

2022年3月26日
第24回春の合宿セミナー（日本行動計量学会）
（統計的因果推論入門）

講義4

重回帰分析による 交絡因子の統制の意味

長崎大学 情報データ科学部 准教授

高橋 将宜

博士（理工学）

m-takahashi@nagasaki-u.ac.jp

概要

- 回帰分析の復習
 - 講義1の復習
 - 三変数のバレンティン・ベン図
 - 三変数の重回帰モデル
 - 共分散分析(再考)
 - 重回帰モデルによる分析
- } 教科書Ch.5
- } 教科書Ch.6



回帰分析の復習

母集団と標本における回帰式

母集団

$$Y_i = \alpha + \beta X_i + \varepsilon_i$$

$$E[Y_i|X_i] = \alpha + \beta X_i$$

$$\varepsilon_i = Y_i - E[Y_i|X_i]$$

標本

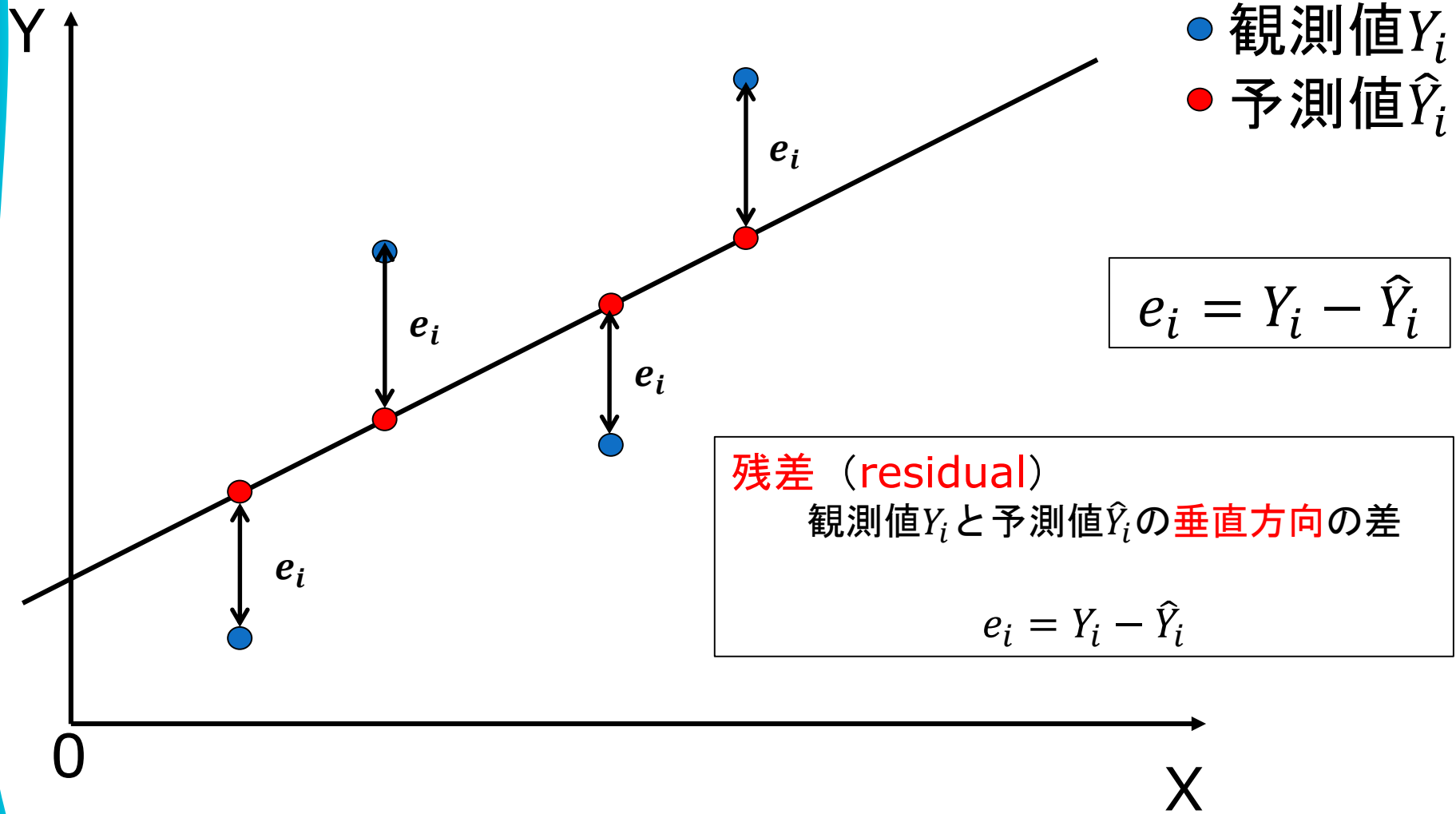
$$Y_i = \hat{\alpha} + \hat{\beta} X_i + e_i$$

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i$$

$$e_i = Y_i - \hat{Y}_i$$

誤差項 ε は偶然によって
生じる確率変数と考える。

残差 e_i



通常の最小二乗法 (OLS: Ordinary Least Squares)

- 残差平方和を最小化する切片と傾きを用いて線を引く方法

$$\hat{\beta} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$

lmの構文

lm(変数1~変数2)

数式の左辺に来る変数が左側、右辺に来る変数が右側

回帰係数に関する仮説

□ 帰無仮説

$$H_0: \beta = 0$$

□ 対立仮説

$$H_A: \beta \neq 0$$

- 回帰式 $Y_i = \alpha + \beta X_i + \varepsilon_i$ において、 $\beta = 0$ は、 X_i から Y_i への直線的な影響がないことを意味する。

- つまり、 X_i が増加しても、それにつれて Y_i が増減する傾向がない。

決定係数 R^2 の定義

□ R^2 は、回帰モデルによって説明できる Y の変動の割合を表す。

Y_i : 観測値
 \bar{Y} : 平均値
 \hat{Y}_i : 予測値

$$R^2 = 1 - \frac{USS}{TSS}$$
$$0 \leq R^2 \leq 1$$

TSS: (総和変動)

Total Sum of Squares

$$\sum (Y_i - \bar{Y})^2$$

ESS: (回帰モデルで説明できる変動)

Explained Sum of Squares

$$\sum (\hat{Y}_i - \bar{Y})^2$$

USS: (回帰モデルで説明できない変動)

