

# Rによる計量分析：データ解析と可視化

## 回帰分析への導入

伊藤 岳

富山大学 経済学部 2017 年度後期



Email: [gito@eco.u-toyama.ac.jp](mailto:gito@eco.u-toyama.ac.jp)

December 18, 2017

# Agenda

- 1 中間課題 (1) の講評
- 2 統計的仮説検定の解釈と統計的過誤 (復習)
- 3 線形回帰分析
- 4 線形回帰分析：R コード

# 中間課題 (1)

## R スクリプトで実行する内容

- 1 連続一様分布  $X \sim U(-10, 10)$  から、大きさ (標本サイズ)  $n = 1,000$  の標本を抽出
    - ▶ 期待値 (母平均) は,  $(-10 + 10)/2 = 0$
  - 2 1 回毎に、標本平均と 99%信頼区間を計算して記録する
  - 3 以上の作業を  $m = 10,000$  回繰り返す
  - 4  $m$  回の結果をまとめて、99%信頼区間をプロットする
    - ▶ ただし、図にする際は見やすいよう、1,000 の標本の結果をランダムに抽出する
- 
- ▶ 90%信頼区間:  $[\bar{X}_n - 1.65SE, \bar{X}_n + 1.65SE]$
  - ▶ 95%信頼区間:  $[\bar{X}_n - 1.96SE, \bar{X}_n + 1.96SE]$
  - ▶ 99%信頼区間:  $[\bar{X}_n - 2.56SE, \bar{X}_n + 2.56SE]$
  - ▶ 95%信頼区間のシミュレーションをしているテンプレートでは「1.96」を使っているの、どこかを変えればよい
    - ▶  $\alpha$  が有意水準のとき、 $100 \times (1 - \alpha)$  が信頼係数

# 中間課題 (1)

## 概要

- ▶ 提出者：4名 (11630806, 11630005, 11430078, 11630043)
- ▶ 配点：20 (/100. シラバス通り)
- ▶ 注意：(1) 学籍番号の誤植, (2) メール本文 (減点対象)

## 確認したかった理解とスキル

- ▶ R スクリプトの作成と実行
  - ▶ 任意のディレクトリに処理結果 (今回は図) を出力する
  - ▶ 一般的なスクリプトの一部を変更することで、異なる処理を行なう
- ▶ 信頼区間の理解と解釈
  - ▶ 信頼区間の意味
  - ▶ 信頼係数と有意水準

# 帰無仮説, 対立仮説と仮説検定の解釈

## 検定結果の解釈と注意

- ▶ 「(有意水準  $\alpha$  で) 統計的に有意 (statistically significant)」: 「 $H_0$  を棄却した」( $H_0$  が正しいとすると, 起こり得ないような稀な現象が起きた)
- ▶ 帰無仮説  $H_0$  が棄却されれば, 対立仮説  $H_1$  の正当性を強く主張できる
- ▶ 帰無仮説  $H_0$  が棄却されなかったからといって,  $H_0$  が真とは限らない

## 例: 母平均の検定

- ▶ 帰無仮説  $H_0$  を 「 $\mu = \mu_0 = 1$ », 対立仮説  $H_1$  を 「 $\mu \neq 1$ 」とした  $t$  検定の結果,  $H_0$  を棄却できなかった
- ▶ しかし,  $\mu$  は  $\mu_0 = 1$  かもしれないが, 1.1 や 0.9 かもしれない (無数の可能性が残る)
- ▶  $\mu_0 = 1$  を棄却できなかったからといって, 「どの帰無仮説 (1, 1.1, 0.9, etc.) が正しいか」は分からない ( $H_0$  「 $\mu = \mu_0 = 1$ 」が真とは限らない)

# 帰無仮説, 対立仮説と仮説検定の解釈

例：有罪か無罪か (副読本・浅野矢内本, 第7章)

