

# Лекция 11. Бутстраппинг

Буре В.М., Грауэр Л.В.

ШАД

Санкт-Петербург, 2013

## Идея бутстраппинга

Рассмотрим случайную величину  $\xi$  с неизвестной функцией распределения  $F(x)$ .

$(x_1, \dots, x_n)$  — выборка объема  $n$  из генеральной совокупности  $\xi$ .

Истинное распределение данных можно хорошо приблизить эмпирическим.

Эмпирическая функция распределения

$$F^*(x) = \frac{1}{n} \sum_{i=1}^n I_{[x_i \leq x]}$$

равномерно почти наверное стремится к  $F(x)$  при  $n \rightarrow \infty$  по лемме Гливленко-Кантелли, где  $I_{[\cdot]}$  — индикатор-функция.

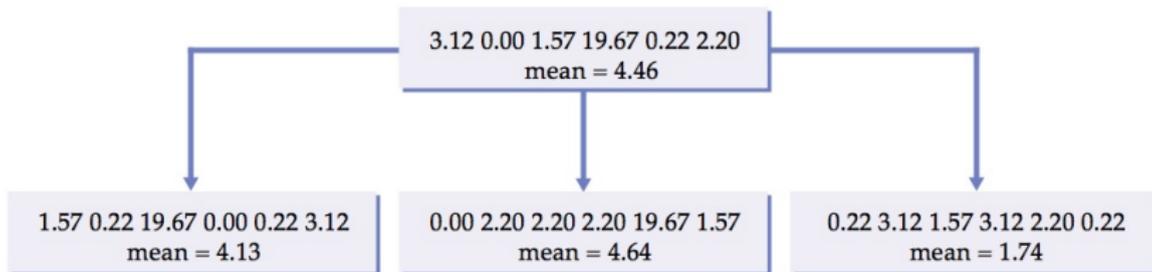
Идея бутстрепа по Б. Эфрону (1989):

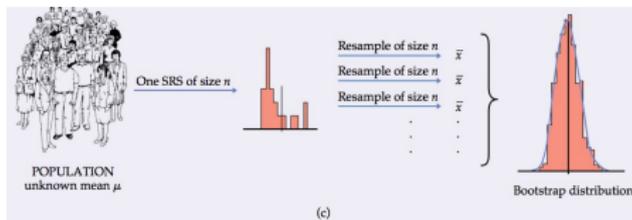
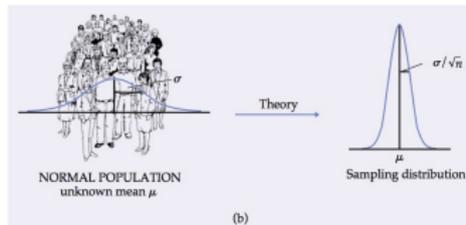
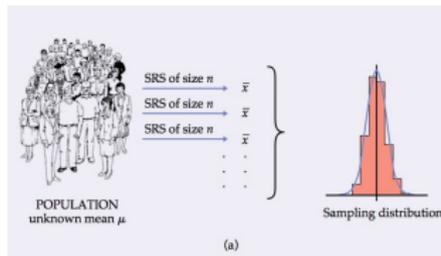
трансформировать приближение для распределения данных в приближенное распределение статистик, используя метод статистических испытаний Монте-Карло, многократно извлекая повторные выборки из эмпирического распределения.

Пусть  $B$  — число бутстрапповских выборок. Для  $b = 1, 2, \dots, B$  построим бутстраповскую выборку

$$\{x_1^*, x_2^*, \dots, x_n^*\}_b \quad (1)$$

А именно: берется конечная совокупность из  $n$  членов исходной выборки  $x_1, x_2, \dots, x_n$ , откуда на каждом шаге из  $n$  последовательных итерации с помощью датчика случайных чисел, равномерно распределенных на интервале  $[1, n]$ , «вытягивается» произвольный элемент  $x_k$ , который снова «возвращается» в исходную выборку.