Unexplained Sum of Squares

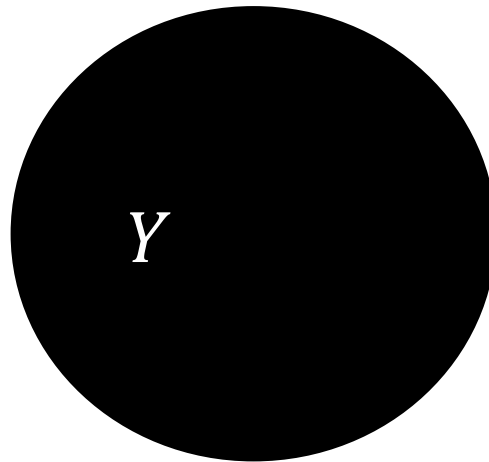
$$\sum (Y_i - \hat{Y}_i)^2$$

全変動（TSS）のバレンティン・ベン図

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_{1i}$$

$$TSS = \sum (Y_i - \bar{Y})^2$$

分散の分子と
同じ

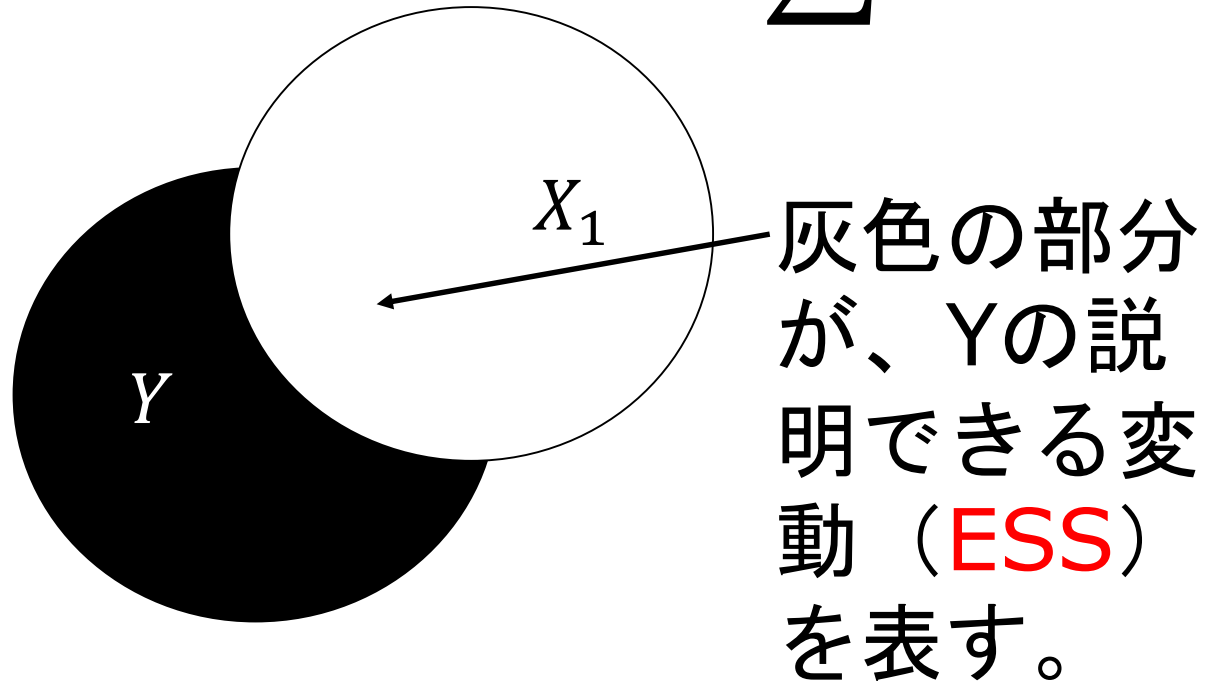


黒丸：Yの全変動（TSS）

説明できる変動（ESS）のバレンティン・ベン図

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_{1i}$$

$$ESS = \sum (\hat{Y}_i - \bar{Y})^2$$

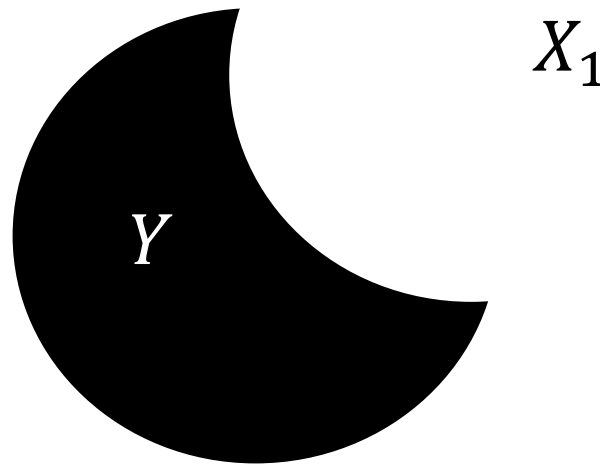


説明できない変動（USS）のバレンティン・ベン図

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_{1i}$$

$$USS = \sum (Y_i - \hat{Y}_i)^2$$

残差 : $e_i = Y_i - \hat{Y}_i$

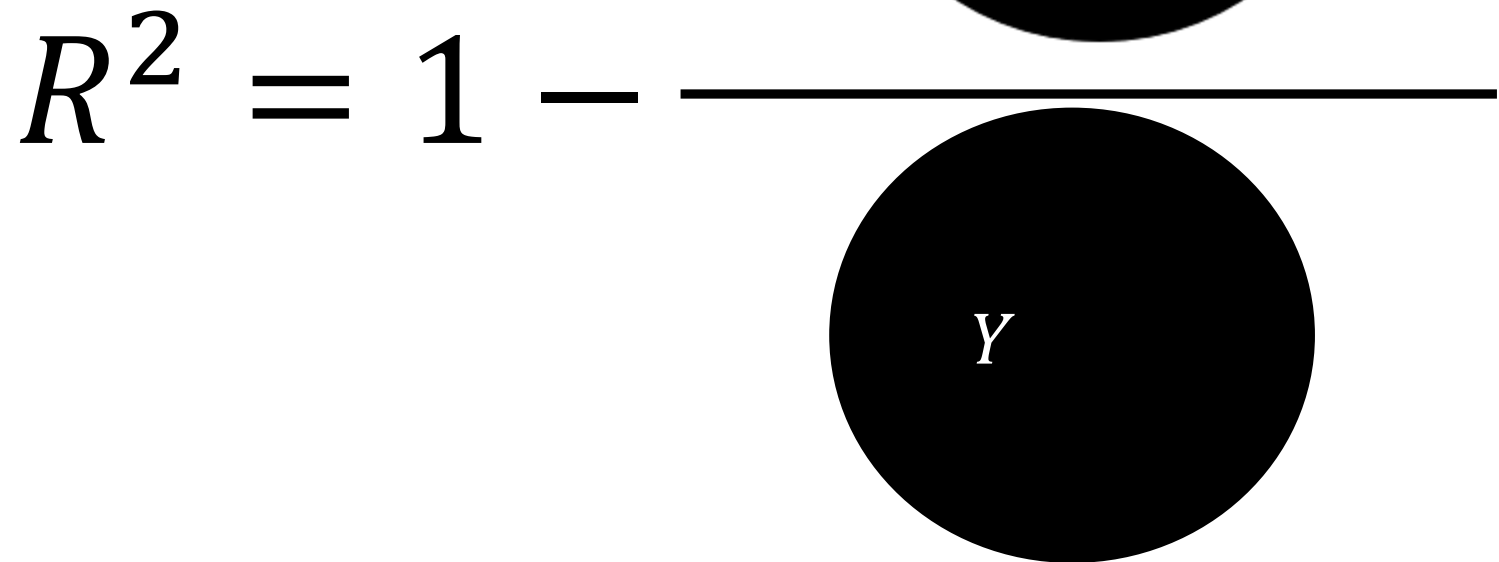


残った黒い部分が、
Yの説明できない変
動（**USS**）を表す。

説明できる変動の割合のベン図

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_{1i}$$

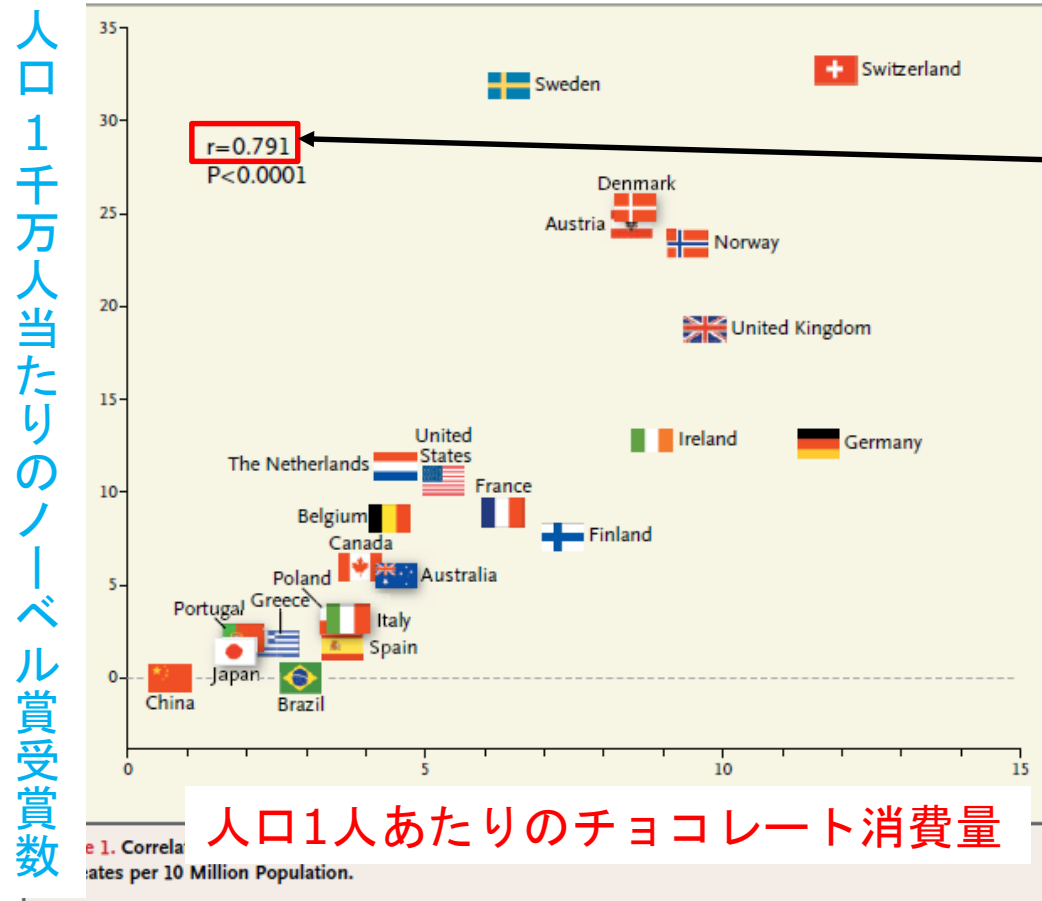
$$R^2 = 1 - \frac{USS}{TSS}$$



講義1の復習

散布図と相関係数

出典：Messerli, F. H. (2012) "Chocolate Consumption, Cognitive Function, and Nobel Laureates," The New England Journal of Medicine, 367 (16), pp.1562-1564.