- ▶ 帰無仮説  $H_0$  : 「容疑者  $X$  はこの事件の犯人ではない (無罪である)」
- ▶ 対立仮説  $H_1$  : 「容疑者  $X$  はこの事件の犯人である (無罪ではない)」
- ▶  $H_0$  が棄却できなかった場合にも, 「 $X$  が完全犯罪を達成し, (推定) 無罪を勝ち取った」可能性は残る
- ▶  $H_0$  が棄却できれば, より強い結論を得られる
  - ▶ ただし, 「 $H_0$  が正しく, かつ非常に稀な現象が起きた」可能性は残る (その確率は有意水準  $\alpha$ )
  - ▶ 確率  $\alpha$  で偽陽性 (この例でいえば「冤罪」) が生じ得るから
  - ▶ 有意水準  $\alpha$  は, 偽陽性と偽陰性のバランスを決める (統計的過誤へ)

# 統計的過誤

## 2つの統計的過誤

- ▶ **第一種過誤** (Type I/ $\alpha$  error; 偽陽性 false positive) : 帰無仮説  $H_0$  が正しいにもかかわらず,  $H_0$  を棄却してしまう誤り
- ▶ **第二種過誤** (Type II/ $\beta$  error; 偽陰性 false negative ) : 帰無仮説  $H_0$  が正しくないにもかかわらず,  $H_0$  を棄却できない誤り
  
- ▶ 有意水準  $\alpha$  を厳しく (小さく) すれば, 偽陽性 (誤検出) を減らせるが, 偽陰性 (検出失敗) を増やしてしまう  
例 有意水準  $\alpha$  を厳しくすれば, 冤罪を減らせるが, 完全犯罪を増やしてしまう
- ▶ 有意水準  $\alpha$  は「偽陽性の確率」, あるいは偽陽性と偽陰性のバランスの (ある学術分野における) 「相場」「伝統的に使う目安」に依存する
  - ▶ 社会科学なら多くの場合  $\alpha = 0.05$  (5%水準)
  - ▶ 他によく使われる水準 : 10%, 1%, 0.1%水準

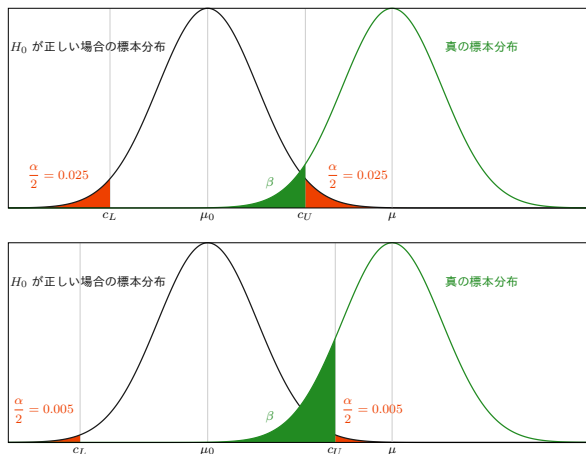
# 統計的過誤

仮説検定の結果	真実	
	$H_0$ が正しい	$H_1$ が正しい
$H_0$ を棄却	第一種過誤 確率: 有意水準 $\alpha$	正しい検定結果 確率: 検出力 $1 - \beta$
$H_0$ を受容	正しい検定結果 確率: $1 - \alpha$	第二種過誤 確率: $\beta$

- ▶ 第一種過誤の確率は、有意水準  $\alpha$  に一致する
- ▶ 第二種過誤の確率  $\beta$  は、(1)  $\alpha$  の大きさと (2) 帰無仮説  $\mu_0$  が真実 ( $\mu$ ) からどれ位離れているかに依存する
  - ▶  $\mu$  を知らないので、標本から  $\beta$  を計算することはできない
  - ▶ (1) を変更するか、(2) を代替的な方法で変更して、 $\beta$  を縮減することはできる
  - ▶ 我々は  $\mu$  を知らないので、 $\mu_0$  を直接変更することは無意味 (できない)
- ▶  $1 - \beta$  を、検定のパワー/検出力 (power) と呼ぶ:  $H_1$  が正しいときに第二種過誤を犯さない確率

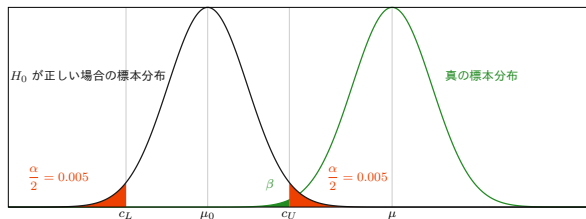
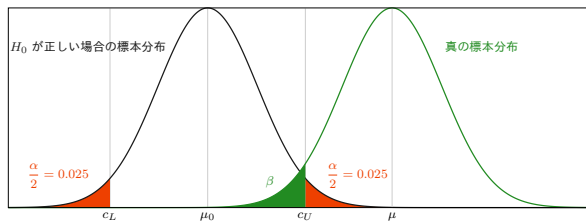


