

回帰分析では`lm()`ではなく
`estimatr::lm_robust()`を使おう

2022/07/23

第100回R勉強会@東京 (#TokyoR)

森下光之助 (@dropout009)

森下光之助

TVISION INSIGHTS株式会社
データサイエンティスト
執行役員（データ・テクノロジー担当）

テレビの視聴行動を分析しています

データの利活用、マネジメント、組織づくり、
因果推論、機械学習の解釈手法などに興味があります

Twitter: @dropout009

Speaker Deck: dropout009

Blog: <https://dropout009.hatenablog.com/>

機械学習を 解釈する技術

予測力と説明力を両立する実践テクニック

著者: 森下光之助



Techniques for Interpreting Machine Learning

そのモデルの振る舞いを 説明できますか？

あらゆる予測モデルを解釈する4つの手法PFI, PD, ICE, SHAP
特徴量の重要度/特徴量と予測値の関係性/インスタンスごとの異質性/予測の理由

技術評論社

Rで回帰分析をするときは

標準の`lm()`だと

不均一分散に対応できないので

`estimatr::lm_robust()`を使おう

という話をします

$\tau_m(\cdot)$ と回帰係数の標準誤差

シミュレーションデータをもとに、`lm()`の挙動を確認していく

以下のようなシミュレーションデータに対して、線形回帰モデルをあてはめる

$$Y = X + U$$

$$X \sim \text{Uniform}(-2, 2)$$

$$U \sim \mathcal{N}(0, 1)$$

誤差項 U の分散が X に依存しない(=均一分散)

```
N <- 50
```

```
df <- tibble(
```

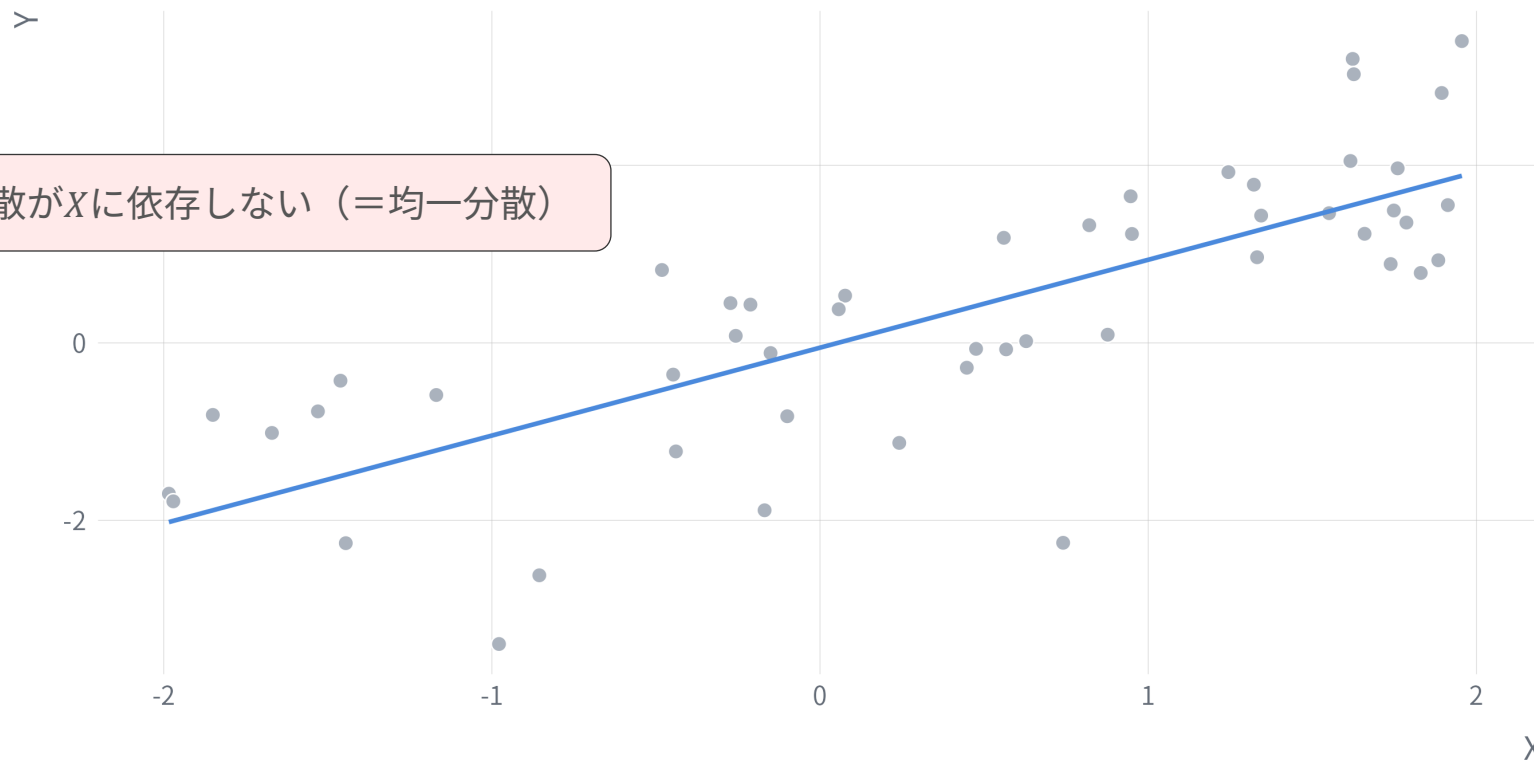
```
  x = runif(N, -2, 2),
```

```
  u = rnorm(N, 0, 1),
```

```
  y = x + e
```

```
)
```

シミュレーションデータと回帰直線



シミュレーションデータをもとに、`lm()`の挙動を確認していく

シミュレーションデータに単回帰モデル

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

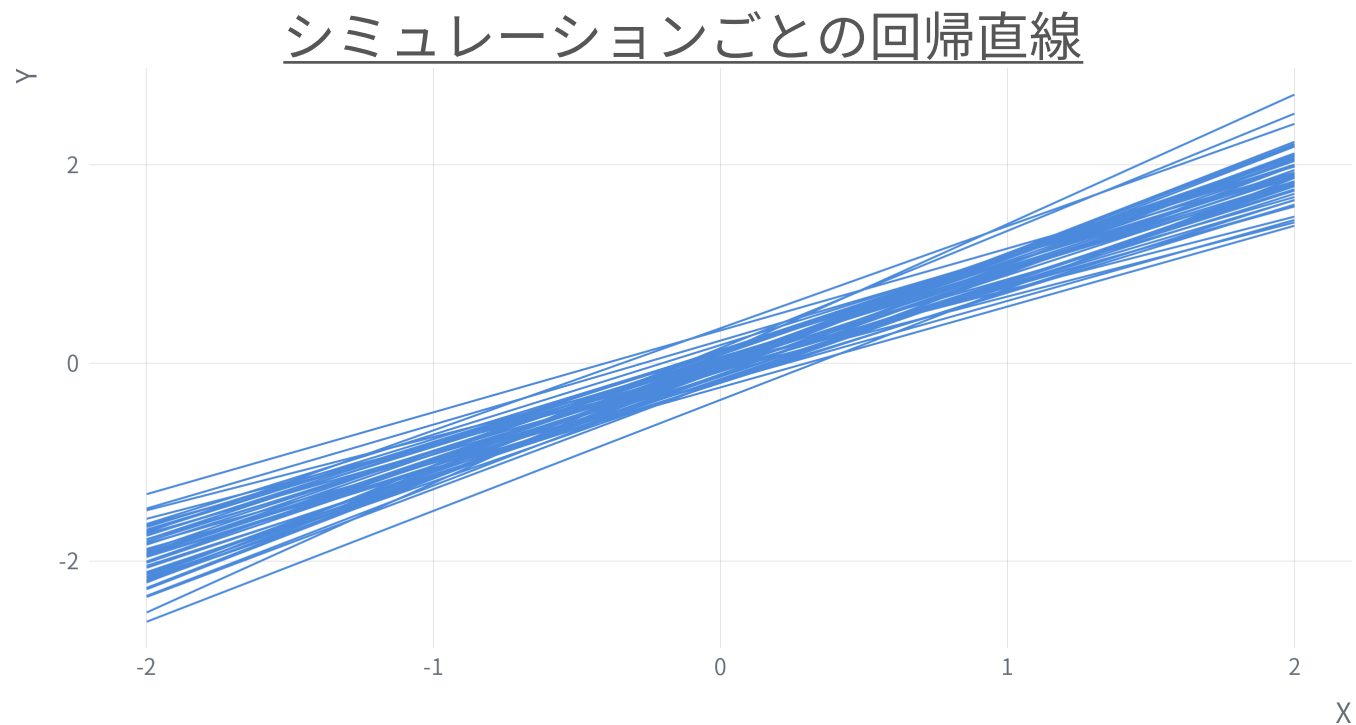
をあてはめる。`lm()`の出力結果から回帰係数の大きさや標準誤差がわかる
信頼区間やp値はこの標準誤差を使って計算される

```
> df %>%  
+   lm(y ~ x, data = .) %>%  
+   tidy()  
  
# A tibble: 2 × 5  
  term          estimate std.error statistic  p.value  
  <chr>         <dbl>    <dbl>    <dbl>    <dbl>  
1 (Intercept) -0.0542   0.146    -0.371  7.12e- 1  
2 x            0.990    0.115     8.58   2.93e-11
```

