

効果の異質性

事例の集計

川田恵介

条件付き平均差

-

$$\tau_P(X) = E_P[Y|D = 1, X] - E_P[Y|D = 0, X]$$

の予測モデル

$$g_\tau(X)$$

を推定

- どのような層において、差が大きいのか？

応用: 医療行為の”個人化”

- 医療行為における既存の統計分析は、平均的な効果に焦点を当てる
 - 有効な医療行為は、個人の体質等に依存している可能性
 - 個人に合わせた医療行為を、“根拠”を持って、行いたい

応用: マンション経営コンサルティング

- 中古マンションの”リノベ”コンサルティングを展開したい
 - 改築を行えば、どの程度市場価値が上がるのか？
 - 根拠を持って、“予測”したい

一般化された部分線形モデル

-

$$E[Y|D, X] = \tau(X) \times D + f(X)$$

- 前回は

$$\tau(X) = \tau$$

R Learner

- 大本のアイデアを示した研究者 (Robinson) にちなんで
1. Y, D の予測モデル $g_Y(X), g_D(X)$ を推定
 2. D についての予測誤差 $D - g_D(X)$ から、 Y についての予測誤差 $Y - g_Y(X)$ を予測するモデルを機械学習を用いて、推定
- 前回は、単回帰

Causal Forest

- 最も代表的な方法
- 以下を最小にするように”深い”決定木を推定

$$E[(Y - g_Y(X) - \tau(X) \times [D - g_D(X)])^2]$$

– 大量の決定木を推定し、平均を予測値とする

例 Causal Forest

```
Fit <- causal_forest(
  X = X,
  Y = Y,
  W = D,
  Y.hat = FitY$SL.predict,
  W.hat = FitD$SL.predict
)
```

例

GRF tree object

Number of training samples: 2086

Variable splits:

- (1) split_variable: Size split_value: 70
- (2) split_variable: RoomNum split_value: 1
- (4) split_variable: Tenure split_value: 42