## 第二種過誤と有意水準 $\alpha$



- ▶ 有意水準  $\alpha$  は「厳しければよい」訳ではない：第一種過誤の確率  $\alpha$  は小さくなるが，第二種過誤の確率  $\beta$  は大きくなる
- ▶ 両者はトレード・オフの関係にある

## 第二種過誤と標本サイズ



- ▶ 我々は  $\mu$  を知らないなので、 $\mu_0$  を直接変更することは無意味 (できない)
- ▶ なら、標本サイズ  $n$  を大きくして分散を小さくすればよい
- ▶  $\alpha$  が一定でも、 $n$  が大きくなれば  $\beta$  は小さくなる (検出力は大きくなる)

## 第二種過誤と標本サイズ

### $t$ 検定の R コードを思い出す

- ▶ 母平均  $\mu = 0$  としたシミュレーションで,
  - ▶  $n = 1,000$  のときには,  $H_0$  「 $\mu = \mu_0 = 0.1$ 」を 5%水準で棄却できた
  - ▶  $n = 100$  のときには,  $H_0$  「 $\mu = \mu_0 = 0.1$ 」を 5%水準で棄却できなかった
- ▶  $n$  が小さいときには第二種過誤 (偽陰性) が生じていた
- ▶  $n$  を 100 から 1,000 に増やしたことで,  $\alpha$  を一定 ( $\alpha = 0.05$ ) にしつつ  $\beta$  を小さくできた

# 統計的過誤：まとめ

- ▶ 第一種過誤は偽陽性，第二種過誤は偽陰性
- ▶ 第一種過誤の確率は有意水準  $\alpha$  に等しく，分析者が設定できる
  - ▶ 分野の「お作法」や「慣習」がある (例: 社会科学での  $\alpha = 0.05$ )
  - ▶  $\alpha$  は検定を始める前に決めておく (結果を見てから都合よく変えてはいけない)
- ▶ 第二種過誤の確率  $\beta$  は標本から計算できないが，(1)  $\alpha$  と (2) 帰無仮説  $H_0$  が主張する  $\mu_0$  と真実 ( $\mu$ ) の距離に依存する
- ▶ (1) 有意水準  $\alpha$  を適切に設定し，(2) 標本サイズ  $n$  を大きくすることで，過誤の可能性を小さくできる
  - ▶  $H_0$  が「真実とかけ離れたこと」を主張していれば  $\beta$  は小さいが，我々は真実を知らないので， $H_0$  を直接調整することはできない
- ▶  $1 - \beta$  を検出力という
  - ▶  $\alpha$  が小さいとき，( $\beta$  が大きくなるので) 検出力は小さくなる
  - ▶  $n$  が大きいとき，( $\beta$  が小さくなるので) 検出力は大きくなる

# 回帰分析

- ▶ 相関関係を考える際には、2つの変数の間に区別はなかった
- ▶ 回帰分析では、一方の変数が他方の変数を説明するという**因果関係を想定**する (仮定する)
- ▶ **回帰 (regression)**: 「興味のある (従属) 変数  $y$  の値・振る舞いを, 他の (独立) 変数  $x$  を用いて予測すること」
- ▶ 線形回帰: 従属変数の平均値の変化を, 独立変数の線形関数で要約する方法
  - ▶ 一般化線形モデル (generalized linear model, GLM)
- ▶ 単回帰分析: 1つの独立変数を含む回帰分析
- ▶ 重回帰分析: 複数の独立変数を含む回帰分析

# 従属変数と独立変数

- ▶ 従属変数 (dependent variable)  $y$ : 説明される変数 (結果)
  - ▶ 応答変数, 被説明変数, 結果変数
- ▶ 独立変数 (independent variable)  $x$ : 説明する変数 (原因)
  - ▶ 説明変数, 予測変数
- ▶ コントロール (調整/統制) 変数
  - ▶ 注目する独立変数以外で, 結果に影響を与える変数
  - ▶ 統計学的/数学的には単に  $x$
- ▶ 従属変数と独立変数は, 分析者が (なんらかの理論的根拠に基づき) 想定・仮定する
  - ▶ 「賃金 (従属変数  $y$ ) を教育年数 (独立変数  $x$ ) に回帰する」
  - ▶ 回帰分析のみによって, 因果関係を示すことができる訳ではない