回帰係数 $\hat{\beta}$ の標準誤差

回帰係数の標準誤差ってなに？

- 観測されたデータはあくまで母集団から抽出されたサンプルで、回帰分析ではそこから母集団のパラメータを推測している
- なので、同じシミュレーション設定で何度もデータを生成し回帰直線を引くと、生成されるデータのばらつきに応じて少しずつ違う回帰直線が引かれる
- この回帰直線のばらつき（標準偏差）を推定しているのが回帰係数の標準誤差



実際に何度もデータを生成して回帰分析の結果を確認していく

```
N <- 50  
M <- 5000
```

```
df <- tibble(  
  m = rep(seq_len(M), each = N),  
  x = runif(N * M, -2, 2),  
  u = rnorm(N * M, 0, 1),  
  y = x + u  
)
```

```
df_result <- df %>%  
  nest(data = c(x, u, y)) %>%  
  mutate(result = map(data, ~ tidy(lm(y ~ x, data = .)))) %>%  
  select(m, result) %>%  
  unnest(c(m, result)) %>%  
  filter(term == "x")
```

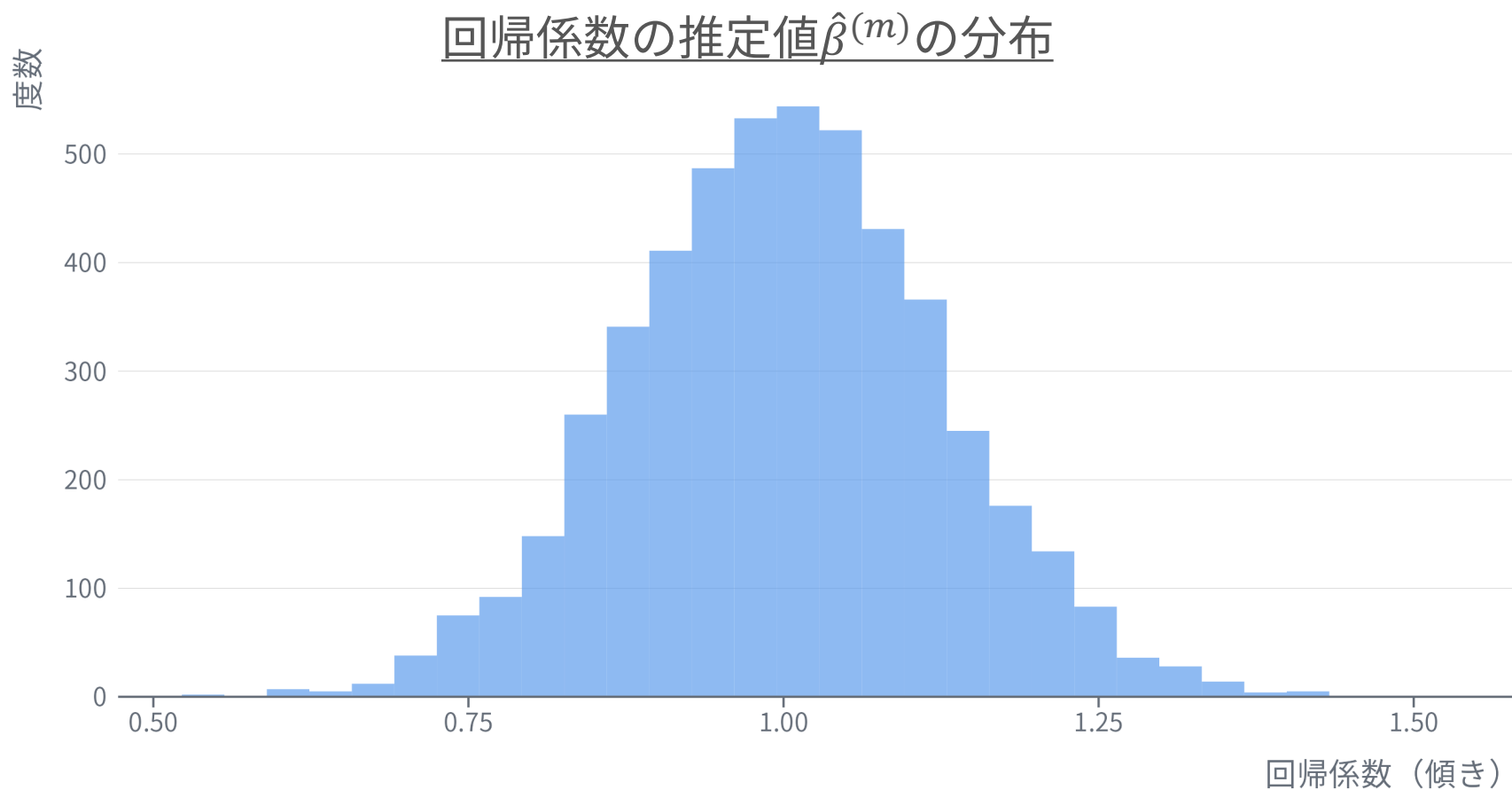
$Y = X + U$
 $X \sim \text{Uniform}(-2, 2)$
 $U \sim \mathcal{N}(0, 1)$

$\{(Y_i^{(m)}, X_i^{(m)})\}_{i=1}^N$ から単回帰で
 $(\hat{\alpha}^{(m)}, \hat{\beta}^{(m)})$ を推定