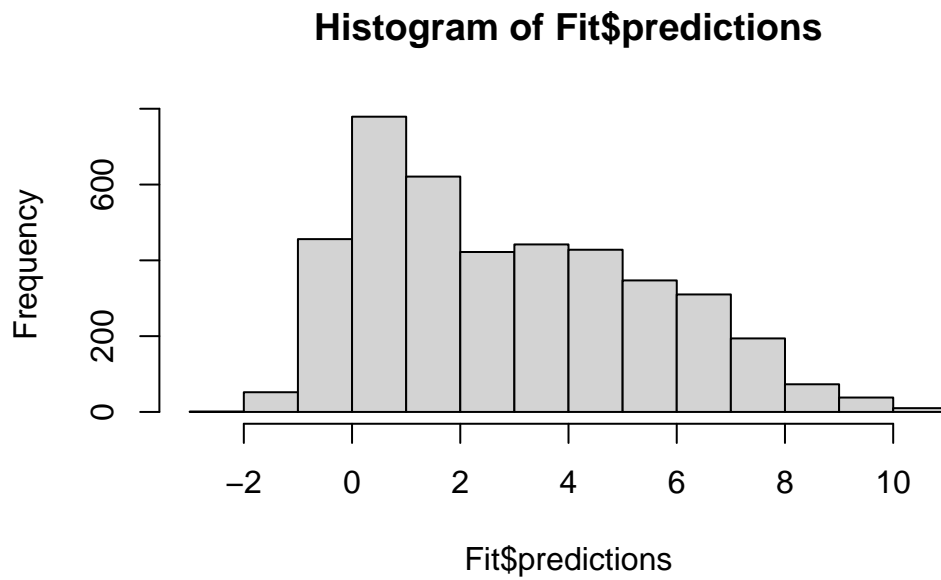
(8) split_variable: Size split_value: 45
 (14) split_variable: RoomD split_value: 0
 (24) split_variable: Size split_value: 25
 (34) split_variable: Tenure split_value: 29
 (46) split_variable: Youseki split_value: 400
 (58) split_variable: ZoneBusiness split_value: 0
 (66) split_variable: Youseki split_value: 300
 (72) split_variable: DistanceStation split_value: 5
 (78) * num_samples: 21 avg_Y: 22.9 avg_W: 0.05
 (79) split_variable: Youseki split_value: 200
 (80) * num_samples: 39 avg_Y: 20.11 avg_W: 0.26
 (81) * num_samples: 39 avg_Y: 23.46 avg_W: 0.15
 (73) * num_samples: 11 avg_Y: 23.64 avg_W: 0.45
 (67) split_variable: DistanceStation split_value: 7
 (74) * num_samples: 61 avg_Y: 23.05 avg_W: 0.03
 (75) * num_samples: 24 avg_Y: 22.17 avg_W: 0.21
 (59) split_variable: DistanceStation split_value: 3
 (68) split_variable: DistanceStation split_value: 2
 (76) * num_samples: 16 avg_Y: 25.44 avg_W: 0
 (77) * num_samples: 22 avg_Y: 24.05 avg_W: 0.09
 (69) * num_samples: 49 avg_Y: 24.04 avg_W: 0.18
 (47) * num_samples: 45 avg_Y: 12.02 avg_W: 0.24
 (35) * num_samples: 17 avg_Y: 32.47 avg_W: 0.12
 (25) split_variable: Tenure split_value: 13
 (36) split_variable: DistanceStation split_value: 4
 (48) * num_samples: 32 avg_Y: 42.16 avg_W: 0.06
 (49) * num_samples: 40 avg_Y: 39.65 avg_W: 0.08
 (37) split_variable: Kenpei split_value: 60
 (50) * num_samples: 22 avg_Y: 33.77 avg_W: 0.32
 (51) split_variable: Tenure split_value: 19
 (60) * num_samples: 29 avg_Y: 38.79 avg_W: 0.24
 (61) * num_samples: 25 avg_Y: 29.4 avg_W: 0.28
 (15) split_variable: DistanceStation split_value: 4
 (26) * num_samples: 12 avg_Y: 62.08 avg_W: 0.08
 (27) * num_samples: 19 avg_Y: 55.21 avg_W: 0.26
 (9) split_variable: Size split_value: 30
 (16) * num_samples: 14 avg_Y: 15.97 avg_W: 0.36
 (17) * num_samples: 10 avg_Y: 25.5 avg_W: 0.6
 (5) split_variable: ZoneHouse split_value: 0
 (10) split_variable: RoomNum split_value: 2

```

(18) split_variable: Tenure split_value: 16
    (28) * num_samples: 66 avg_Y: 58.02 avg_W: 0.12
(29) split_variable: ZoneBusiness split_value: 0
    (38) * num_samples: 22 avg_Y: 37.18 avg_W: 0.45
    (39) split_variable: Tenure split_value: 24
        (52) * num_samples: 16 avg_Y: 51.56 avg_W: 0.5
        (53) split_variable: Size split_value: 45
            (62) * num_samples: 25 avg_Y: 28.38 avg_W: 0.44
            (63) * num_samples: 8 avg_Y: 42.88 avg_W: 0.38
(19) split_variable: Size split_value: 65
    (30) split_variable: Tenure split_value: 10
        (40) * num_samples: 25 avg_Y: 53 avg_W: 0.08
        (41) split_variable: DistanceStation split_value: 9
            (54) split_variable: ZoneBusiness split_value: 0
                (64) split_variable: Tenure split_value: 34
                    (70) * num_samples: 15 avg_Y: 42 avg_W: 0.33
                    (71) * num_samples: 5 avg_Y: 28.4 avg_W: 0.6
                (65) * num_samples: 23 avg_Y: 42.52 avg_W: 0.65
            (55) * num_samples: 28 avg_Y: 31.61 avg_W: 0.43
        (31) * num_samples: 40 avg_Y: 52 avg_W: 0.22
(11) split_variable: Tenure split_value: 18
    (20) * num_samples: 64 avg_Y: 55 avg_W: 0.22
    (21) split_variable: After split_value: 0
        (32) split_variable: DistanceStation split_value: 6
            (42) * num_samples: 5 avg_Y: 35.4 avg_W: 0.6
            (43) split_variable: DistanceStation split_value: 9
                (56) * num_samples: 4 avg_Y: 54 avg_W: 0.75
                (57) * num_samples: 9 avg_Y: 24.11 avg_W: 0.44
        (33) split_variable: RoomNum split_value: 2
            (44) * num_samples: 19 avg_Y: 40.79 avg_W: 0.58
            (45) * num_samples: 17 avg_Y: 40.86 avg_W: 0.41
(3) split_variable: Youseki split_value: 300
(6) split_variable: Size split_value: 75
(12) split_variable: ZoneFactory split_value: 0
    (22) * num_samples: 21 avg_Y: 56.96 avg_W: 0.29
    (23) * num_samples: 8 avg_Y: 54.38 avg_W: 0.38
(13) * num_samples: 32 avg_Y: 77.03 avg_W: 0.12
(7) * num_samples: 44 avg_Y: 110.91 avg_W: 0.3

```

例



例

```
predict(  
  Fit,  
  X[1,],  
  estimate.variance = TRUE  
)  
  
predictions variance.estimates  
1      4.578813      0.9026616
```

例

```
predict(  
  Fit,  
  X[10,],  
  estimate.variance = TRUE  
)
```

```

predictions variance.estimates
1      1.244246      1.298693

```

中間まとめ

- X から $E_P[Y|D=1, X] - E_P[Y|D=0, X]$ を予測することは、(一応) 可能
 - Random Forest であれば、 X の数が少なければ、信頼区間も計算できる
- ただし多くの応用で、精度は高くない

Best Linear Projection

- 線形近似であれば、より高い精度で推定できる
-

$$\tau(X) \sim \beta_0 + \dots + \beta_L X_L$$

例: Best Linear Predictor

Best linear projection of the conditional average treatment effect.
 Confidence intervals are cluster- and heteroskedasticity-robust (HC3):