$r = 0.791$

単回帰モデル

```
> summary(modell<-lm(Nobel~Choco))
```

```
Call:
```

```
lm(formula = Nobel ~ Choco)
```

```
Residuals:
```

	Min	1Q	Median	3Q	Max
	-10.1603	-4.3915	-0.7202	2.5621	16.3355

```
Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-3.4217	3.2274	-1.060	0.301096
Choco	2.7044	0.5985	4.519	0.000188 ***

```
---
```

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

```
Residual standard error: 6.65 on 21 degrees of freedom
```

```
Multiple R-squared:  0.493,    Adjusted R-squared:  0.4689
```

```
F-statistic: 20.42 on 1 and 21 DF,  p-value: 0.000188
```

```
> confint(modell)
```

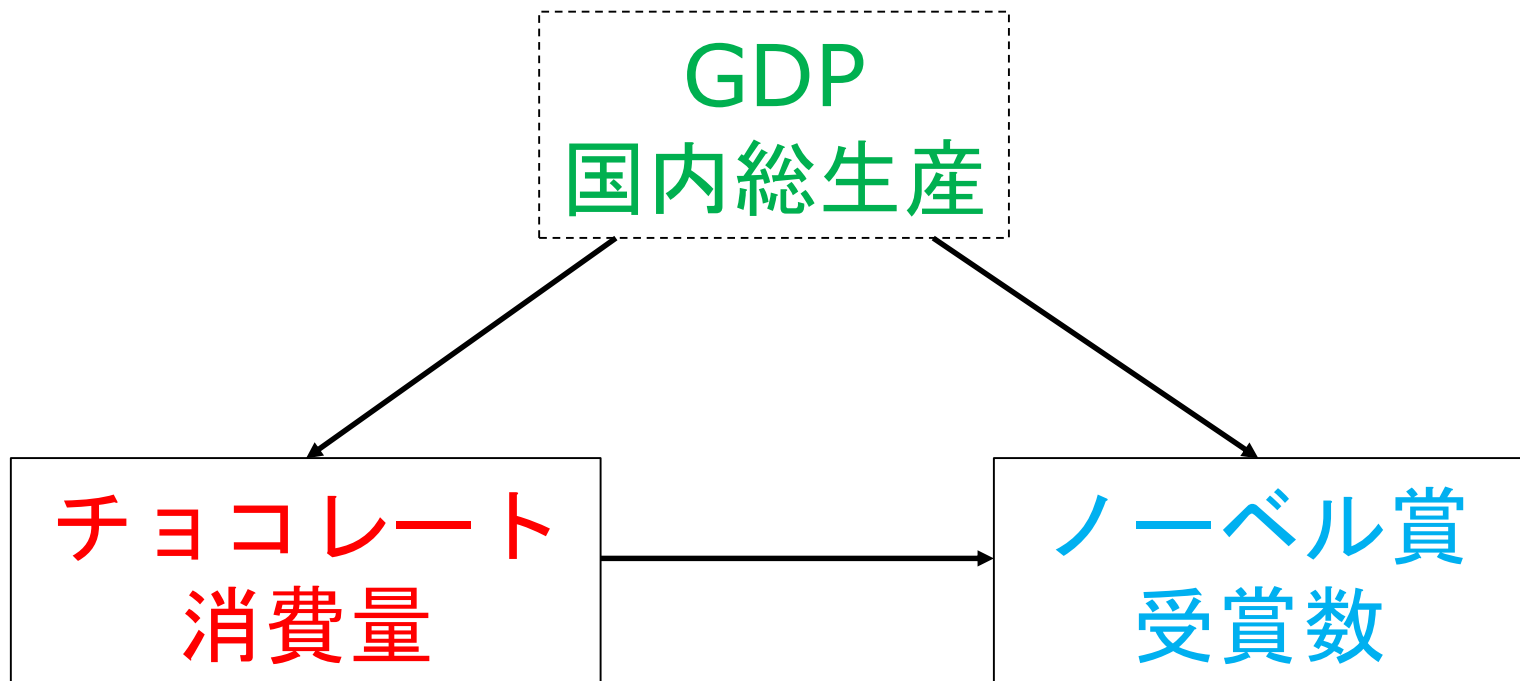
	2.5 %	97.5 %
(Intercept)	-10.133320	3.290018
Choco	1.459837	3.949020

疑似相関

- 相関係数の絶対値が1に近いにもかかわらず、実際には2つの現象に直接的な関係がないこと

第3の変数

- ノーベル賞受賞数とチョコレート消費量の背後に共通の原因であるGDP（国内総生産）があると考えられる.

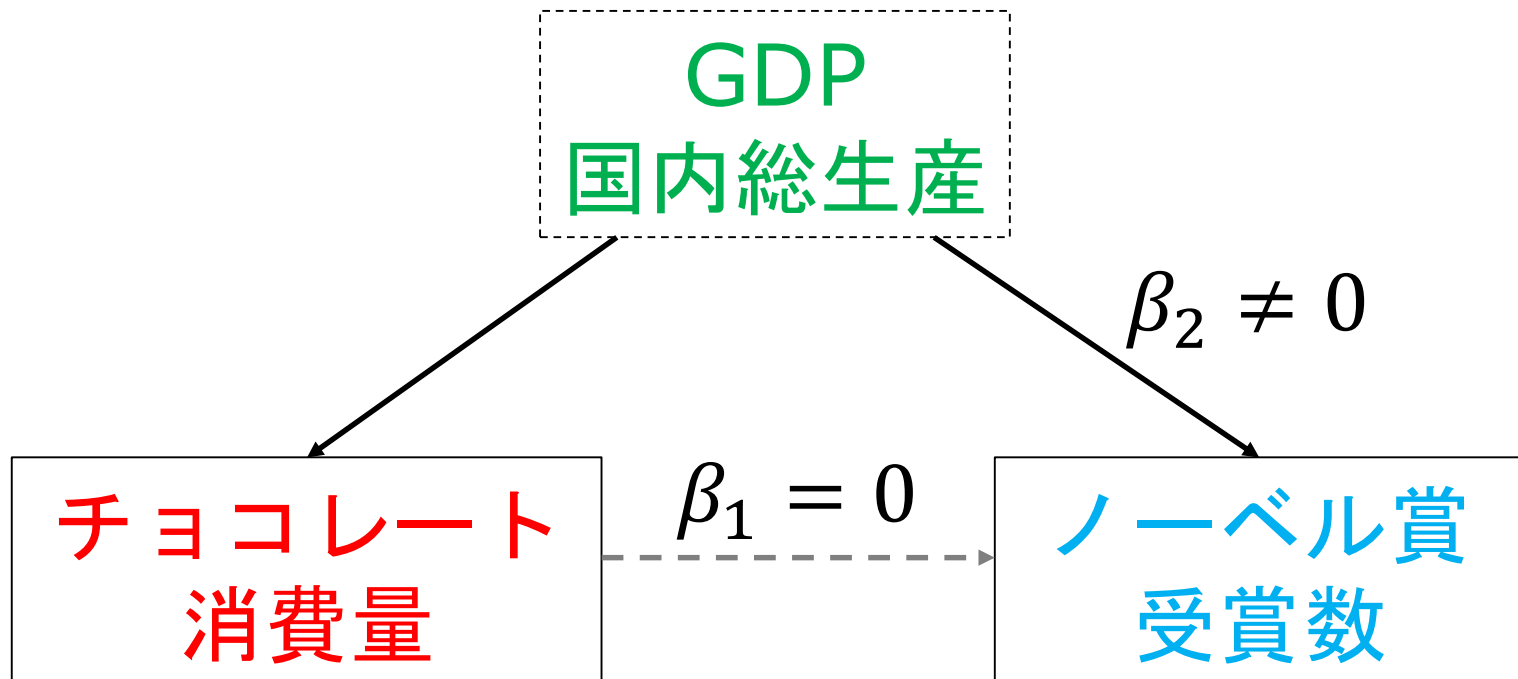


重回帰モデル

$$\hat{Y}_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i}$$

ノーベル賞受賞数_i

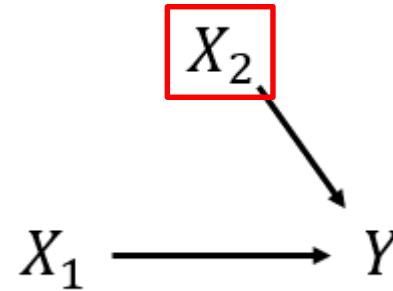
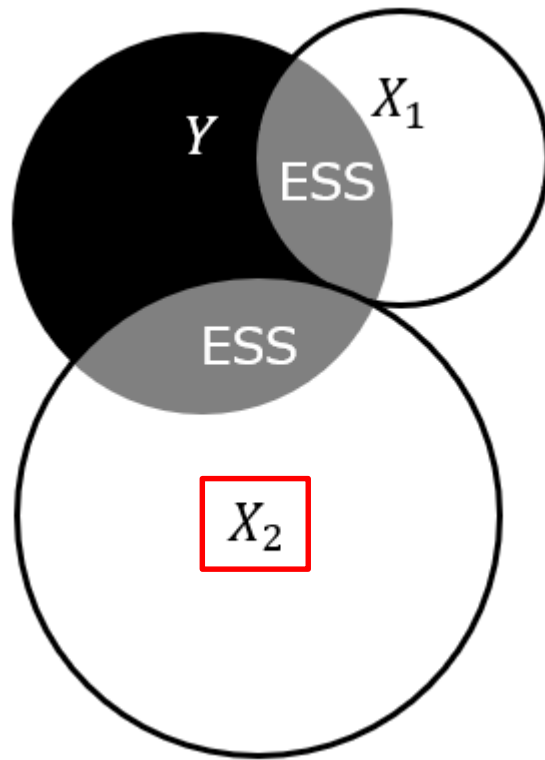
$$= \beta_0 + \beta_1 \text{チョコレート消費量}_i + \beta_2 \text{GDP}_i$$



三変数のバレンティン・ベン図

X_2 は共変量だが、
交絡因子ではない

説明変数の間に相関がない場合（1）



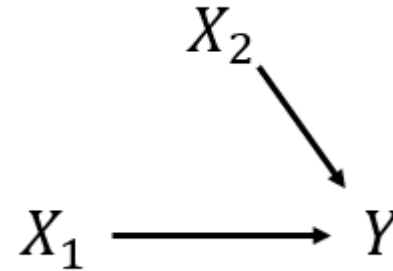
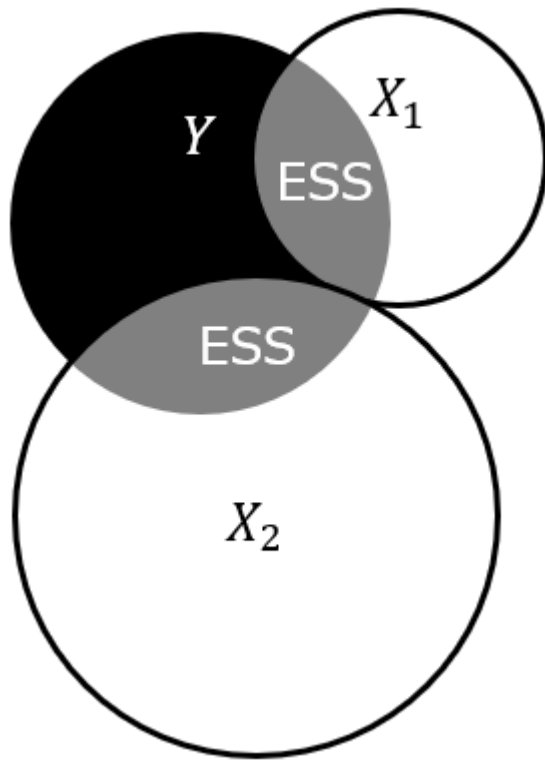
- これまでの2つの変数 Y と変数 X_1 に加えて、3つ目の変数 X_2 があるでしょう。
- 変数 Y と変数 X_2 の重なっている部分がESSである。
- 変数 X_1 と変数 X_2 には重なりがないことから、相関は0である。
- 方向付き非巡回グラフ（DAG）で表すと、 X_1 と X_2 の間に矢印がない状態である。

ESS: (回帰モデルで説明できる変動)
Explained Sum of Squares

$$\sum (\hat{Y}_i - \bar{Y})^2$$

X_2 は共変量だが、
交絡因子ではない

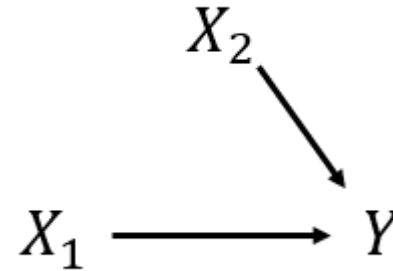