傾き $\hat{\beta}^{(m)}$ に注目

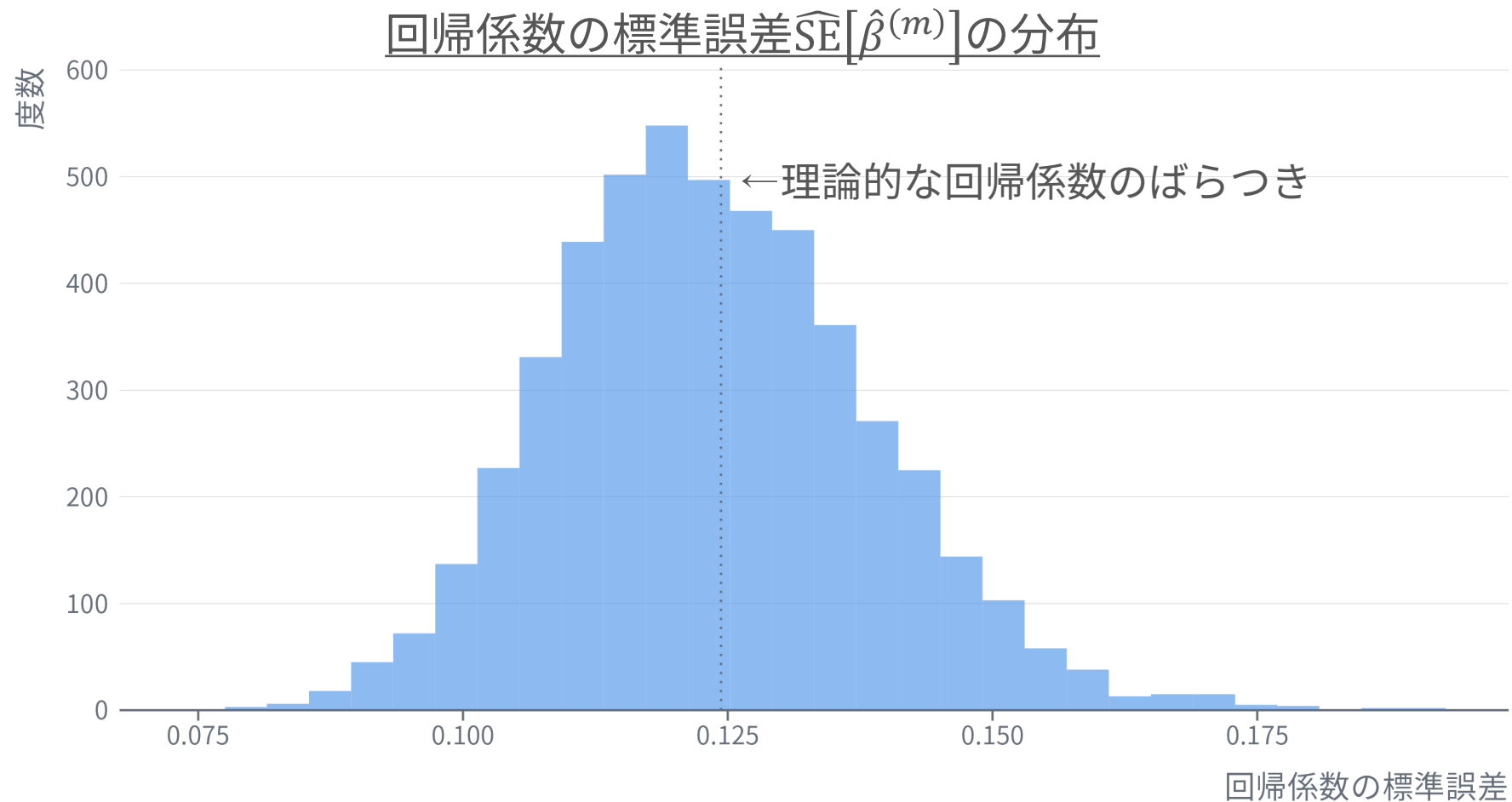
シミュレーションごとに回帰係数の推定値 $\hat{\beta}^{(m)}$ の分布を確認

- 回帰係数はシミュレーション毎に少し異なる値になるので、結果の分布を確認
- この分布の標準偏差 $SD[\hat{\beta}^{(1)}, \dots, \hat{\beta}^{(M)}]$ が回帰係数の理論的なばらつき（のシミュレーションによる近似）
- 回帰係数の標準誤差 $\widehat{SE}[\hat{\beta}^{(m)}]$ は回帰係数は理論的なばらつきを推定している



$\text{lm}()$ の標準誤差は回帰係数のばらつきをうまく推定できている

- 回帰係数の標準誤差はシミュレーション毎に少し異なる値になるので、結果の分布を確認
- 標準誤差は理論的な回帰係数のばらつきを中心に分布しており、ばらつきを概ねうまく推定できていることがわかる（ほんのちょっと小さめに推定してはいる）



不均一分散と $\ln(\cdot)$ の標準誤差

誤差項 U の分散が X に依存している（＝分散が一定でない）場合を考える

誤差項 U の分散が X に依存（＝不均一分散）

$$Y = X + U$$

$$X \sim \text{Uniform}(-2, 2)$$

$$U \sim \mathcal{N}(0, X^2)$$

```
N <- 50
```

```
df_heterogeneous <- tibble(
```

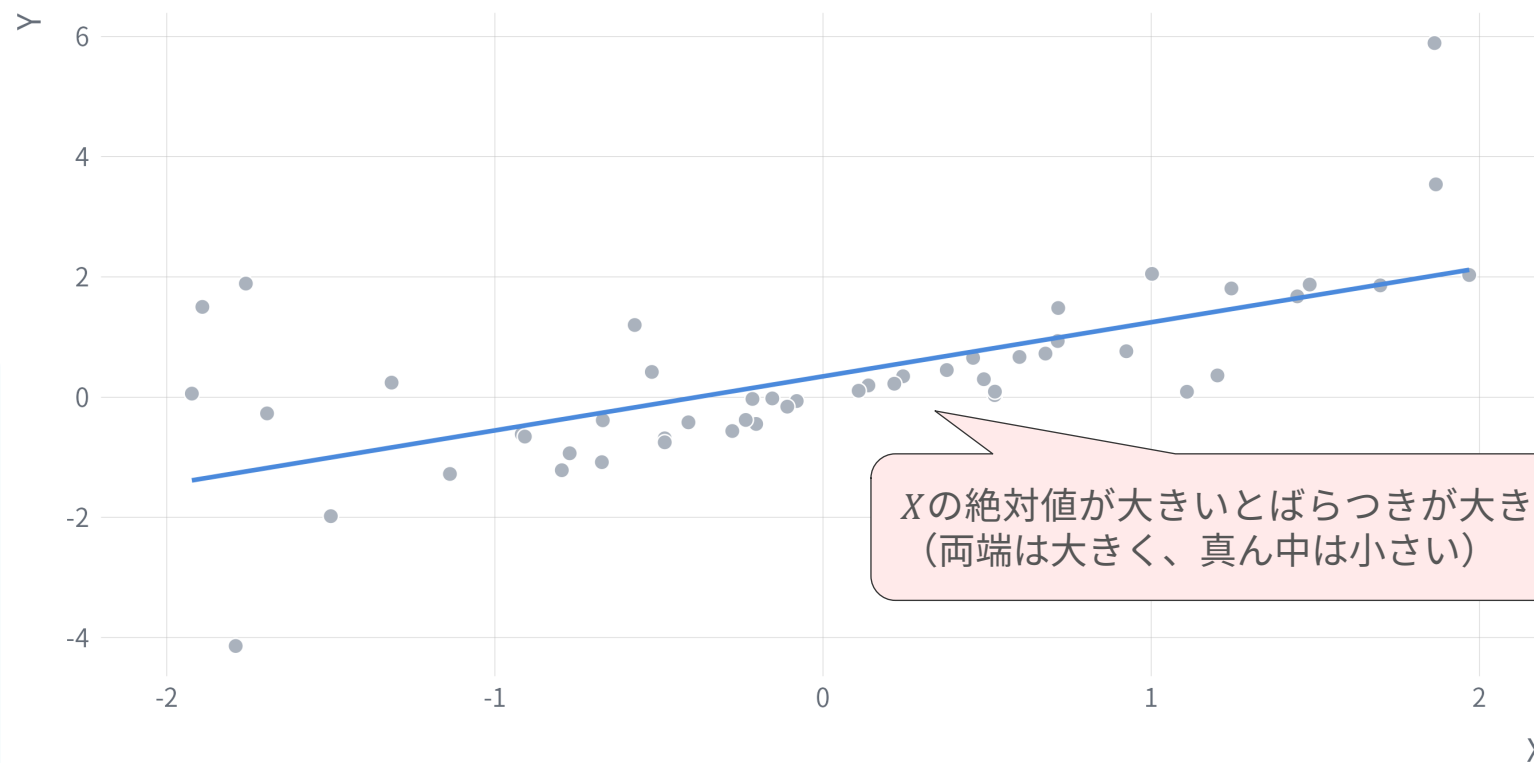
```
  x = runif(N, -2, 2),
```

```
  u = rnorm(N, 0, abs(x)),
```

```
  y = x + u
```

```
)
```