# 回帰式と誤差項

$$\begin{aligned}y_i &= \mathbf{x}_i \boldsymbol{\beta} + \epsilon_i \\ &= \beta_1 x_{i1} + \cdots + \beta_k x_{ik} + \epsilon_i, \text{ for } i = 1, \dots, n,\end{aligned}\tag{1}$$

- ▶ 回帰係数  $\beta_k$ : 独立変数と従属変数の関係を示す
  - ▶ 第一種過誤・第二種過誤の  $\beta$  と混同しないように
- ▶ 誤差項 (error term):  $\epsilon_i = y_i - \hat{y}_i$
- ▶ 従属変数の変化は、独立変数の変化によって完全に説明できる訳ではない
- ▶ 観察不可能な (モデルに取り込まれていない) 「偶然」 のような変数の影響も受ける

# 回帰分析の記法

## 線形回帰モデルの数式表現

- ▶ 線形回帰モデルは、次のように表現できる

$$\begin{aligned}y_i &= \mathbf{x}_i \boldsymbol{\beta} + \epsilon_i \\ &= \beta_1 x_{i1} + \cdots + \beta_k x_{ik} + \epsilon_i, \text{ for } i = 1, \dots, n,\end{aligned}\tag{2}$$

- ▶ また、次のようにも表現できる

$$y_i \sim \mathcal{N}(\mathbf{x}_i \boldsymbol{\beta}, \sigma^2), \text{ for } i = 1, \dots, n.\tag{3}$$

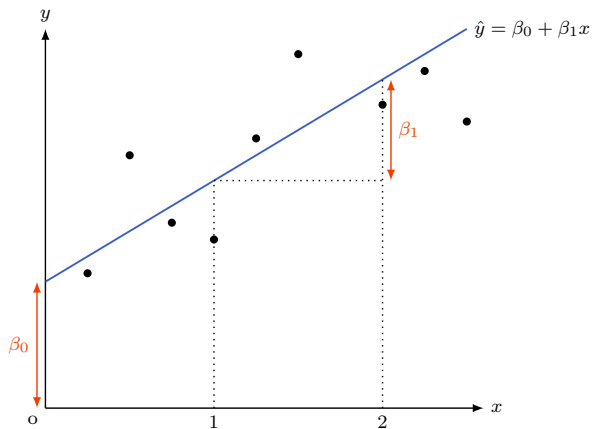
- ▶ あるいは,

$$\mathbf{y} \sim \mathcal{N}(\mathbf{x} \boldsymbol{\beta}, \sigma^2 \mathbf{I}),\tag{4}$$

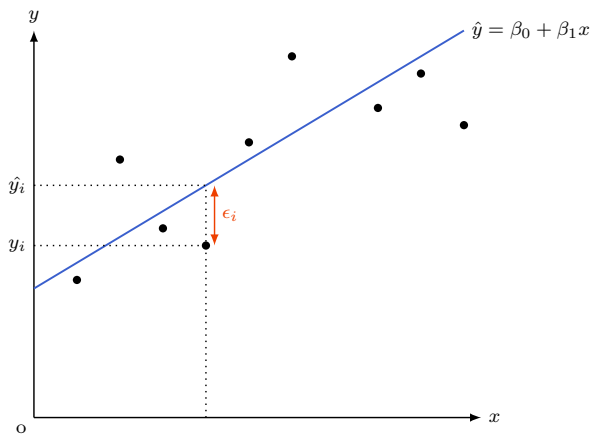
where  $\mathbf{I}$  is an  $n \times n$  identity matrix (単位行列).