Пробутстрапим статистику  $\hat{\theta} = \hat{\theta}(\{x_1, x_2, \dots, x_n\})$ .

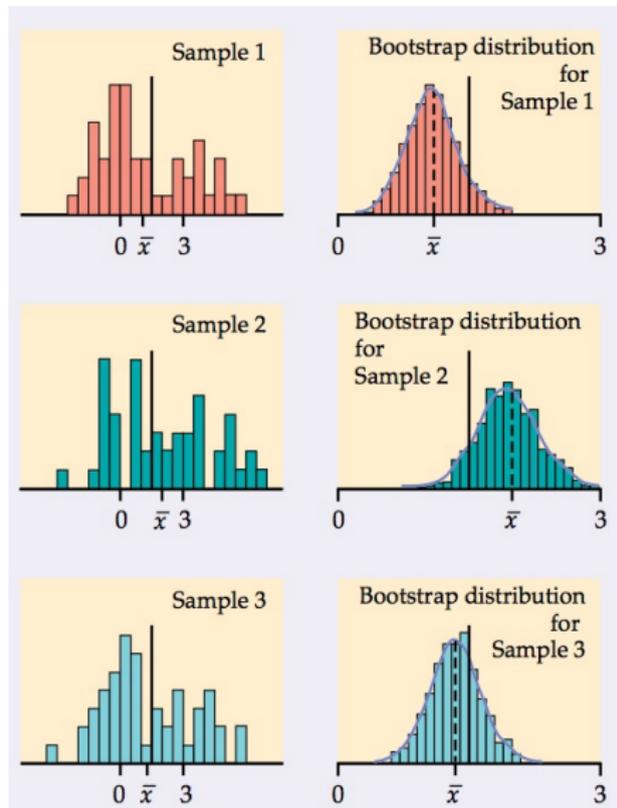
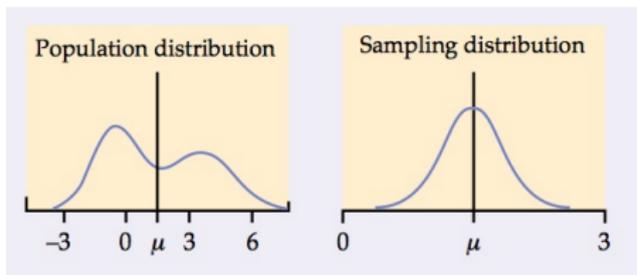
Для каждой выборки (1) вычислим бутстраповскую статистику

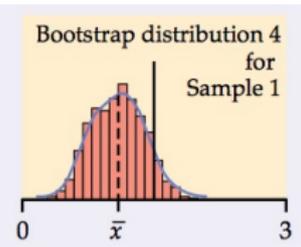
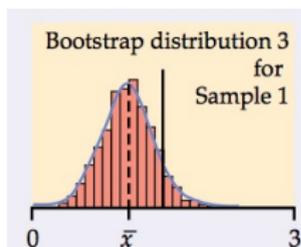
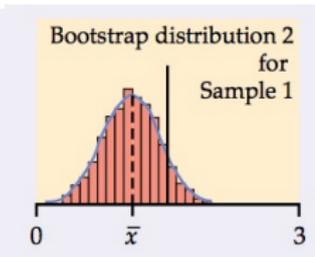
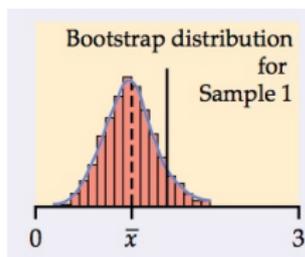
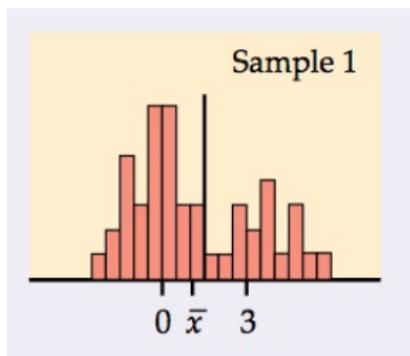
$$\hat{\theta}_b^* = \hat{\theta}(\{x_1^*, x_2^*, \dots, x_n^*\}_b), \quad b = 1, \dots, B. \quad (2)$$