シミュレーションデータと回帰直線



Xの絶対値が大きいとばらつきが大きい
(両端は大きく、真ん中は小さい)

実際に何度もデータを生成して回帰分析の結果を確認

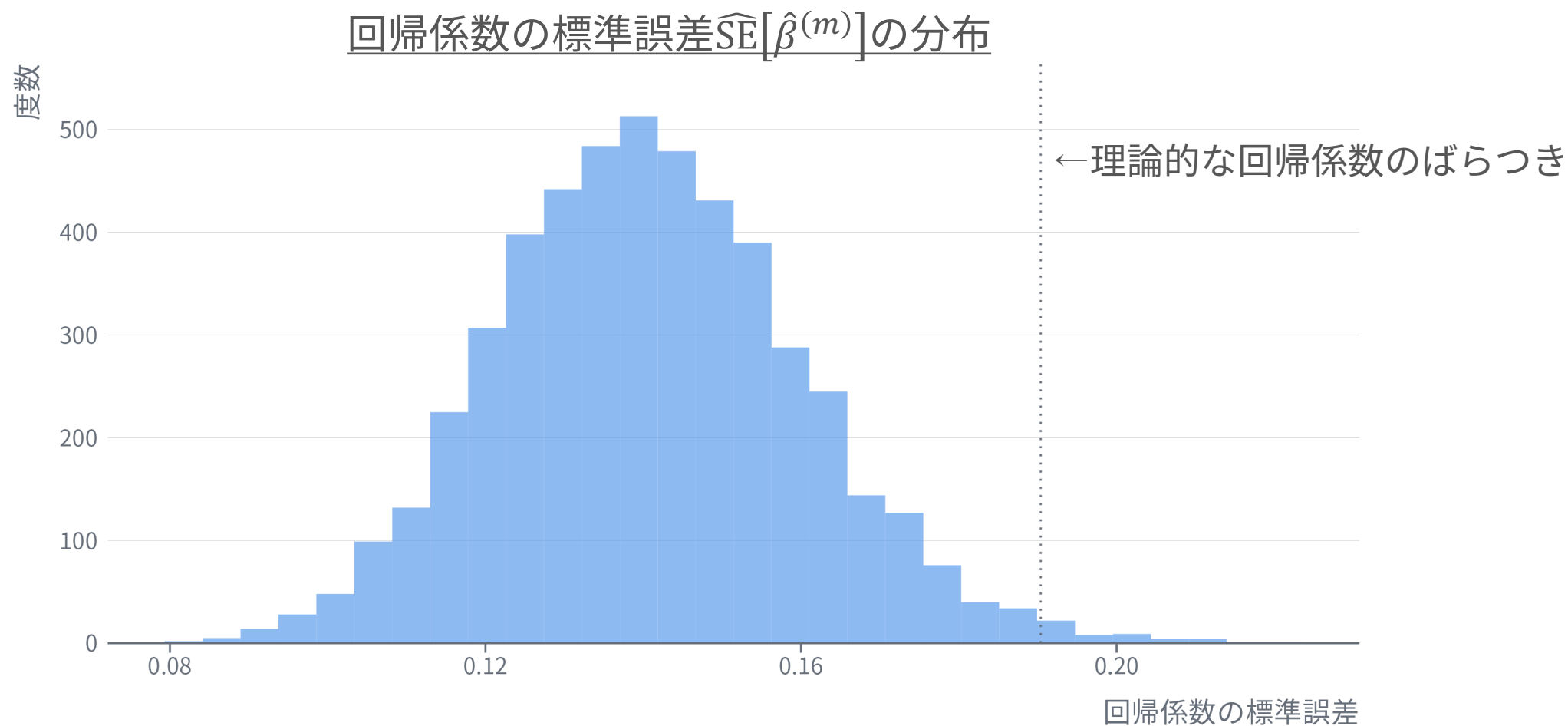
```
N <- 50
M <- 5000

df_simulation_heterogeneous <- tibble(
  m = rep(seq_len(M), each = N),
  x = runif(N * M, -2, 2),
  u = rnorm(N * M, 0, abs(x)),
  y = x + u
)

df_result_heterogeneous <- df_simulation_heterogeneous %>%
  nest(data = c(x, u, y)) %>%
  mutate(result = map(data, ~ tidy(lm(y ~ x, data = .)))) %>%
  select(m, result) %>%
  unnest(c(m, result)) %>%
  filter(term == "x")
```

$\text{lm}()$ の標準誤差は回帰係数のばらつきを過小評価してしまっている

- 回帰係数の標準誤差はシミュレーション毎に少し異なる値になるので、結果の分布を確認
- 標準誤差は理論的な回帰係数のばらつきを過小評価してしまっている



なぜlm()の標準誤差は回帰係数のばらつきを過小評価してしまうのか？

lm()で計算される標準誤差は
均一分散を仮定しているから

lm()の標準誤差は
どうやって計算されているのか

そもそもlm()の標準誤差はどうやって計算されているのか？

回帰係数の標準誤差を考えるため、線形回帰モデル

$$Y = \mathbf{X}'\boldsymbol{\beta} + U$$

を考える

ここで、 \mathbf{X} と $\boldsymbol{\beta}$ は共に $K \times 1$ の行列

$$\mathbf{X} = (X_1, \dots, X_k, \dots, X_K)'$$

$$\boldsymbol{\beta} = (\beta_1, \dots, \beta_k, \dots, \beta_K)'$$

とする

典型的な仮定として $\mathbb{E}[U | \mathbf{X}] = 0$ を置く。このとき、誤差項の条件付き分散 $\sigma^2(\mathbf{X})$ は

$$\sigma^2(\mathbf{X}) = \text{Var}[U | \mathbf{X}] = \mathbb{E}[U^2 | \mathbf{X}]$$

となる

加えて以下も仮定しておく

$$\mathbb{E}[Y^2] < \infty, \mathbb{E}[\|\mathbf{X}\|] < \infty, \mathbb{E}[\mathbf{X}\mathbf{X}'] > 0$$

最小二乗法による回帰係数の推定

ランダムサンプリングされたデータ $\{(Y_1, \mathbf{X}_1), \dots, (Y_i, \mathbf{X}_i), \dots, (Y_N, \mathbf{X}_N)\}$ から最小二乗法で回帰係数を推定する

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin}_{\mathbf{b}} \sum_{i=1}^N (Y_i - \mathbf{X}_i' \mathbf{b})^2 = \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i Y_i \right)$$

この最小二乗推定量 $\hat{\boldsymbol{\beta}}$ の信頼性（標準誤差）を知りたい

後々便利なので、 $\hat{\boldsymbol{\beta}}$ を以下のように変形できることを確認しておく

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i Y_i \right) \\ &= \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i (\mathbf{X}_i' \boldsymbol{\beta} + U_i) \right) \\ &= \boldsymbol{\beta} + \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i U_i \right) \end{aligned}$$

(参考) 回帰係数の条件付き期待値

最小二乗法推定量の信頼性（標準誤差）を見る前に、まずは条件付き期待値を確認しておく

条件付き期待値は

$$\begin{aligned}\mathbb{E}[\hat{\boldsymbol{\beta}} \mid \mathbf{X}_1, \dots, \mathbf{X}_N] &= \mathbb{E} \left[\boldsymbol{\beta} + \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i U_i \right) \mid \mathbf{X}_1, \dots, \mathbf{X}_N \right] \\ &= \boldsymbol{\beta} + \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbb{E}[U_i \mid \mathbf{X}_i] \right) \\ &= \boldsymbol{\beta}\end{aligned}$$