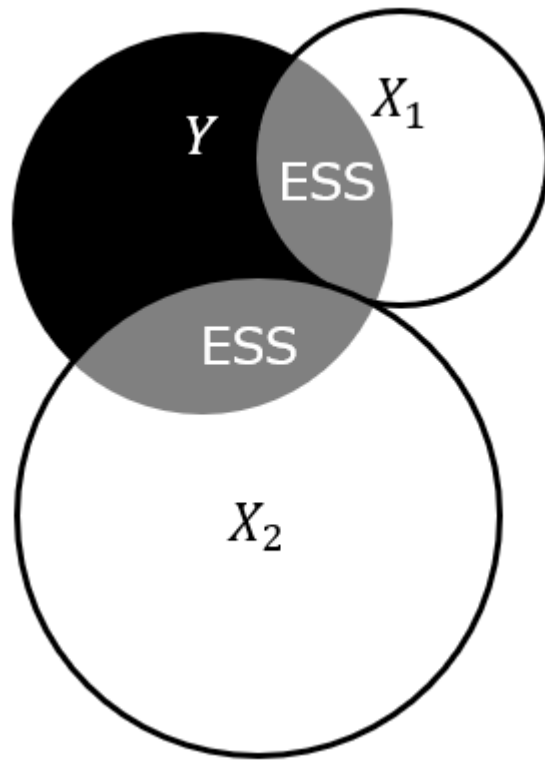
説明変数の間に相関がない場合 (2)



- $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i}$ における $\hat{\beta}_1$ は、 X_2 からの影響を何も受けていないことから、 $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i}$ における $\hat{\beta}_1$ と一致する。
- X_2 をモデルに含めなくても、 $\hat{\beta}_1$ は β_1 の不偏推定量である。
- X_2 は交絡因子ではない。

X_2 は共変量だが、
交絡因子ではない

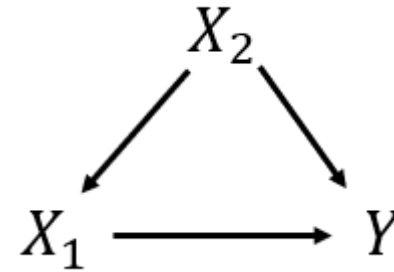
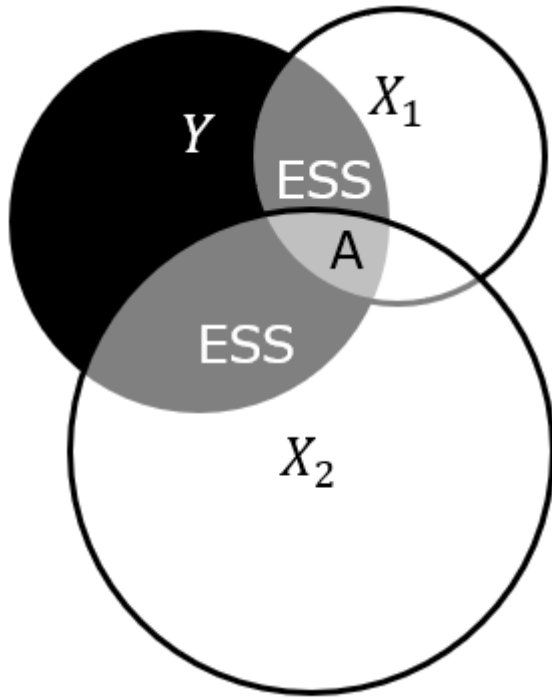
説明変数の間に相関がない場合 (3)



- 単回帰モデルと比べて、重回帰モデルの方が、**ESSの部分が大きくなった**ため、 Y の変動を説明する力が向上している。
- このような変数は、回帰係数の不偏性には影響を与えないが、**標準誤差の大きさに影響を及ぼす**ため、データセット内に含まれていて利用できる状況ならば、モデルに取り込むことが望ましいが、モデルに入れなくても問題ではない。

X_2 は共変量であり,
交絡因子である.

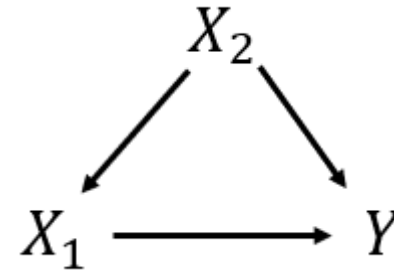
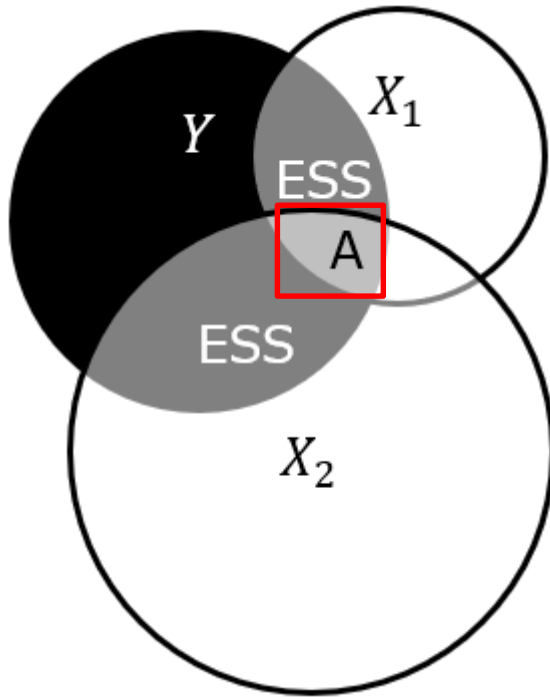
説明変数の間に相関がある場合 (1)



- 変数 X_1 と変数 X_2 に重なりがあるとしよう.
- DAGで表すと, X_2 から X_1 へ矢印がある状態である.

X_2 は共変量であり,
交絡因子である.

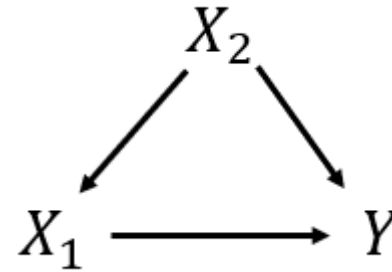
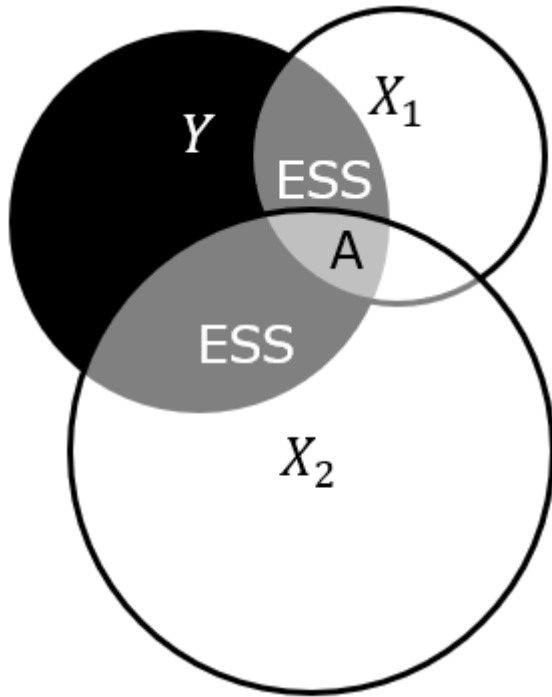
説明変数の間に相関がある場合 (2)



- X_1 から Y への因果効果を $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i}$ の単回帰モデルの $\hat{\beta}_1$ で測ろうとすると, その中には, A で表されている部分が含まれている.
- この A の部分は, X_2 からの間接的な効果であり, これが交絡である.
- X_1 と X_2 が重なっている部分 A を取り除いて分析する必要がある.

三変数の重回帰モデル

説明変数の間に相関がある場合の続き



- $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \varepsilon_i$ における β_1 の偏りのない推定を行うことを考える.
- A の部分を取り除いて、 X_1 から Y への純粋な効果を測る方法を考察する.

三変数の重回帰モデル

$$\square Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \varepsilon_i$$

$$\square \hat{\beta}_1 = \frac{\sum (X_{1i} - \hat{X}_{1i})(Y_i - \bar{Y})}{\sum (X_{1i} - \hat{X}_{1i})^2}$$

