

Бутстреп

а) Моделирование Монте-Карло, бутстреп

Бутстреп – метод оценивания доверительных интервалов и подсчета ошибок

Пусть задана i.i.d. выборка $X_1, \dots, X_n \sim F$, $T_n = g(X_1, \dots, X_n)$ – статистика от нее

Необходимо оценить $\mathbb{V}_F(T_n)$

Пример. $T_n = \bar{X}_n$ – оценка среднего значения, тогда $\mathbb{V}_F(T_n) = \sigma^2/n$,
где $\sigma^2 = \int (x - \mu)^2 dF(x)$ и $\mu = \int x dF(x)$

Основная идея:

Шаг 1. Оценить $\mathbb{V}_F(T_n)$ с помощью $\mathbb{V}_{\hat{F}_n}(T_n)$

2 Шаг 2. Аппроксимировать значение $\mathbb{V}_{\hat{F}_n}(T_n)$ используя моделирование Монте-Карло

Пример. В случае $T_n = \overline{X}_n$ получаем, что $\mathbb{V}_{\hat{F}_n}(T_n) = \hat{\sigma}^2/n$, где $\hat{\sigma}^2 = n^{-1} \sum_{i=1}^n (X_i - \overline{X}_n)^2$, при этом шаг 2 не требуется. Шаг 2 требуется в более сложных случаях, в которых не удастся выписать явной формулы для $\mathbb{V}_{\hat{F}_n}(T_n)$.

Основная идея метода Монте-Карло состоит в следующем.

Пусть задана i.i.d. выборка Y_1, \dots, Y_B , распределение элементов которой равно G . По закону больших чисел при $B \rightarrow \infty$

$$\overline{Y}_n = \frac{1}{B} \sum_{j=1}^B Y_j \xrightarrow{P} \int y dG(y) = \mathbb{E}(Y)$$

3

В общем случае для функции h , у которой среднее значение относительно распределения G конечно, получаем, что

$$\frac{1}{B} \sum_{j=1}^B h(Y_j) \xrightarrow{P} \int h(y) dG(y) = \mathbb{E}(h(Y))$$

В частности,

$$\frac{1}{B} \sum_{j=1}^B (Y_j - \bar{Y})^2 = \frac{1}{B} \sum_{j=1}^B Y_j^2 - \left(\frac{1}{B} \sum_{j=1}^B Y_j \right)^2$$

4

$$\xrightarrow{P} \int y^2 dF(y) - \left(\int y dF(y) \right)^2 = \mathbb{V}(Y)$$

**=> можно использовать выборочную дисперсию сгенерированной выборки
для оценки дисперсии $\mathbb{V}(Y)$**

б) Оценка дисперсии на основе бутстрепа

Итак, значение $\mathbb{V}_{\hat{F}_n}(T_n)$ можно оценить с помощью моделирования

5 $\mathbb{V}_{\hat{F}_n}(T_n)$ = дисперсия T_n в случае, когда распределение элементов выборки равно $\hat{F}_n \Rightarrow$ надо моделировать выборки, в которых каждый из элементов имеет распределение \hat{F}_n , после чего оценить $\mathbb{V}_{\hat{F}_n}(T_n)$

$$\begin{array}{llll} F & \Longrightarrow & X_1, \dots, X_n & \Longrightarrow & T_n = g(X_1, \dots, X_n) \\ \hat{F}_n & \Longrightarrow & X_1^*, \dots, X_n^* & \Longrightarrow & T_n^* = g(X_1^*, \dots, X_n^*) \end{array}$$

$\hat{F}_n \Leftrightarrow$ наблюдения могут принимать одно из значений X_1, \dots, X_n , при этом вероятность получить любое из этих значений равна $1/n$

6

Таким образом, чтобы получить псевдовыборку $\hat{F}_n \Rightarrow X_1^*, \dots, X_n^*$, надо случайным образом выкинуть номера $\{i_1, \dots, i_n\}$, где $P(i_j = k) = 1/n, k = 1, \dots, n$, тогда $\{X_1^*, \dots, X_n^*\} = \{X_{i_1}, \dots, X_{i_n}\}$ (выборка с возвращением)

Оценка дисперсии с помощью бутстрепа:

1. С моделировать $X_1^*, \dots, X_n^* \sim \hat{F}_n$

2. Подсчитать $T_n^* = g(X_1^*, \dots, X_n^*)$

3. Повторить шаги 1 и 2, получить $T_{n,1}^*, \dots, T_{n,B}^*$

4. Подсчитать оценку дисперсии по формуле

$$v_{\text{boot}} = \frac{1}{B} \sum_{b=1}^B \left(T_{n,b}^* - \frac{1}{B} \sum_{r=1}^B T_{n,r}^* \right)^2$$