となり、最小二乗法推定量は（条件付き）不偏性を満たすことがわかる

回帰係数の条件付き分散

回帰係数の標準誤差は、回帰係数の条件付き分散の平方根をとったものなので、とりあえず条件付き分散を考える

$$\begin{aligned}\text{Var}[\hat{\boldsymbol{\beta}} \mid \mathbf{X}_1, \dots, \mathbf{X}_N] &= \text{Var} \left[\boldsymbol{\beta} + \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i U_i \right) \mid \mathbf{X}_1, \dots, \mathbf{X}_N \right] \\ &= \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \text{Var}[U_i \mid \mathbf{X}_i] \right) \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \\ &= \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \mathbb{E}[U_i^2 \mid \mathbf{X}_i] \right) \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \\ &= \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \sigma_i^2 \right) \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1}\end{aligned}$$

$i \neq j$ の場合は
 $\text{Cov}[U_i, U_j \mid \mathbf{X}_1, \dots, \mathbf{X}_N]$
 $= \mathbb{E}[U_i U_j \mid \mathbf{X}_1, \dots, \mathbf{X}_N]$
 $= \mathbb{E}[U_i \mid \mathbf{X}_i] \mathbb{E}[U_j \mid \mathbf{X}_j]$
 $= 0$

ここで、 $\sigma_i^2 = \sigma^2(\mathbf{X}_i) = \mathbb{E}[U_i^2 \mid \mathbf{X}_i]$ としている

回帰係数の分散の推定：均一分散を仮定する場合

N 個の σ_i^2 を推定するのは大変そうなので均一分散 $\sigma_i^2 = \sigma^2$ を仮定すると話が単純になる

$$\begin{aligned}\text{Var}[\hat{\beta} \mid \mathbf{X}_1, \dots, \mathbf{X}_N] &= \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \sigma_i^2 \right) \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \\ &\approx \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \sigma^2 \right) \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \\ &= \sigma^2 \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1}\end{aligned}$$

均一分散の仮定による近似
(成り立たない場合は推定
がうまくいかない)

分散 σ^2 を

$$s^2 = \frac{1}{N-K} \sum_{i=1}^N \hat{U}_i^2 = \frac{1}{N-K} \sum_{i=1}^N (Y_i - \mathbf{X}_i' \hat{\beta})^2$$

で推定して置き換えると、（均一分散の仮定のもとでの）回帰係数の分散が推定できる

$$\widehat{\text{Var}}[\hat{\beta} \mid \mathbf{X}_1, \dots, \mathbf{X}_N] = s^2 \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1}$$

あとは平方根を取れば標準誤差になる。lm()ではこうして標準誤差を計算している

不均一分散だとなぜ標準誤差をうまく推定できないのか

話をわかりやすくするために、単回帰モデル

$$Y = \alpha + \beta X + U$$

を考える。単回帰モデルの最小二乗推定量は

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

$$\hat{\beta} = \frac{\sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

なので、傾き $\hat{\beta}$ の分散は以下になる。

$$\text{Var}[\hat{\beta} \mid X_1, \dots, X_N] = \frac{\sum_{i=1}^N (X_i - \bar{X})^2 \sigma_i^2}{\left(\sum_{i=1}^N (X_i - \bar{X})^2\right)^2} = \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right) \frac{\sum_{i=1}^N (X_i - \bar{X})^2 \sigma_i^2}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

一方、均一分散での近似は

$$\text{Var}[\hat{\beta} \mid X_1, \dots, X_N] \approx \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right) \sigma^2$$

なので、平均からの乖離で σ_i^2 を加重平均するべきところをワンナンバー σ^2 で近似している。

実際に均一分散が成り立つなら問題ないが、不均一分散の場合は乖離が生まれてしまう。

分散を $(X_i - \bar{X})^2$ で重みをつけて平均していると解釈できる（平均から乖離しているほど重くする）

不均一分散が標準誤差に与える影響をシミュレーションで確認

均一分散、平均周りがばらついている不均一分散①、平均から離れた部分がばらついている不均一分散②の3種類を考える。全体としてのばらつきは同じ。

$$Y = X + U$$

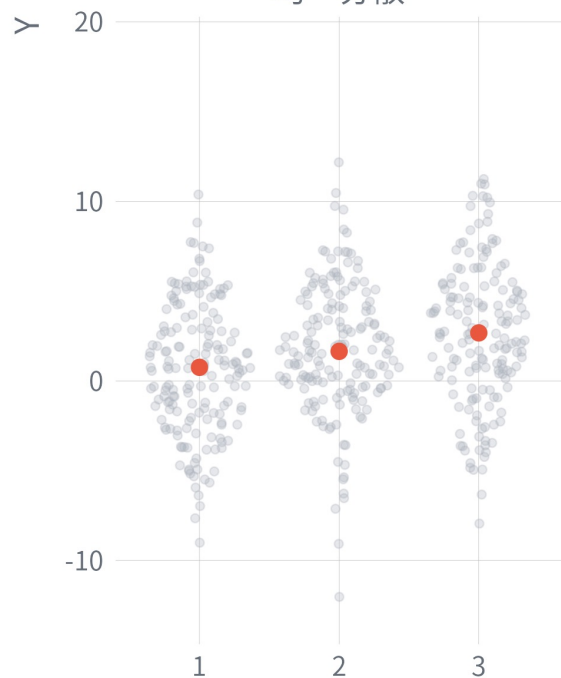
$$X \sim \text{Uniform}(\{1, 2, 3\})$$

$$U \sim \mathcal{N}(0, \sigma^2(X))$$

誤差項 U の分散が X に依存

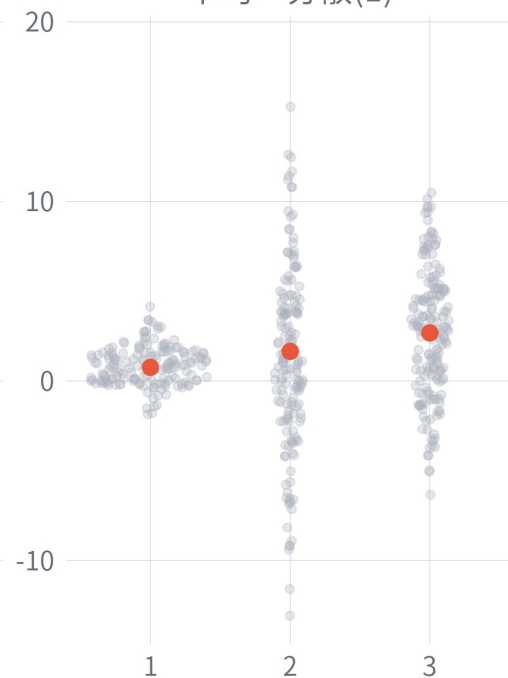
$$\sigma^2(X) = \begin{cases} 16^2 & \text{if } X = 1 \\ 16^2 & \text{if } X = 2 \\ 16^2 & \text{if } X = 3 \end{cases}$$

均一分散



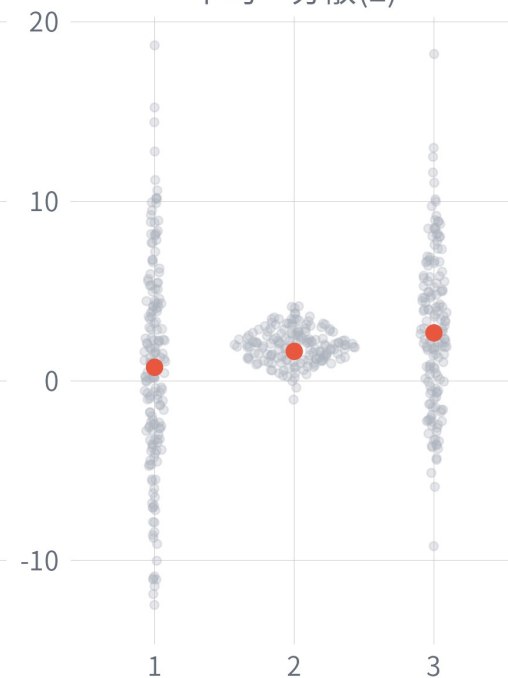
$$\sigma^2(X) = \begin{cases} 1^2 & \text{if } X = 1 \\ 31^2 & \text{if } X = 2 \\ 16^2 & \text{if } X = 3 \end{cases}$$