$$\square \hat{\beta}_2 = \frac{\sum (X_{2i} - \hat{X}_{2i})(Y_i - \bar{Y})}{\sum (X_{2i} - \hat{X}_{2i})^2}$$

$$\square \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}_1 - \hat{\beta}_2 \bar{X}_2$$

- β_1 と β_2 は、特に偏回帰係数（partial regression slope）と呼ばれる。
- ここでは、 $\hat{\beta}_1$ を使って、具体的にどのようにして交絡を取り除くことができるのかを確認する。
- この考え方は、統計的因果推論において極めて重要である。

単回帰モデルと重回帰モデルの違い

□ 単回帰モデル

$$\square \hat{Y}_i = \hat{\beta}_0 + \hat{\beta} X_{1i}$$

$$\square \hat{\beta} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$$

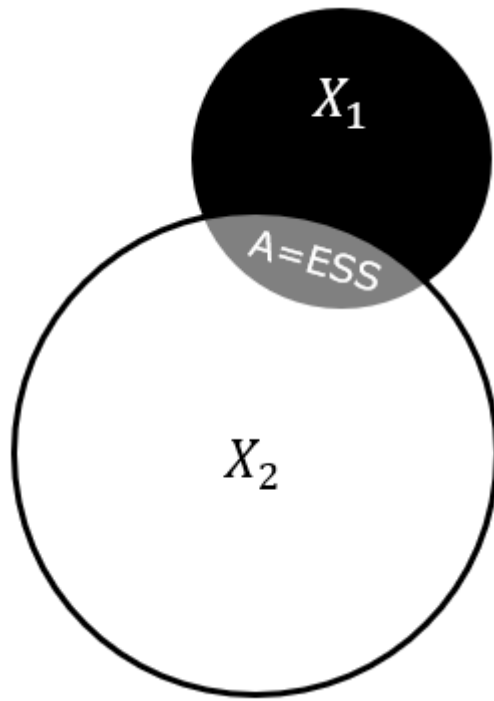
□ 重回帰モデル

$$\square \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i}$$

$$\square \hat{\beta}_1 = \frac{\sum (X_{1i} - \hat{X}_{1i})(Y_i - \bar{Y})}{\sum (X_{1i} - \hat{X}_{1i})^2}$$

- 重回帰モデルの $\hat{\beta}_1$ は、一見すると単回帰モデルの $\hat{\beta}$ とよく似ている。
- $\hat{\beta}$ は X の平均値からの偏差を計算している。
- $\hat{\beta}_1$ は X_1 の残差を計算している。
- 重回帰モデルでは、 $\hat{X}_{1i} = a_1 + b_1 X_{2i}$ である点に注意しよう。
- すなわち、三変数の重回帰モデルは、**二段階で分析**を行っている。
- まず、単回帰モデル $\hat{X}_{1i} = a_1 + b_1 X_{2i}$ を作り、 X_{2i} から予測値 \hat{X}_{1i} を計算する。
- 次に、残差 $X_{1i} - \hat{X}_{1i}$ を計算して、この残差部分から Y への回帰を行う。

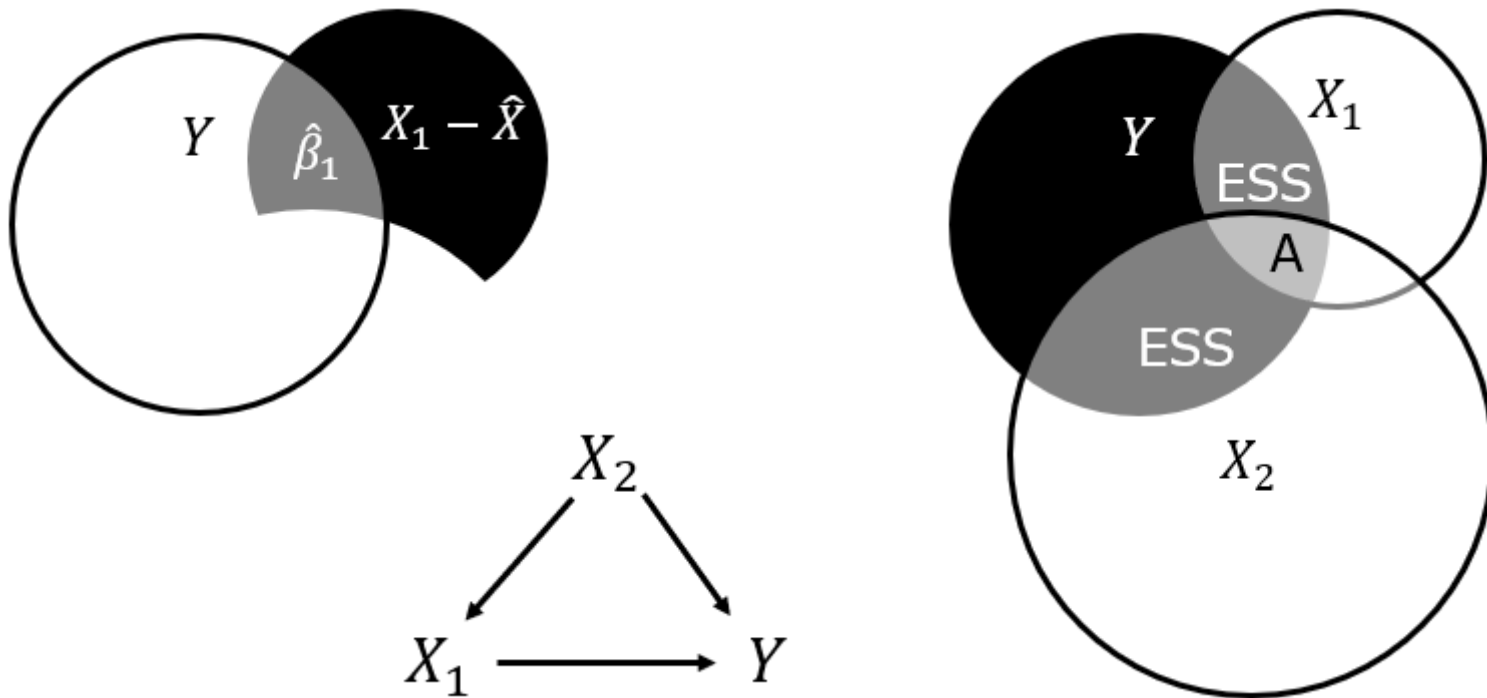
重回帰モデル：第1段階目



- まず，図から一時的に変数 Y を取り除く．
- すると，変数 X_1 と変数 X_2 の単回帰モデル $\hat{X}_{1i} = a_1 + b_1 X_{2i}$ と同じ状況である．
- A の部分は，ESSである．
- そして，図の黒い部分がUSSである．
- このとき，USSは以下の式である．

$$USS = \sum_{i=1}^n (X_{1i} - \hat{X}_{1i})^2$$

重回帰モデル：第2段階目

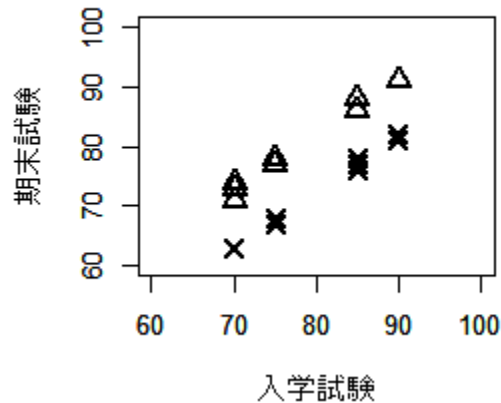


図に変数 Y を戻すと，一段階目の手順により，変数 X_2 からの交絡 A を取り除くことができる様子が分かる．

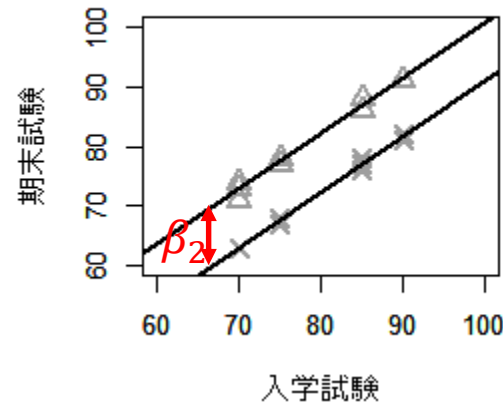
共分散分析(再考)

復習：共分散分析 (ANCOVA)

C. 散布図(群ごと)



D. 回帰直線(群ごと)