# 回帰分析の記法



# 回帰分析の記法



# 推定値と推定法

- ▶ 推定量 (estimator): 推定の結果を, 標本  $\{(y_i, x_i) : i = 1, \dots, n\}$  の関数として表現したもの

- ▶ OLS 推定量  $\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2}$  (次スライド)

- ▶ 推定値 (estimate): (推定量に) 実際に標本を代入して計算した数値のこと

# 最小二乗法

- ▶ どのようにして (上の図の) 直線を引けば (回帰直線を推定すれば) よい?
- ▶ **最小二乗法 (ordinary least squares, OLS)**: 予測誤差 (残差 residual) の平方和 (二乗の和) を最小化することで, 事で, 観測値と予測値のズレを最小化する方法
- ▶ 予測誤差は観測値 (点) と予測値 (直線) のズレなので, そのズレが最も小さくなる回帰直線を見つける
- ▶ 平均二乗誤差 (mean squared error, MSE):

$$\text{MSE} = \frac{1}{n}(\epsilon_1^2, \dots, \epsilon_n^2) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

# 最小二乗法

- ▶ 回帰式  $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$  を考える
- ▶ 予測誤差  $\epsilon_i$  について,

$$\frac{1}{n} \epsilon_i^2 = \frac{1}{n} (y_i - \hat{y}_i)^2 = \frac{1}{n} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \quad (5)$$

$$\frac{1}{n} \sum_{i=1}^n \epsilon_i^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \quad (6)$$

- ▶ これを最小化する回帰係数の値を、**最小二乗推定量 (OLS 推定量)** と呼ぶ
- ▶ MSE を最小化するには、回帰係数で偏微分して「=0」と置いた連立方程式を解けばよい (最適化の1階条件)

$$\frac{\partial \text{MSE}}{\partial \hat{\beta}_0} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \quad (7)$$

$$\frac{\partial \text{MSE}}{\partial \hat{\beta}_1} = -\frac{2}{n} \sum_{i=1}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \quad (8)$$

- ▶ この連立方程式を解けば、OLS 推定量が得られる

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \quad (9)$$

# 最小二乗法

▶ 行列を使えば、OLS 推定量は  $\hat{\beta} = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{y}$  と書ける

▶ OLS 推定量と相関係数

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2}$$

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)}{\sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y}_n)^2}}$$

▶ **違いは分母**：独立変数の分散 (回帰係数) か両者の標準偏差 (相関係数) か  
▶ 「回帰分析のみで実証できること」は因果関係ではない (とは限らない)

▶ 回帰係数にせよ相関係数にせよ、**2つの変数の「変化 (variation) の関係」**を見るもの

▶ 「変化の関係」をみるのだから、両者に十分な変化がなければならない！  
▶ 副読本森田本の例：市区町村の失業者数が、1つだけ500でそれ以外は1,000の場合？

# 回帰係数の検定 (詳しくは次回)

- ▶  $\frac{\beta_k - \hat{\beta}_k}{SE}$  は自由度  $n - k$  の  $t$  分布に従う
  - ▶  $n$  は標本サイズ,  $k$  は切片  $\beta_0$  を含む回帰係数の数
- ▶ 自由度  $n - k$  の  $t$  分布を使った  $t$  検定を行なう
  - ▶ 帰無仮説  $H_0 : \beta_k = 0$
- ▶ 前回扱った平均値の検定と同様に, 仮説検定ができる
- ▶ 有意水準  $\alpha = 0.05$  なら,  $\pm 1.96SE$  が目安になる
- ▶ ( $H_0$  は  $\beta_k = 0$  なので, ) おおよそ  $\hat{\beta}_k \pm 1.96SE$  (あるいは約  $2SE$ ) の範囲に "0" が含まれなければ, その回帰係数の効果は統計的に有意 (statistically significant)
- ▶  $\hat{\beta}_k$  の 95%信頼区間:  $[\hat{\beta}_k - t_{0.025, n-k}SE, \hat{\beta}_k + t_{0.025, n-k}SE]$

# 回帰係数と限界効果

## 線形回帰における (偏) 回帰係数の解釈

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon \quad (10)$$

という回帰式を考える

- ▶ この回帰式において、独立変数  $x_k$  の偏回帰係数  $\beta_k$  は、他の独立変数の値を一定に保ったとき、 $x_k$  が 1 (単位) 増加したときの従属変数 (の予測値)  $\hat{y}$  の増加量
- ▶ 重回帰分析における「偏回帰係数」の「偏」は偏微分の「偏 (partial)」と同じ意味
- ▶ 「 $x_k$  が 1 (単位) 増加したときの従属変数 (の予測値)  $\hat{y}$  の増加量 (影響の大きさ)」のことを、**限界効果 (marginal effect)** という
  - ▶ 回帰係数と同じに思えるが、より高度な一般化線形モデル (本講義では扱わない) の理解や、後述の交互作用項を含むモデルの理解には限界効果が非常に重要