不均一分散(1)



$$\sigma^2(X) = \begin{cases} 31^2 & \text{if } X = 1 \\ 1^2 & \text{if } X = 2 \\ 16^2 & \text{if } X = 3 \end{cases}$$

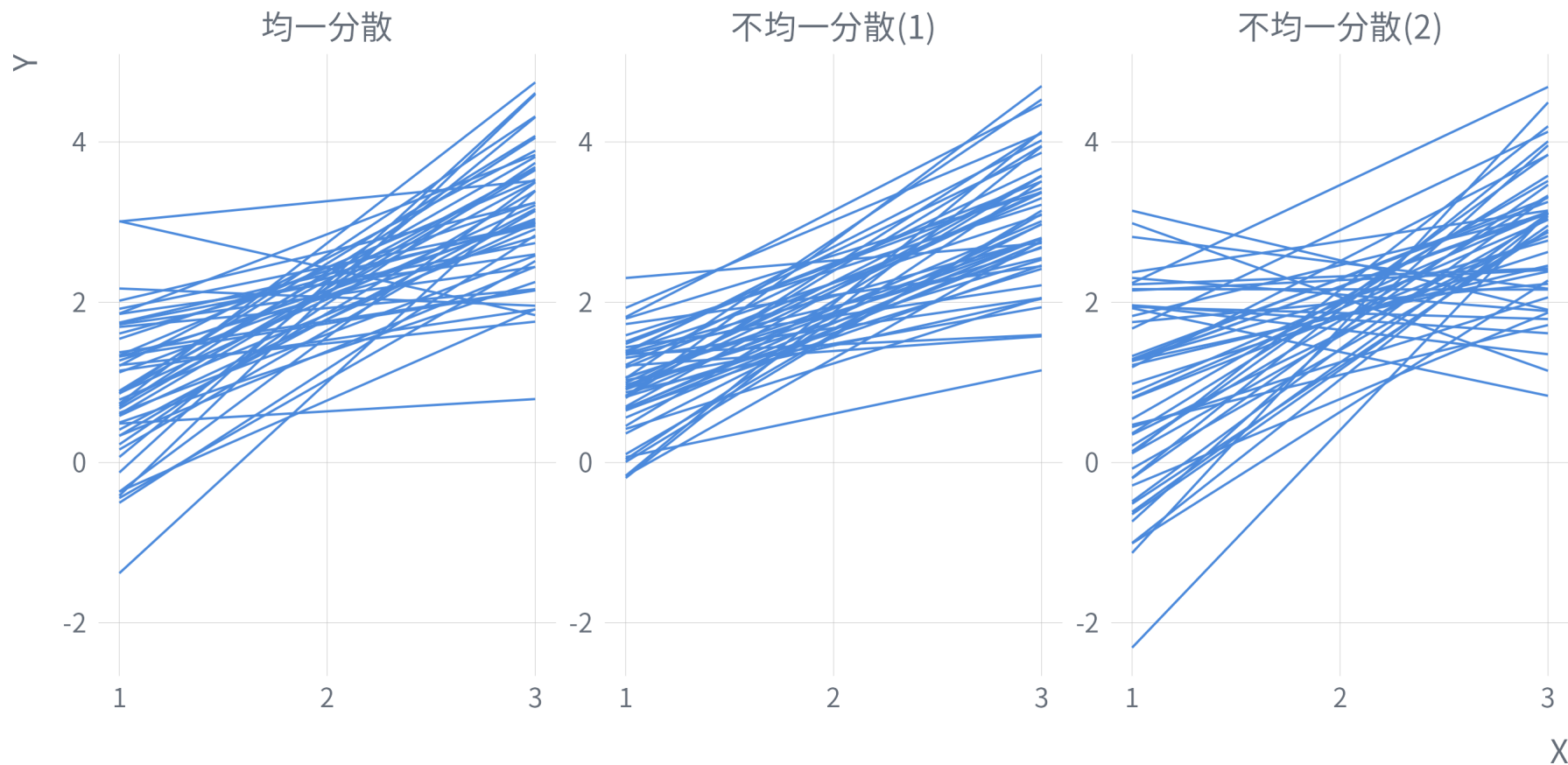
不均一分散(2)



均一分散と不均一分散で回帰直線のばらつきを比較

均一分散と比較して、不均一分散①（平均周りの分散大、左端の分散小）の場合は回帰直線のばらつきが小さく、不均一分散②（左端の分散大、平均周りの分散小）の場合はばらつきが大きい

シミュレーションごとの回帰直線



均一分散を仮定した標準誤差は理論的なばらつきと乖離する

全体としては同程度の分散でも、不均一分散がある場合は回帰係数の理論的なばらつきは変わってくる

lm()の推定では均一分散を仮定しているので、全体の分散が同程度なら同じ標準誤差を返してしまう

分散タイプごとのlm()の出力

分散タイプ	回帰係数の理論的なばらつき	lm()で推定された標準誤差の平均
均一分散	0.70	0.70
不均一分散①	0.52	0.70
不均一分散②	0.84	0.69

Robust Standard Errorsで 不均一分散に対処する

不均一分散を考慮した標準誤差の推定を考える

回帰係数の分散

$$\text{Var}[\hat{\beta} \mid \mathbf{X}_1, \dots, \mathbf{X}_N] = \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \mathbb{E}[U_i^2 \mid \mathbf{X}_i] \right) \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1}$$

を推定するにあたって、均一分散を仮定する場合は誤差項の平均的な分散を推定した s^2 で置き換えていた

$$s^2 = \frac{1}{N - K} \sum_{i=1}^N \hat{U}_i^2$$

ここで、 $\mathbb{E}[U_i^2 \mid \mathbf{X}_i]$ をワンナンバー s^2 で近似するのではなく、直接残差の二乗 \hat{U}_i^2 で置き換えるやり方も考えられる

$$\widehat{\text{Var}}[\hat{\beta} \mid \mathbf{X}_1, \dots, \mathbf{X}_N] = \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \hat{U}_i^2 \right) \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \right)^{-1}$$

個別の \hat{U}_i^2 はサンプルサイズ1なので誤差があるが、 N が十分に大きくなれば $\frac{1}{N} \sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \hat{U}_i^2$ はうまく推定できることを利用

上記の平方根をとったものを Heteroskedasticity-Consistent Standard Errors や Robust Standard Errors と言う。Robust SE はいくつか手法が提案されていて、これは「HC0」と呼ばれているもの

RでRobust SEを計算するには？

`estimatr::lm_robust()`を

使えば計算できる

(再掲) 不均一分散のシミュレーションデータ

誤差項 U の分散が X に依存 (=不均一分散)