Псевдокод для оценки дисперсии оценки медианы:

Входные данные $X = (X(1), \dots, X(n))$:

```
T ← median(X)
```

8 Tboot ← **вектор** длины B

```
for(i in 1:B){  
    Xstar ← выборка объема n с возвращением из X  
    Tboot[i] ← median(Xstar)  
}
```

```
se ← sqrt(variance(Tboot))
```


$$\mathbb{V}_F(T_n) \approx \mathbb{V}_{\hat{F}_n}(T_n) \approx v_{boot}$$

Пример (данные о моментах времени между последовательными импульсами вдоль нервного волокна):

9 $\theta = T(F) = \int (x - \mu)^3 dF(x) / \sigma^3$ - коэффициент асимметрии

$$\hat{\theta} = T(\hat{F}_n) = \frac{\int (x - \mu)^3 d\hat{F}_n(x)}{\hat{\sigma}^3} = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^3}{\hat{\sigma}^3} = 1.76$$

Применяя метод оценки дисперсии на основе бутстрепа при $B = 1,000$ получаем, что стандартная ошибка оценки коэффициента асимметрии равна 0.16

с) Оценка доверительных интервалов на основе бутстрепа

Метод 1: нормальный интервал.

$$T_n \pm z_{\alpha/2} \hat{se}_{boot},$$

где $\hat{se}_{boot} = \sqrt{v_{boot}}$ - оценка стандартной ошибки на основе бутстрепа.

10 Метод хорошо работает, если распределение данных близко к нормальному распределению

Метод 2: центральный интервал. Пусть $\theta = T(F)$ и $\hat{\theta}_n = T(\hat{F}_n)$.

Положим $R_n = \hat{\theta}_n - \theta$. Обозначим через $\hat{\theta}_{n,1}^*, \dots, \hat{\theta}_{n,B}^*$ - повторную

выборку значений $\hat{\theta}_n = T(\hat{F}_n)$ на основе бутстрепа, а через $H(r)$ -

распределение величины R_n , то есть

$$H(r) = \mathbb{P}_F(R_n \leq r)$$

Определим доверительный интервал согласно формуле $C_n^* = (a, b)$, где

$$a = \hat{\theta}_n - H^{-1} \left(1 - \frac{\alpha}{2} \right) \quad \text{и} \quad b = \hat{\theta}_n - H^{-1} \left(\frac{\alpha}{2} \right).$$

Очевидно, что

$$\begin{aligned} \mathbb{P}(a \leq \theta \leq b) &= \mathbb{P}(a - \hat{\theta}_n \leq \theta - \hat{\theta}_n \leq b - \hat{\theta}_n) \\ &= \mathbb{P}(\hat{\theta}_n - b \leq \hat{\theta}_n - \theta \leq \hat{\theta}_n - a) \\ &= \mathbb{P}(\hat{\theta}_n - b \leq R_n \leq \hat{\theta}_n - a) \\ &= H(\hat{\theta}_n - a) - H(\hat{\theta}_n - b) \\ &= H \left(H^{-1} \left(1 - \frac{\alpha}{2} \right) \right) - H \left(H^{-1} \left(\frac{\alpha}{2} \right) \right) \\ &= 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha \end{aligned}$$

Однако, распределение $H(r)$ неизвестно. Будем использовать оценку распределения $H(r)$ на основе бутстрепа, то есть положим

$$\hat{H}(r) = \frac{1}{B} \sum_{b=1}^B I(R_{n,b}^* \leq r)$$

12

где $R_{n,b}^* = \hat{\theta}_{n,b}^* - \hat{\theta}_n$. Пусть r_β^* обозначает β выборочную квантиль, подсчитанную по выборке значений $(R_{n,1}^*, \dots, R_{n,B}^*)$, а θ_β^* обозначает β выборочную квантиль, подсчитанную по выборке значений $(\hat{\theta}_{n,1}^*, \dots, \hat{\theta}_{n,B}^*)$. Несложно показать, что $r_\beta^* = \theta_\beta^* - \hat{\theta}_n$.