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-8.12001899	16.88852052	-0.4808	0.6306833
Tenure	0.18776568	0.05535516	3.3920	0.0007003 ***
DistanceStation	0.01700855	0.16167882	0.1052	0.9162225
Size	0.06706808	0.15020357	0.4465	0.6552489
Youseki	0.00016679	0.00624013	0.0267	0.9786772
Kenpei	0.06950402	0.25514394	0.2724	0.7853196
RoomNum	-1.03064318	1.82439243	-0.5649	0.5721560
ZoneHouse	-2.21206287	1.71596586	-1.2891	0.1974330
ZoneBusiness	-0.69148008	5.22979253	-0.1322	0.8948172
RoomL	0.09706019	2.48781843	0.0390	0.9688810
RoomD	2.78727653	2.18615683	1.2750	0.2023926
StructureSRC	-0.04674066	1.52815268	-0.0306	0.9756009
After	-0.05751743	1.16172250	-0.0495	0.9605149

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

定数項の解釈

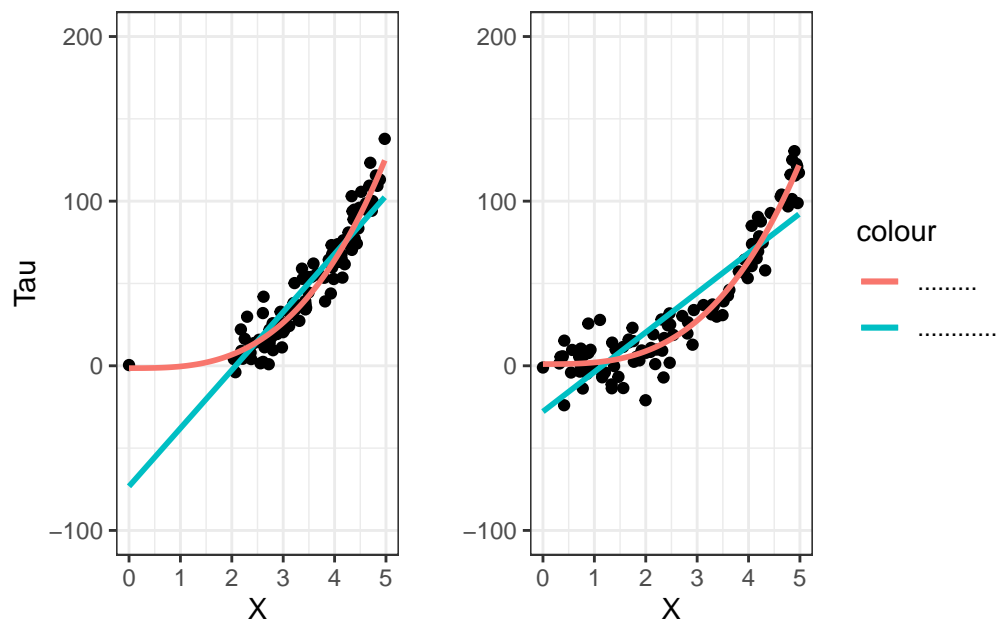
- 線形近似モデルの定数項は、通常解釈困難

•

$$E[\tau|X] \simeq \underbrace{\beta_0}_{?} + \dots + \beta_L X_L$$

- まっすぐな解釈は、 $E[\tau|X=0] = \beta_0$
 - 多くのデータで、 $X=0$ の付近には事例がない
 - 近似の際に”無視”されている

例



中心化

- 元の変数を平均 0 に変換

•

$$Z = X - E[X]$$

- 発展: 標準化

•

$$Z = \frac{X - E[X]}{SD[X]}$$

例: Best Linear Predictor

Best linear projection of the conditional average treatment effect.

Confidence intervals are cluster- and heteroskedasticity-robust (HC3):

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	2.35003725	0.62820167	3.7409	0.0001858	***
Tenure	0.18776568	0.05535516	3.3920	0.0007003	***
DistanceStation	0.01700855	0.16167882	0.1052	0.9162225	
Size	0.06706808	0.15020357	0.4465	0.6552489	
Youseki	0.00016679	0.00624013	0.0267	0.9786772	
Kenpei	0.06950402	0.25514394	0.2724	0.7853196	
RoomNum	-1.03064318	1.82439243	-0.5649	0.5721560	
ZoneHouse	-2.21206287	1.71596586	-1.2891	0.1974330	
ZoneBusiness	-0.69148008	5.22979253	-0.1322	0.8948172	
RoomL	0.09706019	2.48781843	0.0390	0.9688810	
RoomD	2.78727653	2.18615683	1.2750	0.2023926	
StructureSRC	-0.04674066	1.52815268	-0.0306	0.9756009	
After	-0.05751743	1.16172250	-0.0495	0.9605149	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1					

補論: 信頼区間の計算

- 多くの関数において、推定誤差 (Std.Error) は報告される
- 95% 信頼区間はざっくり
-

$$\text{推定値} \pm 2 * \text{Std.Error}$$

まとめ

- 機械学習は、条件付き平均差をそのまま推定するための活路
 - ただ依然として困難が多い

- 母集団における平均差を推測するツールの方が、現状より確立されている
 - 典型的な方法は、線形近似モデルの推定
 - 定式化に注意