Полученный набор статистик  $\hat{\theta}_1^*, \hat{\theta}_2^*, \dots, \hat{\theta}_B^*$  с приписанием каждой веса  $1/B$  составляет приближенное бутстраповское распределение статистики  $\hat{\theta}$  с характеристиками

$$E(\hat{\theta}^*) = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b^*$$

$$SE_{boot}(\hat{\theta}^*) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\hat{\theta}_b^* - E(\hat{\theta}^*))^2}$$





## Рецентрирование

Когда речь идет о разностях или расстояниях между выборочными и популяционными объектами, необходимо рецентрирование. Правильным бутстраповским аналогом разности  $\hat{\theta} - \theta$  является  $\hat{\theta}^* - \hat{\theta}$ .

### Пример 1

Бутстрапируя  $t$ -статистику при нулевой гипотезе  $H_0 : \theta = 0$

$$t = \frac{\hat{\theta}}{se(\hat{\theta})},$$

следует использовать бутстрап аналог

$$t^* = \frac{\hat{\theta}^* - \hat{\theta}}{se_{boot}(\hat{\theta}^*)}$$

## Пример 2

Бутстраповским аналогом статистики отношения правдоподобия

$$LR = 2 \left( \max_{\theta} \ln L(\theta) - \max_{\theta: g(\theta)=0} \ln L(\theta) \right)$$

будет

$$LR^* = 2 \left( \max_{\theta} \ln L(\theta) - \max_{\theta: g(\theta)=\hat{\theta}} \ln L(\theta) \right)$$

Рецентрировать  $se_{boot}(\hat{\theta}^*)$  нет необходимости так, как это нормирующий множитель.

## Бутстраповская корректировка смещения

Бутстрап позволяет скорректировать смещение, связанное с конечностью выборки.

Пусть  $\hat{\theta}$  — состоятельная смещенная оценка неизвестного параметра  $\theta$  ( $E(\hat{\theta}) \neq \theta$ ).

Смещение оценки  $\hat{\theta}$  равно  $C(\hat{\theta}) = E(\hat{\theta}) - \theta$ .

Вычислив бутстраповский аналог этого смещения

$$C^*(\hat{\theta}) = E(\hat{\theta}^*) - \hat{\theta},$$

можно скорректировать исходную статистику

$$\hat{\theta}_{BC} = \hat{\theta} - C^*(\hat{\theta}) = 2\hat{\theta} - E(\hat{\theta}^*) = 2\hat{\theta} - \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b^*.$$

Бутстрап не способен справиться с асимптотическим смещением.

## Пример 3

Пусть  $\mu$  — неизвестное математическое ожидание случайной величины  $\xi$ . Рассмотрим статистику  $\mu^2$ . В качестве оценки  $\mu$  возьмем  $\bar{x}_n$ , в качестве оценки  $\mu^2$  —  $\bar{x}_n^2$ .

Бутстраповский аналог смещения  $C^*(\bar{x}_n) = E^*(\bar{x}_n^*) - \bar{x}_n = 0$ , так как  $E^*(\bar{x}_n^*) = \bar{x}_n$ .

Смещение  $\bar{x}_n^2$  равно

$$C(\bar{x}_n^2) = E(\bar{x}_n^2) - \mu^2 = D(\bar{x}_n) = \frac{1}{n}D\xi$$

Бутстраповский аналог смещения в данном случае есть

$$C^*(\bar{x}_n^2) = \frac{1}{n} \left( \frac{1}{n} \sum_{i=1}^n x_i - \bar{x}_n^2 \right).$$

Тогда оценкой  $\mu^2$ , скорректированной на смещение, будет

$$\bar{x}_n^2 - C^*(\bar{x}_n^2) = \frac{n+1}{n} \bar{x}_n^2 - \frac{1}{n^2} \sum_{i=1}^n x_i.$$

## Доверительные интервалы

Пусть  $\hat{\theta}$  — состоятельная оценка параметра  $\theta$ . Построим доверительные интервалы параметра  $\theta$ .