# 回帰係数と限界効果

## 線形回帰における (偏) 回帰係数の解釈

- ▶ 解析では、統計的有意性 (statistical significance) と同程度以上に、実質的有意性 (substantial significance; 効果の大小) が重要
- ▶ 実質的有意性を解釈する上では、限界効果が鍵になる
  - ▶ 一般化線形モデル (identity link 以外) では線形回帰のような直感的な解釈が困難なため、限界効果の計算・作図がより重要になる

## 限界効果

- ▶ 「偏」の字から直感的にわかる通り、 $x_k$  の限界効果は (10) 式を  $x_k$  で偏微分すれば求められる。

$$\frac{\partial \hat{y}}{\partial x_k} = \beta_k \quad (11)$$

# 回帰係数と限界効果

## 限界効果, 実質的有意性の重要さ

- ▶ よく聞く表現: 「この変数はトウケイテキニユウイだった. 大事だ」
- ▶ 気をつけるべき点: 統計的に有意だからといって, 実質的に有意とは限らない (「統計的に有意でも, 大して影響のない」変数もある)

## 「他の独立変数が一定のとき」の意味をよく考える

- ▶ 常に「一定に保つ」ことができる訳ではない
  - ▶ 例: 二乗項や交互作用項/交差項 (interaction term) が含まれる回帰式 (次スライド)
  - ▶ 交互作用 (interaction): ある独立変数の効果が, 他の独立変数の値によって異なること
- ▶ いずれの場合でも, ある独立変数  $x_k$  「だけ」を動かして, その限界効果を提示することは不可能 (掛け算してるから!)
- ▶ さらに, 観察データ (observational data) では, 往々にして独立変数間に相関がある (Hanmer & Ozan Kalkan, 2013)

# 限界効果の計算：Interaction term がある場合

## 交互作用項とは

- ▶  $x_1 \times x_2$  のような、「かけ算」で表現される (独立) 変数
- ▶  $x_1 \times x_2$  を “product term,”  $x_1$  と  $x_2$  を “constitutive term(s)” (main effect) と呼ぶ
  - ▶ 理論的議論：Berry et al. (2010, 2012); Brambor et al. (2006); Clark et al. (2006); Kam & Franzese (2007); Rainey (2015)

## 交互作用項の解釈と限界効果

- ▶ よくある間違い (1)： $x_1$  (e.g., 投薬の有無) や  $x_2$  (e.g., 性別) の交互作用項  $x_1 \times x_2$  を回帰式に投入したから、 $x_1$  と  $x_2$  をバラバラに回帰式に投入する必要はない！
  - ▶ 誤： $y = \beta_0 + \beta_3 x_1 x_2 + \epsilon$
  - ▶ 正： $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon$
- ▶ よくある間違い (2)：交互作用項の回帰係数  $\beta_3$  が有意だった！だから、 $x_1$  や  $x_2$  単独ではなく、その組み合わせが重要なんだ！（もうちょっと作業が必要）

## 限界効果の計算：Interaction term がある場合

- ▶ 交互作用項がある場合， product term も constitutive term も， 単一の回帰係数のみでは解釈できない
- ▶ たとえば，  $y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_1x_2 + \epsilon$  の回帰式を考える
- ▶  $x_1$  の限界効果を計算するには回帰式を  $x_1$  で偏微分すればよいので，

$$\frac{\partial y}{\partial x_1} = \beta_1 + \beta_3x_2 \quad (12)$$

- ▶ すなわち，  $x_1$  の限界効果は  $\beta_1$  ではなく，  $\beta_1 + \beta_3x_2$
- ▶  $x_2$  の限界効果も同様に求められる

$$\frac{\partial y}{\partial x_2} = \beta_2 + \beta_3x_1 \quad (13)$$

- ▶  $x_2$  の限界効果は  $\beta_2$  ではなく，  $\beta_2 + \beta_3x_1$
- ▶  $y = \beta_0 + \beta_1x_1 + \beta_2x_1^2 + \epsilon$  のような， 二乗項 (“2” である必要はない。 3 以上でもよい) を含む回帰式についても同様