$$Y = X + U$$

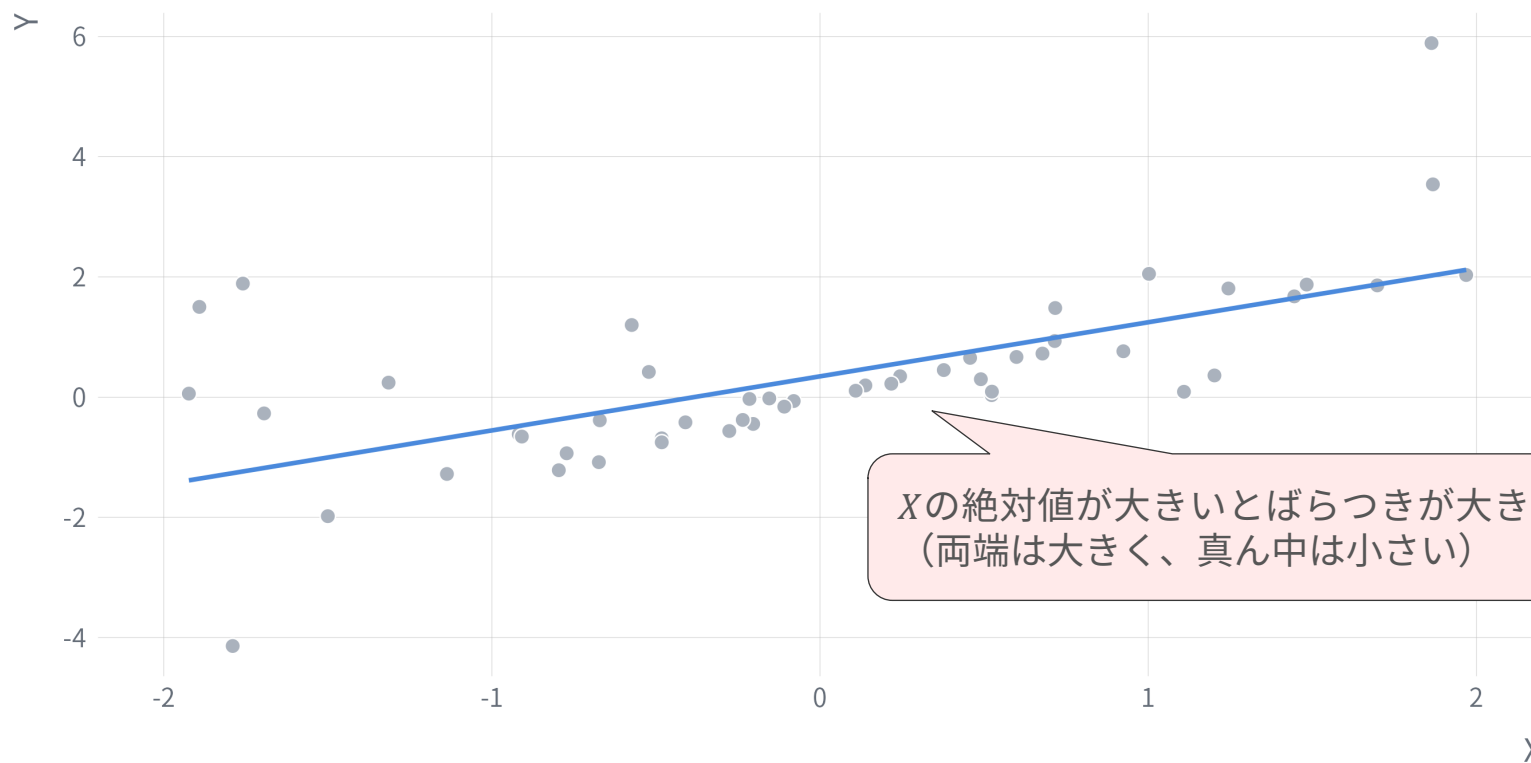
$$X \sim \text{Uniform}(-2, 2)$$

$$U \sim \mathcal{N}(0, X^2)$$

```
N <- 50
```

```
df_heterogeneous <- tibble(  
  x = runif(N, -2, 2),  
  u = rnorm(N, 0, abs(x)),  
  y = x + u  
)
```

シミュレーションデータと回帰直線



X の絶対値が大きいとばらつきが大きい
(両端は大きく、真ん中は小さい)

estimatr::lm_robust()の挙動を確認

estimatr::lm_robust()でse_typeを指定することでRobust SEを計算できる
標準誤差を確認すると、Robust SEは誤差を大きめに見積もっている

```
> df_heterogeneous %>%  
+   estimatr::lm_robust(y ~ x, data = ., se_type = "classical") %>%  
+   tidy()  
      term estimate std.error statistic      p.value  conf.low conf.high df outcome  
1 (Intercept) 0.3456885 0.1594783  2.167621 3.517795e-02 0.02503593 0.666341 48      y  
2           x 0.8999665 0.1531487  5.876423 3.869962e-07 0.59204038 1.207893 48      y
```

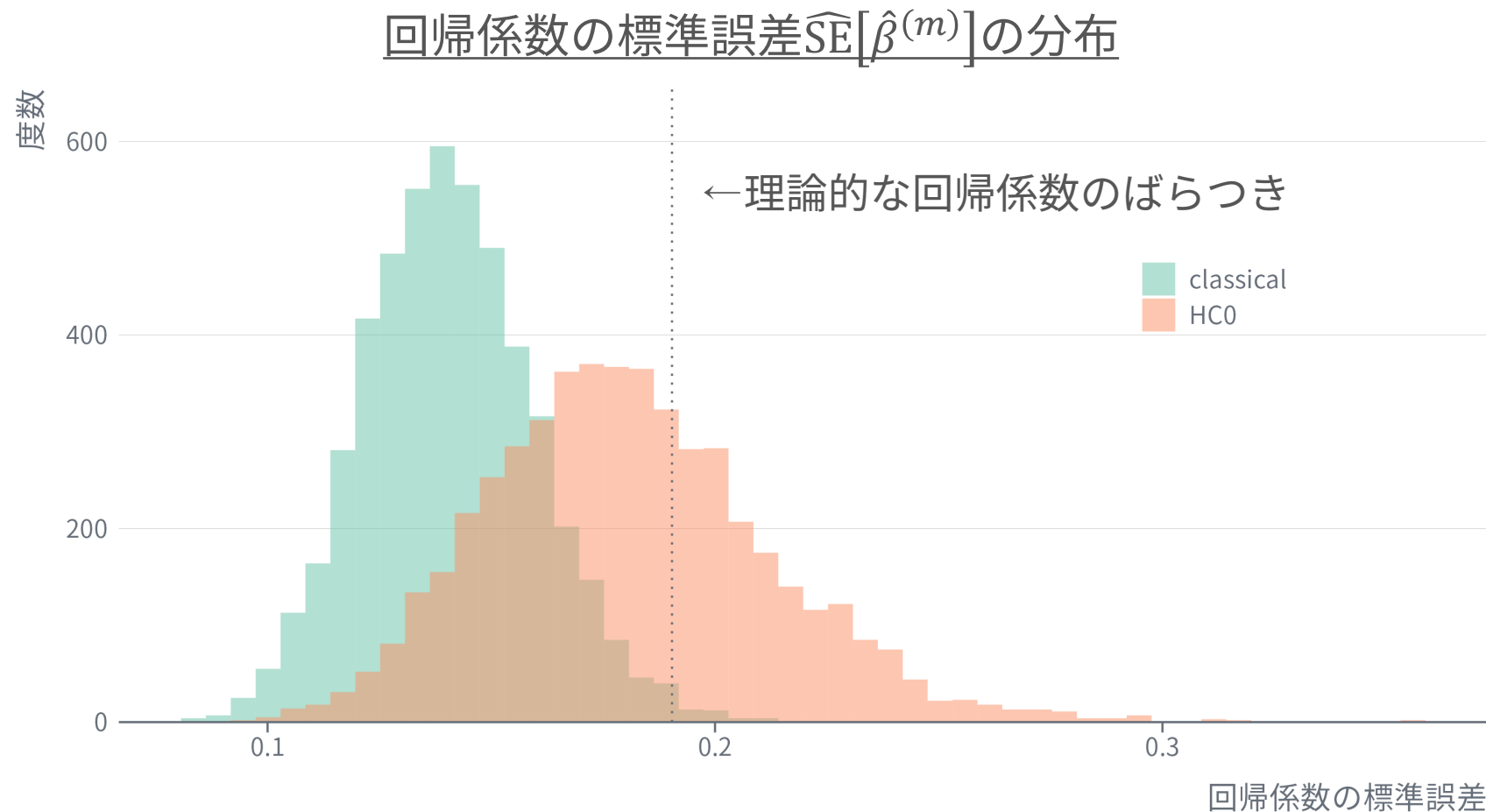
均一分散を仮定した標準誤差
(lm()の出力と同じ)

```
> df_heterogeneous %>%  
+   estimatr::lm_robust(y ~ x, data = ., se_type = "HC0") %>%  
+   tidy()  
      term estimate std.error statistic      p.value  conf.low conf.high df outcome  
1 (Intercept) 0.3456885 0.1563591  2.210862 0.031843538 0.03130734 0.6600696 48      y  
2           x 0.8999665 0.2389852  3.765784 0.000453233 0.41945455 1.3804784 48      y
```