### Эфронов доверительный интервал

Используем бутстраповское распределение статистики  $\hat{\theta}$  (2) для построения доверительного интервала.

В качестве квантилей  $q_{\alpha_1}^*$ ,  $q_{1-\alpha_2}^*$  возьмем порядковые статистики  $\hat{\theta}_{([B\alpha_1])}^*$  и  $\hat{\theta}_{([B(1-\alpha_2)]+1)}^*$ .

Эфронов доверительный интервал с покрытием  $1 - \alpha_1 - \alpha_2$  равен

$$CI_{efr} = \left( \hat{\theta}_{([B\alpha_1])}^*, \hat{\theta}_{([B(1-\alpha_2)]+1)}^* \right).$$

Так как не используется рецентрирование, в данном случае смещение исходной выборки лишь усиливается. Применим для симметричных распределений или в случае малых смещений.

## Квантильный доверительный интервал (рецентрированный)

Пробутстрапим статистику  $\hat{\theta} - \theta$ , бутстраповский аналог которой  $\hat{\theta}^* - \hat{\theta}$ .

В качестве квантилей  $q_{\alpha_1}^{*\%}$ ,  $q_{1-\alpha_2}^{*\%}$  возьмем соответствующие порядковые статистики порядков  $[B\alpha_1]$  и  $[B(1 - \alpha_2)] + 1$  и построим квантильный доверительный интервал с покрытием  $1 - \alpha_1 - \alpha_2$

$$CI_{\%} = \left( \hat{\theta} - q_{1-\alpha_2}^{*\%}, \hat{\theta} - q_{\alpha_1}^{*\%} \right).$$

## T-квантильный доверительный интервал

Пробутстрапим статистику  $\frac{\hat{\theta} - \theta}{se(\hat{\theta})}$ , бутстраповский аналог которой  $\frac{\hat{\theta}^* - \hat{\theta}}{se^*(\hat{\theta})}$ . Получим бутстраповские квантили  $q_{\alpha_1}^{*\%t}$ ,  $q_{1-\alpha_2}^{*\%t}$  и построим t-квантильный доверительный интервал с покрытием  $1 - \alpha_1 - \alpha_2$

$$CI_{\%t} = \left( \hat{\theta} - se(\hat{\theta})q_{1-\alpha_2}^{*\%t}, \hat{\theta} - se(\hat{\theta})q_{\alpha_1}^{*\%t} \right).$$

Если  $\alpha_1 = \alpha_2 = \alpha/2$ , то можно бутстрапировать не центрированную оценку или  $t$ -статистику, а их модули  $|\hat{\theta} - \theta|$  и  $|\hat{\theta} - \theta|/se(\hat{\theta})$ , соответственно.

**Симметричный квантильный доверительный интервал**

$$CI_{|\%|} = \left( \hat{\theta} - q_{1-\alpha/2}^{*|\%|}, \hat{\theta} + q_{1-\alpha/2}^{*|\%|} \right).$$

**Симметричный  $t$ -квантильный доверительный интервал**

$$CI_{|\%|t} = \left( \hat{\theta} - se(\hat{\theta})q_{1-\alpha/2}^{*|\%|t}, \hat{\theta} + se(\hat{\theta})q_{1-\alpha/2}^{*|\%|t} \right).$$

# Асимптотическое рафинирование

## Определение 1

Пусть  $\hat{\varphi}$  — некоторая статистика, истинное распределение которой  $F_{\hat{\varphi}}(x)$ ,  $F_{\hat{\varphi}}^*(x)$  — бутстраповское распределение этой статистики. Говорят, что с помощью бутстрапа достигается асимптотическое рафинирование, если ошибка аппроксимации истинного распределения  $F_{\hat{\varphi}}(x)$  бутстраповским  $F_{\hat{\varphi}}^*(x)$  большего порядка малости, чем ошибка аппроксимации асимптотическим распределением при стремлении объема выборки к бесконечности.