# 限界効果の計算：Interaction term がある場合

## 含意

- (1)  $x_1$  と  $x_2$  の「掛け算」が入るのだから、 $x_1$  を「だけ」を変化させ、「他の変数を一定に」することはできない
- (2)  $x_2 = 0$  (かつ/あるいは  $\beta_3 = 0$ ) でない限り、 $x_1$  の限界効果は  $x_2$  の水準に依存する
  - (2')  $x_2$  の水準を固定しなければ  $x_1$  の限界効果を計算できないが、 $x_2$  の水準によって  $x_1$  の限界効果が変化する
    - ▶ 実際問題、「様々な  $x_2$  の水準における  $x_1$  の限界効果を示す図」がなければ解釈できない!
    - ▶ 典型的には、 $x_2$  を  $x$  軸に、 $x_2$  がある値のときの  $x_1$  の限界効果を  $y$  軸にとったグラフを使う
- (3) 一般化線形モデル (identity link 以外) の場合，“compression effect”が生じるため、さらに大変 (「他の全ての独立変数」の値に依存する; Berry et al., 2010; Rainey, 2015)

# Rによる回帰分析とシミュレーション

## 現実のデータを用いた解析

- ▶ カナダのセンサスデータ (1971) Prestige
- ▶ car パッケージから読み込めるサンプルデータ

## シミュレーション

- ▶ **疑問**：モデルを正しく特定できている (かつ種々の仮定も満たされている) として、回帰分析は「真の」影響の大きさとその不確実性を、本当に教えてくれるのだろうか
- ▶ 検討：「真のモデル」が分かる仮想データを大量につくり、解析
- ▶ 手順：「真のモデル (データ生成プロセス data generating process) は同じだが、微妙に異なる」 $n = 1,000$  の仮想データを 3,000 個つくり、3,000 回の回帰分析を回し、その結果を俯瞰する
- ▶ **問題**：「ポチポチ」してたら日が暮れる (どころじゃ済まない. 3,000 どころか 3,000,000 回かもしれない)

# Rによる回帰分析とシミュレーション

## Rによる実行

- ▶ サンプルデータ Prestige の回帰分析
- ▶ シミュレーションによる回帰分析の理解
  - ▶ 3,000 回だろうが 3,000,000 回だろうが一箇所変えるだけ
- ▶ R コード: [http://cfes-project.eco.u-toyama.ac.jp/wp-content/uploads/5\\_regression\\_fall2017.html](http://cfes-project.eco.u-toyama.ac.jp/wp-content/uploads/5_regression_fall2017.html)
  - ▶ 「演習資料」ページの「5. データの解析」にある「回帰分析のシミュレーション」

# 次回講義と課題

## ▶ 次回講義

- ▶ 今日の続き (回帰分析)
- ▶ データの自動取得 (API, ウェブスクレイピング)

## ▶ 文献と課題

**必須** 星野・田中 『R による実証分析』 第 4–6 章 (教科書)

推奨 Gelman & Hill. *Data analysis*. Chaps. 3–4 (教科書)

推奨 浅野・矢内 『Stata による計量政治学』 第 7–12 章 (副読本)

推奨 石田ほか 『R によるスクレイピング入門』 (副読本)

推奨 森田 『実証分析入門』 第 4–8 章 (副読本)

**課題** (1) 講義資料「R 言語の基礎, オブジェクトとその要素へのアクセス」と「R によるデータの読み込みと書き出し」を RStudio で練習しておくこと. (2) シミュレーションを再度行なってみること



# 補足・再掲：R スクリプト作成と実行

RStudio で次の作業をする

- 1 RStudio を開き、`command + shift + N (m)/control + shift + N (w)` して、新しい R スクリプトを開く (作成する)
- 2 `command + S (m)/control + S (w)` して、(真っ白な) R スクリプトを分かりやすい場所に保存
  - ▶ **ファイル名は半角文字のみ**にすること！
  - ▶ 例：`clt_simulation.R`
- 3 演習の際は、ブラウザで開いているテンプレートの中身を、今作成した R スクリプトにコピーする。自分で全く新しいスクリプトを作る際は、R の関数や処理を上から (実行したい順番に) 書いていく
- 4 `command + S (m)/control + S (w)` して、R スクリプトを (上書き) 保存
- 5 実行したい部分を選択し、`command + enter (m)/control + enter (w)` で実行
- 6 (エラーが出れば、コードに間違いがないか確認する)

