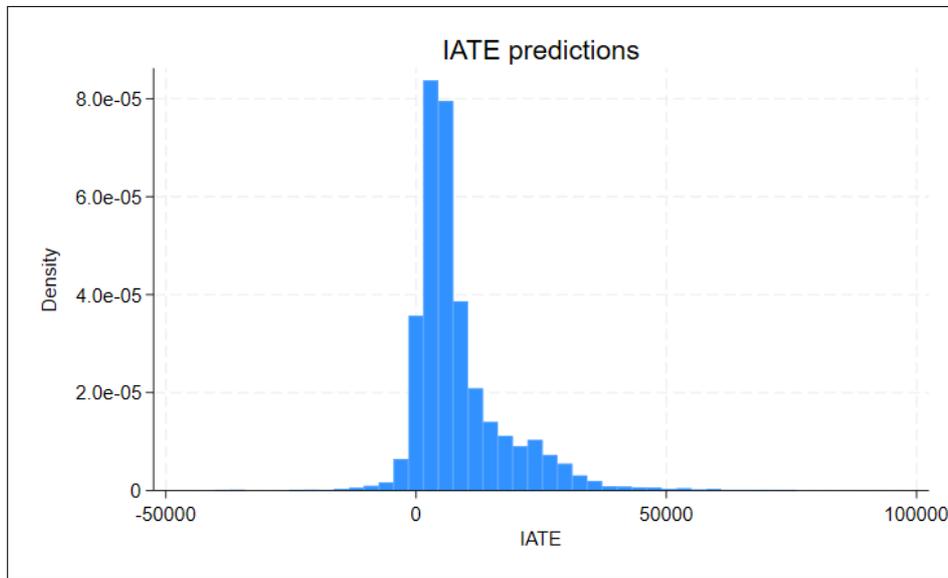


図 2 categraph histogram ダイアログ - Main タブ

```
. categraph histogram
(bin=39, start=-40204.13, width=2975.4332)
```

*3 計算には 1 分前後を要します。



ヒストグラムからは処置効果の大半がポジティブではあるが、右側の tail が太めであることがわかります。従って ATE では 401(k) 資格の効果が過少に評価されてしまうグループの存在が予想されます。

評価版では割愛しています。

▷ Example 2: 制御変数の追加

評価版では割愛しています。

▷ Example 3: グループ別 ATE の推定

評価版では割愛しています。

▷ Example 4: 連続変数に関する ATE の推定

評価版では割愛しています。

▷ Example 5: AIPW 推定法の使用

評価版では割愛しています。

▷ Example 6: データドリブンなグループ間仮説検定

評価版では割愛しています。

▷ Example 7: モデルの選択

評価版では割愛しています。

▷ Example 8: 処置割当てポリシーの評価

評価版では割愛しています。

補足 1 – table コマンド操作

評価版では割愛しています。