Рассмотрим асимптотически нормальную статистику

$$\hat{\varphi} = \frac{\hat{\theta} - \theta}{se(\hat{\theta})} \xrightarrow{d} N(0, 1).$$

Разложения Эджворта истинного и бутстраповского распределений вокруг асимптотического выглядят следующим образом

$$F_{\hat{\varphi}}(x) = \Phi(x) + \frac{h_1(x, F)}{\sqrt{n}} + \frac{h_2(x, F)}{n} + O\left(\frac{1}{n\sqrt{n}}\right),$$

$$F_{\hat{\varphi}^*}(x) = \Phi(x) + \frac{h_1(x, F^*)}{\sqrt{n}} + \frac{h_2(x, F^*)}{n} + O\left(\frac{1}{n\sqrt{n}}\right),$$

где  $h_1(x, F)$  — четная по  $x$ , непрерывная по  $F$  функция,  $h_2(x, F)$  — нечетная по  $x$ , непрерывная по  $F$  функция

Ошибки аппроксимации точного распределения асимптотическим и бутстраповским, соответственно, равны

$$\Phi(x) - F_{\hat{\varphi}}(x) = \frac{h_1(x, F)}{\sqrt{n}} + O\left(\frac{1}{\sqrt{n}}\right) = O\left(\frac{1}{\sqrt{n}}\right),$$

$$F_{\hat{\varphi}}^*(x) - F_{\hat{\varphi}}(x) = \frac{h_1(x, F^*) - h_1(x, F)}{\sqrt{n}} + O\left(\frac{1}{n}\right) = O\left(\frac{1}{n}\right),$$

Сравнивая порядки ошибок аппроксимации, можно сделать вывод, что в данном случае использование бутстрапа приводит к асимптотическому рафинированию.

На практике это означает, что в достаточно больших выборках ошибка бутстраповского приближения, намного меньше, чем ошибка асимптотического приближения.

Аналогично можно показать, что  
если статистика

$$\hat{\varphi} = \sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, D_\theta),$$

то бутстрап не приводит к асимптотическому рафинированию

$$\Phi(x, D_\theta) - F_{\hat{\varphi}}(x) = O\left(\frac{1}{\sqrt{n}}\right), \quad F_{\hat{\varphi}}^*(x) - F_{\hat{\varphi}}(x) = O\left(\frac{1}{\sqrt{n}}\right);$$

если статистика

$$\hat{\varphi} = \frac{|\hat{\theta} - \theta|}{se(\hat{\theta})} \xrightarrow{d} |N(0, 1)|,$$

то бутстрап приводит к асимптотическому рафинированию, причем  
ошибка большего порядка малости

$$\Phi(x) - F_{\hat{\varphi}}(x) = O\left(\frac{1}{n}\right), \quad F_{\hat{\varphi}}^*(x) - F_{\hat{\varphi}}(x) = O\left(\frac{1}{n\sqrt{n}}\right).$$

# Проверка статистических гипотез

$$H_0 : \theta = \theta_0$$

Для проверки гипотезы  $H_0 : \theta = \theta_0$  с помощью бутстрапирования в случае

- двусторонней альтернативной гипотезы  $H_1 : \theta \neq \theta_0$ , пробустранировав модуль  $t$ -статистики  $|\hat{\theta} - \theta|/se(\hat{\theta})$ , получим бутстраповское распределение  $|\hat{\theta}^* - \hat{\theta}|/se^*(\hat{\theta})$  и ее квантиль  $q_{1-\alpha/2}^{*|\%|t}$ . Гипотеза  $H_0$  отвергается на уровне значимости  $\alpha$ , если  $|\hat{\theta} - \theta_0|/se(\hat{\theta}) > q_{1-\alpha/2}^{*|\%|t}$ .
- односторонней альтернативной гипотезы  $H_1 : \theta > \theta_0$ , пробустранировав  $t$ -статистику  $(\hat{\theta} - \theta)/se(\hat{\theta})$ , получим бутстраповское распределение  $(\hat{\theta}^* - \hat{\theta})/se^*(\hat{\theta})$  и ее квантиль  $q_{1-\alpha}^{*\%t}$ . Гипотеза  $H_0$  отвергается на уровне значимости  $\alpha$ , если  $(\hat{\theta} - \theta_0)/se(\hat{\theta}) > q_{1-\alpha}^{*\%t}$ .