# 補足：両側検定と片側検定

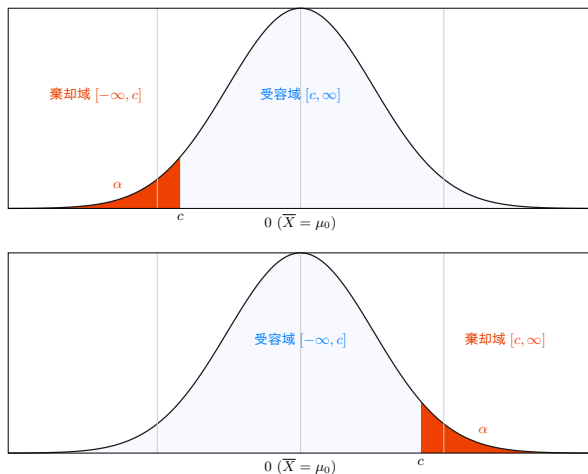
## 両側検定と片側検定

- ▶ 両側検定 (two-sided/-tailed test)：分布の両側に棄却域を設定する検定
- ▶ 片側検定 (one-sided/-tailed test)：分布の片側に棄却域を設定する検定

## 対立仮説の性質

- ▶ 帰無仮説  $H_0$  が「 $\mu = \mu_0$ 」のとき、対立仮説  $H_1$  には3つのバリエーションがあり得る
  - 1  $H_{1a}$  「 $\mu \neq \mu_0$ 」
  - 2  $H_{1b}$  「 $\mu < \mu_0$ 」
  - 3  $H_{1c}$  「 $\mu > \mu_0$ 」
- ▶ 対立仮説の性質に応じ、両側検定 ( $H_{1a}$ ) か片側検定 ( $H_{1b}$ ,  $H_{1c}$ ) かが決まる
- ▶ 「 $\mu < \mu_0$ 」 ( $H_{1b}$ ) や 「 $\mu > \mu_0$ 」 ( $H_{1c}$ ) を設定する特段の理由がなければ、 $H_{1a}$  と両側検定を用いる (一般に、片側検定の方が  $H_0$  を棄却しやすい)
- ▶ 副読本・浅野矢内本 pp.132-134 を参照のこと

## 補足：両側検定と片側検定



- ▶  $H_{1b}$  「 $\mu < \mu_0$ 」 (上図) と  $H_{1c}$  「 $\mu > \mu_0$ 」 (下図) を想定した片側検定の棄却域と受容域

- Berry, William D; Jacqueline H R DeMeritt & Justin Esarey (2010) Testing for Interaction Effects in Binary Logit and Probit Models: Is an Interaction Term Essential? *American Journal of Political Science* 54(1): 248–266.
- Berry, William D; Matt Golder & Daniel Milton (2012) Improving Tests of Theories Positing Interaction. *Journal of Politics* 74(3): 653–671.
- Brambor, Thomas; William Clark & Matt Golder (2006) Understanding interaction models: Improving empirical analyses. *Political Analysis* 14(1): 63–82.
- Clark, William R; Michael J Gilligan & Matt Golder (2006) A simple multivariate test for asymmetric hypotheses. *Political Analysis* 14(3): 311–331.
- Hanmer, Michael J & Kerem Ozan Kalkan (2013) Behind the Curve: Clarifying the Best Approach to Calculating Predicted Probabilities and Marginal Effects from Limited Dependent Variable Models. *American Journal of Political Science* 57(1): 263–277.
- Kam, Cindy D & Robert J Franzese (2007) *Modeling and Interpreting Interactive Hypotheses in Regression Analysis: A Refresher and Some Practical Advice*. Ann Arbor, MI: University of Michigan Press.
- Rainey, Carlisle (2015) Compression and Conditional Effects: A Product Term Is Essential When Using Logistic Regression to Test for Interaction. *Political Science Research and Methods* forth.