$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 T_i + \varepsilon_i$$

Y_i : 結果変数 (期末試験)

$\beta_0, \beta_1, \beta_2$: 回帰の母数 (パラメータ)

X_i : 共変量 (入学試験)

T_i : 処置を表す二値変数 (補習授業)

ε_i : 誤差項 (error term)

共分散分析 (ANCOVA)

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 T_i + \varepsilon_i$$

- **ダミー変数**を説明変数として持つ重回帰モデルである.
 - ダミー変数とは, **0または1の二値をとる変数**のことである.

復習：処置の割付け変数

□ T_i

- 個体 i が**処置 (treatment)** に割付けられたかどうかを表す二値変数

- $T_i \in \{0, 1\}$

□ $T_i = 0$

- 個体 i が処置に割付けられていないこと
- 統制群： $T_i = 0$ となる集団

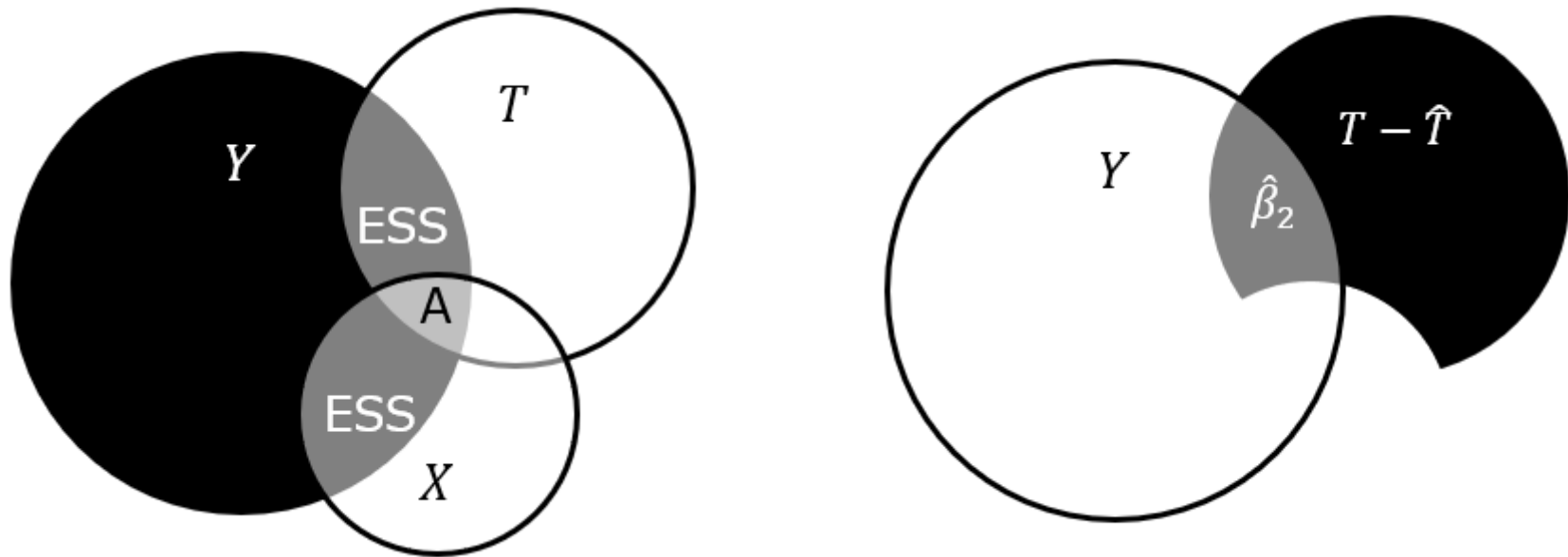
□ $T_i = 1$

- 個体 i が処置に割付けられたこと
- 処置群： $T_i = 1$ となる集団

表2.2

ID	入学試験 x1	処置 t1	期末 試験0 y0	期末 試験1 y1	潜在的 結果0 y0t	潜在的 結果1 y1t	潜在的 結果の 差 y1t - y0t
1	74	1		76	68	76	8
2	82	0	75		75	84	9
3	72	1		75	65	75	10
4	96	0	84		84	97	13
5	83	0	75		75	84	9
6	72	1		74	65	74	9
7	85	0	76		76	87	11
8	87	0	77		77	89	12
9	86	0	77		77	87	10
10	77	1		80	70	80	10
11	95	0	87		87	96	9
12	84	0	75		75	85	10
13	74	1		77	67	77	10
14	58	1		61	52	61	9
15	91	0	81		81	93	12
16	80	0	72		72	84	12
17	80	0	72		72	82	10
18	89	0	70		70	89	19
19	88	0	70		70	90	20
20	86	0	78		78	87	9

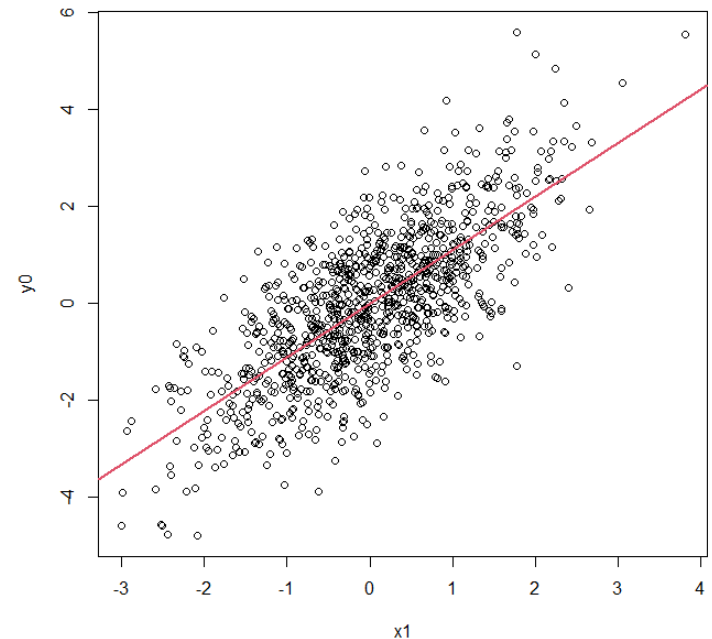
二段階の推定により交絡を取り除く



- 共分散分析のメカニズムも、これまでと同様に考えることができる.
- β_2 は結果変数 Y の変動を説明する際に、共変量に対して、ダミー変数 T の純粋な貢献度合いを表していると考えることができる

回帰分析 (regression)

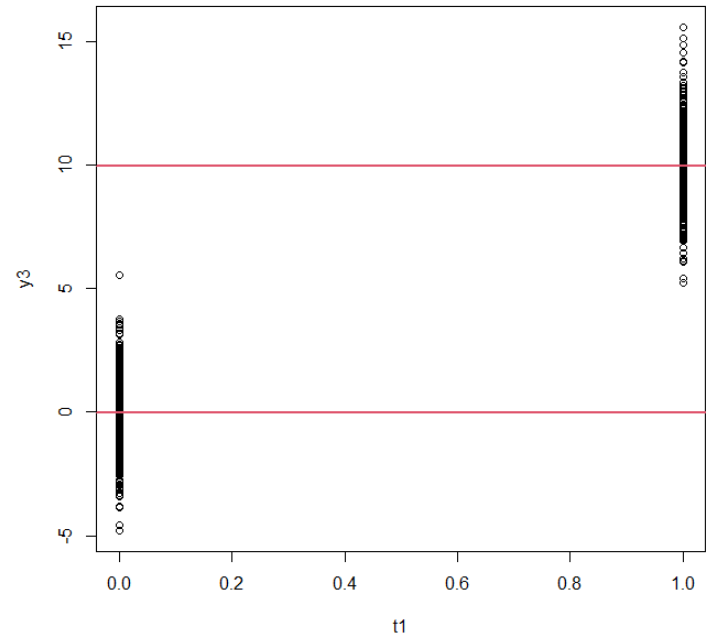
- $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$
- $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$
- すべての説明変数が**連続変数**のみから構成されている



分散分析 (ANOVA: analysis of variance)

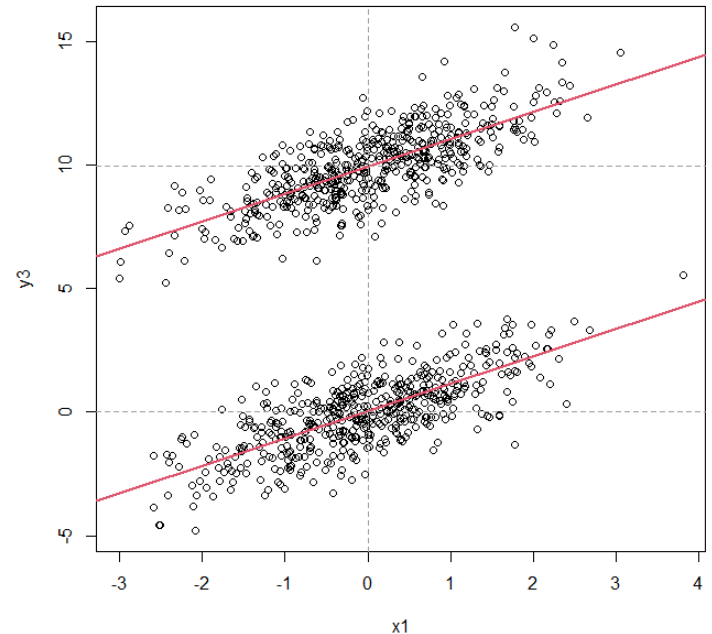
- $Y_i = \beta_0 + \beta_2 T_i + \varepsilon_i$