## Проверка статистических гипотез

$$H_0 : \theta_1 = \theta_2$$

Рассмотрим выборки  $(x_1, \dots, x_n)$  и  $(y_1, \dots, y_m)$  из генеральных совокупностей  $\xi$  и  $\eta$ . Проверим гипотезу о равенстве параметров двух распределений  $H_0 : \theta_1 = \theta_2$  при двусторонней альтернативной гипотезе  $H_1 : \theta_1 \neq \theta_2$ .

### Бутстрапирование

Пробутстрапируем статистики  $\hat{\theta}_1(\{x_1, \dots, x_n\})$  и  $\hat{\theta}_2(\{y_1, \dots, y_m\})$ . Для каждого  $b = 1, \dots, B$  получим значение бутстрап статистики

$$\hat{\theta}_1^*(\{x_1^*, \dots, x_n^*\}_b) - \hat{\theta}_2^*(\{y_1^*, \dots, y_m^*\}_b).$$

Построим распределение бутстрап статистики  $\hat{\theta}_1^* - \hat{\theta}_2^*$  и ее квантили  $q_{\alpha/2}^*$  и  $q_{1-\alpha/2}^*$ .

Гипотеза  $H_0$  отвергается на уровне значимости  $\alpha$ , если  $0 > q_{1-\alpha/2}^*$  либо  $0 < q_{\alpha/2}^*$ .

## Пермутейшн

Сформируем единую выборку

$$(x_1, \dots, x_n, y_1, \dots, y_m). \quad (3)$$

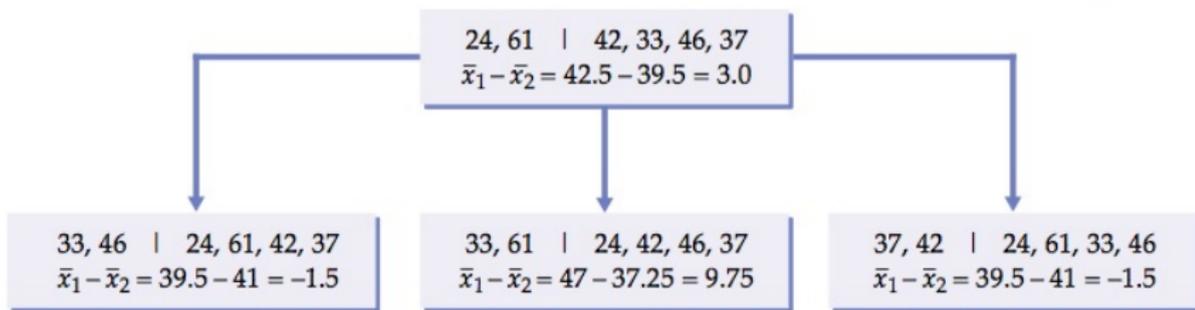
Случайным образом без возвращения извлечем  $n$  наблюдений из выборки (3)  $x_1^*, \dots, x_n^*$ , остальные элементы обобщенной выборки обозначим  $y_1^*, \dots, y_m^*$ . Получим значение статистики

$$\hat{\theta}_1^*({x_1^*, \dots, x_n^*}_b) - \hat{\theta}_2^*({y_1^*, \dots, y_m^*}_b).$$

Повторим данный ресэмплинг  $B$  раз.

Построим распределение пермутейшн статистики  $\hat{\theta}_1^* - \hat{\theta}_2^*$  и ее квантили  $q_{\alpha/2}^*$  и  $q_{1-\alpha/2}^*$ .