不均一分散対してRobustな標準誤差

Robust SEを使うことで標準誤差の過小評価を軽減できる

- 均一分散を仮定した標準誤差 (classical) と比べて、Robust SE (HC0) は標準誤差を大きく見積もっている
- Robust SEは理論的な回帰係数のばらつきよりも少し小さい傾向。これはサンプルサイズが50と小さいことによる



回帰係数の標準誤差をさらにうまく推定する

- HC0を改良した推定量がいくつか提案されている
- サンプルサイズが十分に大きければHC0 - HC3はどれも似たような値をとるが、小さい場合はHC2やHC3がベター。HC0 < HC2 < HC3なので、HC3がもっとも保守的
- `estimatr::lm_robust()`のデフォルトはHC2（ちなみにStataのデフォルトはHC1）

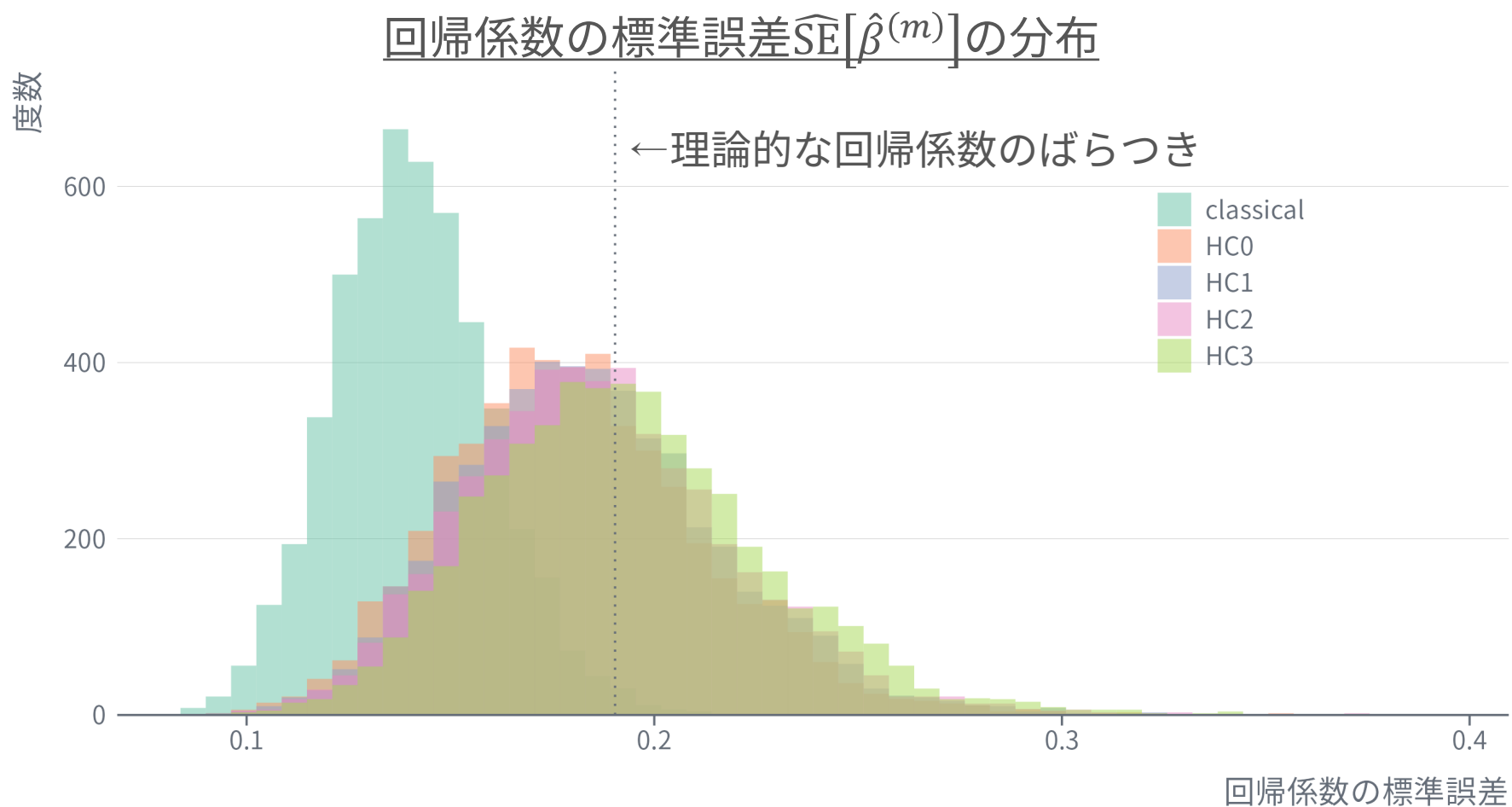
名称	計算式
HC0	$\left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i'\right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \hat{U}_i^2\right) \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i'\right)^{-1}$
HC1	$\frac{N}{N-K} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i'\right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \hat{U}_i^2\right) \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i'\right)^{-1}$
HC2	$\left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i'\right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \frac{\hat{U}_i^2}{1-h_{ii}}\right) \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i'\right)^{-1}$
HC3	$\left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i'\right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i' \frac{\hat{U}_i^2}{(1-h_{ii})^2}\right) \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i'\right)^{-1}$

$$h_{ii} = \mathbf{X}_i' \left(\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i'\right)^{-1} \mathbf{X}_i$$

はleverageと呼ばれていて、観測値 x_i の他の観測値と比較してどのくらい外れているかを表している。 $0 \leq h_{ii} \leq 1$ なので、HC0 < HC2 < HC3。詳細はHansen(2022)を参照

Robust SEの比較

- HC0-HC3は似たような傾向を示す
- HC1-HC3はHC0よりもやや標準誤差を大きめに見積もる。このデータ例では、より保守的な標準誤差を出すHC3がやや優勢か



まとめ

まとめ

- Rは線形回帰分析ができる関数`lm()`をデフォルトで備えているが、`lm()`で計算される回帰係数の標準誤差は均一分散を仮定している。よって、不均一分散の状況下では標準誤差の推定にバイアスが生じてしまう
- 不均一分散を考慮した標準誤差としてRobust SEがある。実際のところ均一分散の仮定が成り立つことはあまりないと考えられるので、とりあえずRobust SEを使っておくことが推奨されている
- Rでは、`estimatr`パッケージの`lm_robust()`を使えばRobust SEを計算できる。`sandwich`パッケージの`vcovHC()`でも計算できるが、`lm_robust()`なら`lm()`を置き換えるだけで済むので便利
- Robust SEにはHC0, HC1, HC2, HC3のようにいくつかタイプがある。サンプルサイズが十分に大きければほとんど同じ値になるが、サンプルサイズが小さい場合はHC2やHC3がベターと言われている

参考文献

参考文献

- Hansen, Bruce E. "Econometrics." (2022). <https://www.ssc.wisc.edu/~bhansen/econometrics/>.
- Blair G, Cooper J, Coppock A, Humphreys M, Sonnet L. (2022). estimatr: Fast Estimators for Design-Based Inference. <https://declaredesign.org/r/estimatr/>. <https://github.com/DeclareDesign/estimatr>.

R

Rを使えるデータサイエンティスト、エンジニアを大募集

T»VISION
INSIGHTS

弊社では以下の職種を募集中です

- ☑ バックエンドエンジニア
- ☑ データエンジニア
- ☑ データサイエンティスト

エントリーはこちらから！

TVISION INSIGHTS 募集