- $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_2 T_i$

- すべての説明変数が**ダミー変数**のみから構成されている回帰分析の特殊版である.
 - T_i が二値の場合, t検定, 分散分析, 回帰分析はすべて同じである.
 - 説明変数にダミー変数のみを用いた回帰分析は, 2標本t検定と同じである.



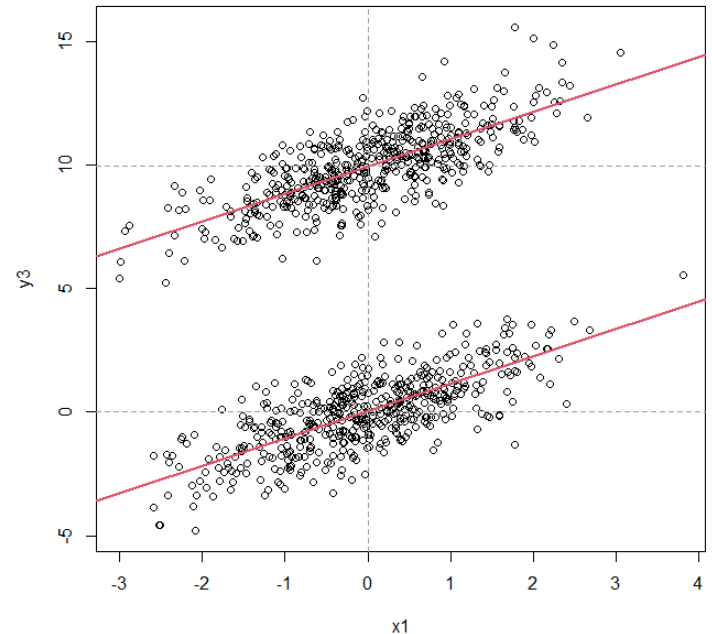
共分散分析 (ANCOVA: analysis of covariance)

- $Y_i = \beta_0 + \beta_1 X_i + \beta_2 T_i + \varepsilon_i$
- $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{\beta}_2 T_i$
- 回帰分析と分散分析の長所を併せ持つように工夫されたものである.



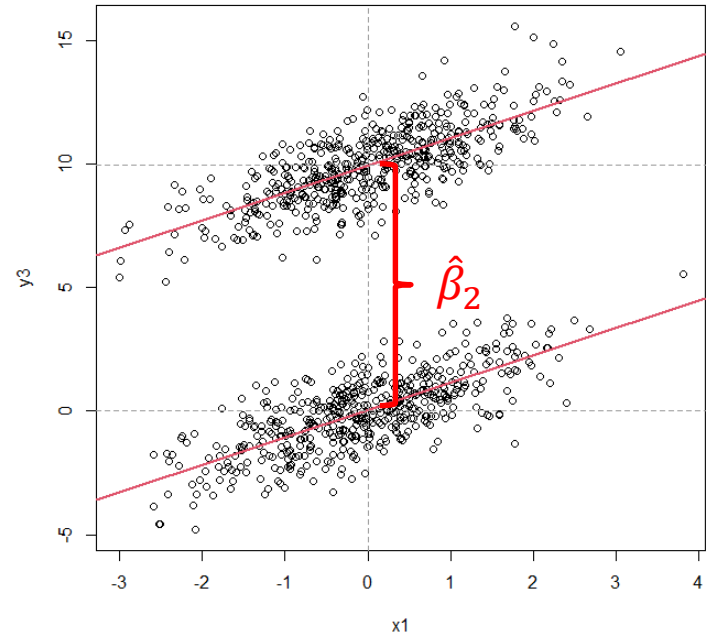
$T_i = 0$ のとき

- $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{\beta}_2 T_i$
- $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$
- 統制群における Y -切片は,
 $\hat{\beta}_0$ である.
- 統制群では, X_i の値が0の
とき, Y_i の値は $\hat{\beta}_0$ である.



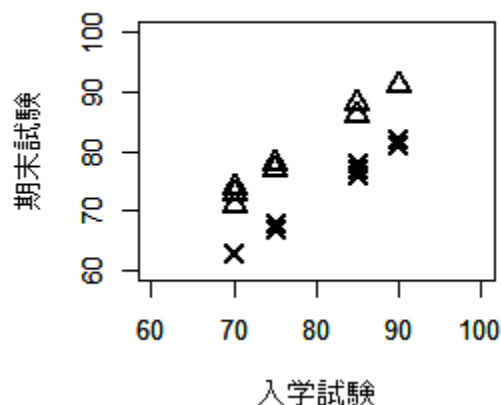
$T_i = 1$ のとき

- $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{\beta}_2 T_i$
- $\hat{Y}_i = (\hat{\beta}_0 + \hat{\beta}_2) + \hat{\beta}_1 X_i$
- 処置群における Y -切片は,
 $\hat{\beta}_0 + \hat{\beta}_2$ である.
- 処置群では, X_i の値が0の
とき, Y_i の値は $\hat{\beta}_0 + \hat{\beta}_2$ であ
る. したがって, $\hat{\beta}_2$ は2つ
の集団における切片の変化
を表している

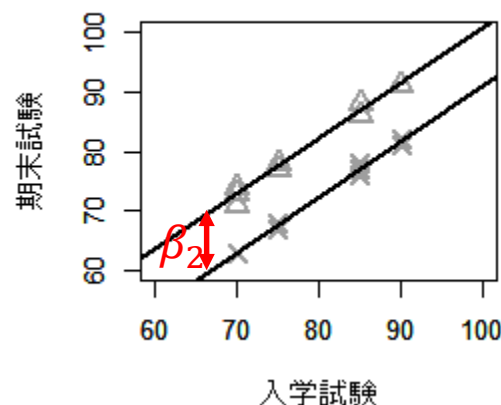


復習：共分散分析 (ANCOVA)

C. 散布図(群ごと)



D. 回帰直線(群ごと)



$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 T_i + \varepsilon_i$$

Y_i : 結果変数 (期末試験)

$\beta_0, \beta_1, \beta_2$: 回帰の母数 (パラメータ)

X_i : 共変量 (入学試験)

T_i : 処置を表す二値変数 (補習授業)

ε_i : 誤差項 (error term)

重回帰モデルによる分析

データの読み込み（教科書pp.77-78）

- ❑ `data06 <- read.csv(file.choose())`
- ❑ `attach(data06)`
- ❑ `summary(data06)`

y1 : ノーベル賞受賞数
x1 : チョコレート消費量
x2 : 国内総生産（GDP）