Гипотеза  $H_0$  отвергается на уровне значимости  $\alpha$ , если  $\hat{\theta}_1 - \hat{\theta}_2 > q_{1-\alpha/2}^*$  либо  $\hat{\theta}_1 - \hat{\theta}_2 < q_{\alpha/2}^*$ .



# Проверка статистических гипотез

$H_0 : m_1 = m_2$ , связанные выборки

Рассмотрим выборку  $(x_1, y_1), \dots, (x_n, y_n)$  из двумерной генеральной совокупности  $(\xi, \eta)$ . Проверим гипотезу о равенстве математических ожиданий двух распределений  $H_0 : m_1 = m_2$  при двусторонней альтернативной гипотезе  $H_1 : m_1 \neq m_2$ .

## Бутстрапирование

Перейдем к выборке  $(d_1, \dots, d_n)$ , где  $d_i = x_i - y_i$  и нулевой гипотезе  $H_0 : [m_1 - m_2 =] \theta = 0$ .

Построим на основе выборки  $(d_1, \dots, d_n)$   $t$ -статистику и пробутстрапируем модуль  $t$ -статистики  $|\hat{\theta}|/se(\hat{\theta})$ , получим бутстраповское распределение  $|\hat{\theta}^* - \hat{\theta}|/se^*(\hat{\theta})$  и ее квантиль  $q_{1-\alpha/2}^{*|\%|t}$ .

Гипотеза  $H_0$  отвергается на уровне значимости  $\alpha$ , если

$$|\hat{\theta}|/se(\hat{\theta}) > q_{1-\alpha/2}^{*|\%|t}$$

## Пермутейшн

Для каждого парного наблюдения  $(x_i, y_i)$  случайным образом переставим местами  $x_i$  и  $y_i$  внутри пары.

Получим ресэмплированную выборку  $(x_1^*, y_1^*), \dots, (x_n^*, y_n^*)$ .

Повторим действия  $B$  раз.

Построим распределение пермутейшн статистики  $\hat{\theta}_1^* - \hat{\theta}_2^*$  и ее квантили  $q_{\alpha/2}^*$  и  $q_{1-\alpha/2}^*$ .

Гипотеза  $H_0$  отвергается на уровне значимости  $\alpha$ , если  $\hat{\theta}_1 - \hat{\theta}_2 > q_{1-\alpha/2}^*$  либо  $\hat{\theta}_1 - \hat{\theta}_2 < q_{\alpha/2}^*$ .

# Ресэмплинг в регрессиях

Рассмотрим линейную регрессионную модель

$$y_i = X_i\beta + \varepsilon_i.$$

Существует два основных вида ресэмплинга в регрессионных моделях

- ресэмплинг парных наблюдений  $(y_1, x_1), (y_2, x_2), \dots, (y_n, x_n)$ , при котором бутсрапированные данные  $(y_1, x_1)^*, (y_2, x_2)^*, \dots, (y_n, x_n)^*$  извлекаются независимо с равными вероятностями  $1/n$ , и
- ресэмплинг остатков:  
Построив регрессионную модель  $X\hat{\beta}$ , извлекаем случайным образом  $\epsilon_i^*$  из набора центрированных стандартизированных остатков  $e_1, \dots, e_n$  и устанавливаем  $y_i^*$  равными

$$y_i^* = \hat{\beta}x_i + \epsilon_i^*, \quad i = 1, \dots, n.$$

# Литература

Эфрон Б. Нетрадиционные методы многомерного статистического анализа. М.: Финансы и статистика, 1988

Efron B., Tibshirani R.J. An introduction to the bootstrap. N.Y.: Chapman & Hall, 1993

Moore D. Bootstrap methods and permutation tests The practice of business statistics.

Эконометрический ликбез: бутстрап. Квантиль, 2007, №3, стр. 1-66

Chibara L., Hesterberg T. Mathematical statistics with resampling and R. Wiley