```
> summary(data06)
```

country	y1	x1	x2
Length:23	Min. : 0.000	Min. : 0.10	Min. : 8.755
Class :character	1st Qu.: 1.961	1st Qu.: 4.00	1st Qu.: 31.453
Mode :character	Median : 8.644	Median : 4.90	Median : 46.232
	Mean : 9.748	Mean : 4.87	Mean : 44.704
	3rd Qu.: 13.857	3rd Qu.: 6.20	3rd Qu.: 53.908
	Max. : 30.763	Max. : 8.80	Max. : 85.135

単回帰モデル（教科書p.80）

- `model1 <- lm(y1 ~ x1)`
- `summary(model1)`
- `confint(model1, level=0.95)`

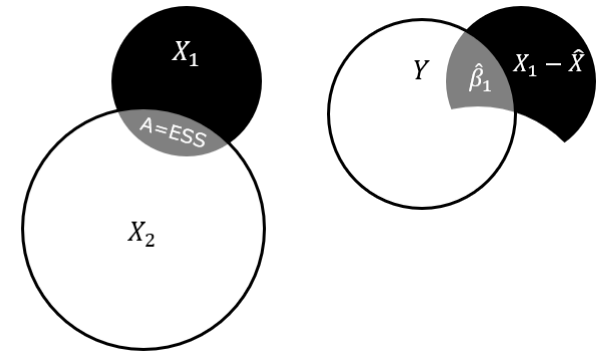
```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  -3.4217      3.2274  -1.060 0.301096
x1             2.7044      0.5985   4.519 0.000188 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.65 on 21 degrees of freedom
Multiple R-squared:  0.493,    Adjusted R-squared:  0.4689
F-statistic: 20.42 on 1 and 21 DF,  p-value: 0.000188

> confint(model1, level=0.95)
              2.5 %    97.5 %
(Intercept) -10.13332  3.290018
x1           1.459837  3.949020
```

二段階推定による重回帰モデル（教科書p.87）

```
□ model2a <- lm(x1 ~ x2)
□ ex1 <- resid(model2a)
□ model2b <- lm(y1 ~ ex1)
□ summary(model2b)
```



Coefficients:

	Estimate
(Intercept)	9.748
ex1	1.505

重回帰モデル（教科書p.87）

- ❑ `model3<-lm(y1~x1+x2)`
- ❑ `summary(model3)`
- ❑ `confint(model3, level=0.95)`

```

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  -6.32035    3.20331  -1.973   0.0625 .
x1             1.50477    0.75619   1.990   0.0604 .
x2             0.19552    0.08531   2.292   0.0329 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.064 on 20 degrees of freedom
Multiple R-squared:  0.5985,    Adjusted R-squared:  0.5583
F-statistic: 14.9 on 2 and 20 DF,  p-value: 0.0001089

> confint(model3, level=0.95)
              2.5 %    97.5 %
(Intercept) -13.00233992  0.3616436
x1          -0.07260349  3.0821496
x2           0.01757188  0.3734626

```