Тогда приблизительный доверительный интервал с доверительной вероятностью $1 - \alpha$ имеет вид $C_n = (\hat{a}, \hat{b})$, где

$$\hat{a} = \hat{\theta}_n - \hat{H}^{-1} \left(1 - \frac{\alpha}{2} \right) = \hat{\theta}_n - r_{1-\alpha/2}^* = 2\hat{\theta}_n - \theta_{1-\alpha/2}^*$$

$$\hat{b} = \hat{\theta}_n - \hat{H}^{-1} \left(\frac{\alpha}{2} \right) = \hat{\theta}_n - r_{\alpha/2}^* = 2\hat{\theta}_n - \theta_{\alpha/2}^*$$

13 Таким образом получаем, что центральный доверительный интервал, построенный на основе бутстрепа, имеет вид

$$C_n = \left(2\hat{\theta}_n - \hat{\theta}_{1-\alpha/2}^*, 2\hat{\theta}_n - \hat{\theta}_{\alpha/2}^* \right)$$

Отметим, что $\mathbb{P}_F(T(F) \in C_n) \rightarrow 1 - \alpha$ при $n \rightarrow \infty$

Метод 3: интервал на основе процентилей. Доверительный интервал, построенный на основе процентилей, имеет вид

$$C_n = \left(\theta_{\alpha/2}^*, \theta_{1-\alpha/2}^* \right)$$

Действительно, допустим, что существует монотонное преобразование

$$U = m(T), \quad \text{для которого} \quad U \sim N(\phi, c^2), \quad \text{где} \quad \phi = m(\theta)$$

14 (предполагается только, что такое преобразование существует; конкретный вид преобразования может быть неизвестен). Пусть $U_b^* = m(\theta_{n,b}^*)$.

Обозначим через u_{β}^* - β выборочную квантиль последовательности чисел U_b^* , $b = 1, \dots, B$. Так как монотонное преобразование сохраняет квантили, то $u_{\alpha/2}^* = m(\theta_{\alpha/2}^*)$. Также $U \sim N(\phi, c^2)$, поскольку $\alpha/2$

квантиль величины U равна $\phi - z_{\alpha/2}c$. Таким образом,
 $u_{\alpha/2}^* = \phi - z_{\alpha/2}c$ и, аналогично, $u_{1-\alpha/2}^* = \phi + z_{\alpha/2}c$.

Итак,

$$\begin{aligned}
 15 \quad \mathbb{P}(\theta_{\alpha/2}^* \leq \theta \leq \theta_{1-\alpha/2}^*) &= \mathbb{P}(m(\theta_{\alpha/2}^*) \leq m(\theta) \leq m(\theta_{1-\alpha/2}^*)) \\
 &= \mathbb{P}(u_{\alpha/2}^* \leq \phi \leq u_{1-\alpha/2}^*) \\
 &= \mathbb{P}(U - cz_{\alpha/2} \leq \phi \leq U + cz_{\alpha/2}) \\
 &= \mathbb{P}(-z_{\alpha/2} \leq \frac{U - \phi}{c} \leq z_{\alpha/2}) \\
 &= 1 - \alpha
 \end{aligned}$$

Точное «нормализующее» преобразование существует редко, однако существует много преобразований, которые позволяют делать распределения близкими к нормальному. Например, одно- и двухпараметрические семейства преобразований Бокса-Кокса:

$$y = \begin{cases} \frac{x^\lambda - 1}{\lambda}, \lambda \neq 0 \\ \log(x), \lambda = 0 \end{cases}$$
$$y = \begin{cases} \frac{(x + \lambda_2)^{\lambda_1} - 1}{\lambda_1}, \lambda_1 \neq 0 \\ \log(x + \lambda_2), \lambda_1 = 0 \end{cases}$$

Пример (данные о моментах времени между последовательными импульсами вдоль нервного волокна):

95% интервал для коэффициента асимметрии

17 Нормальный интервал: (1.44, 2.09)
 Центральный интервал: (1.48, 2.11)
 Интервал на основе процентилей: (1.42, 2.03)

Все три типа интервалов имеют сравнимую точность

Пример (данные о значениях холестерина в плазме): построим доверительный интервал на основе бутстрепа для разности медиан

```
x1 <- первая выборка; x2 <- вторая выборка; n1 <- length(x1)
n2 <- length(x2); th.hat <- median(x2) - median(x1)
B <- 1000; Tboot <- вектор длины B
for(i in 1:B){
18   xx1 <- выборка длины n1 с возвращением из x1
   xx2 <- выборка длины n2 с возвращением из x2
   Tboot[i] <- median(xx2) - median(xx1) }
se <- sqrt(variance(Tboot))
Normal <- (th.hat - 2*se, th.hat + 2*se)
percentile <- (quantile(Tboot,.025), quantile(Tboot,.975))
pivotal <- ( 2*th.hat-quantile(Tboot,.975),
            2*th.hat-quantile(Tboot,.025) )
```

95% интервал для разности медиан имеет вид

Нормальный интервал: (3.7, 33.3)

Центральный интервал: (5.0, 34.0)

Интервал на основе процентилей: (5.0, 33.3)

Пример: данные о LSAT - Law School Admissible Test и GPA - Grade Point Average.

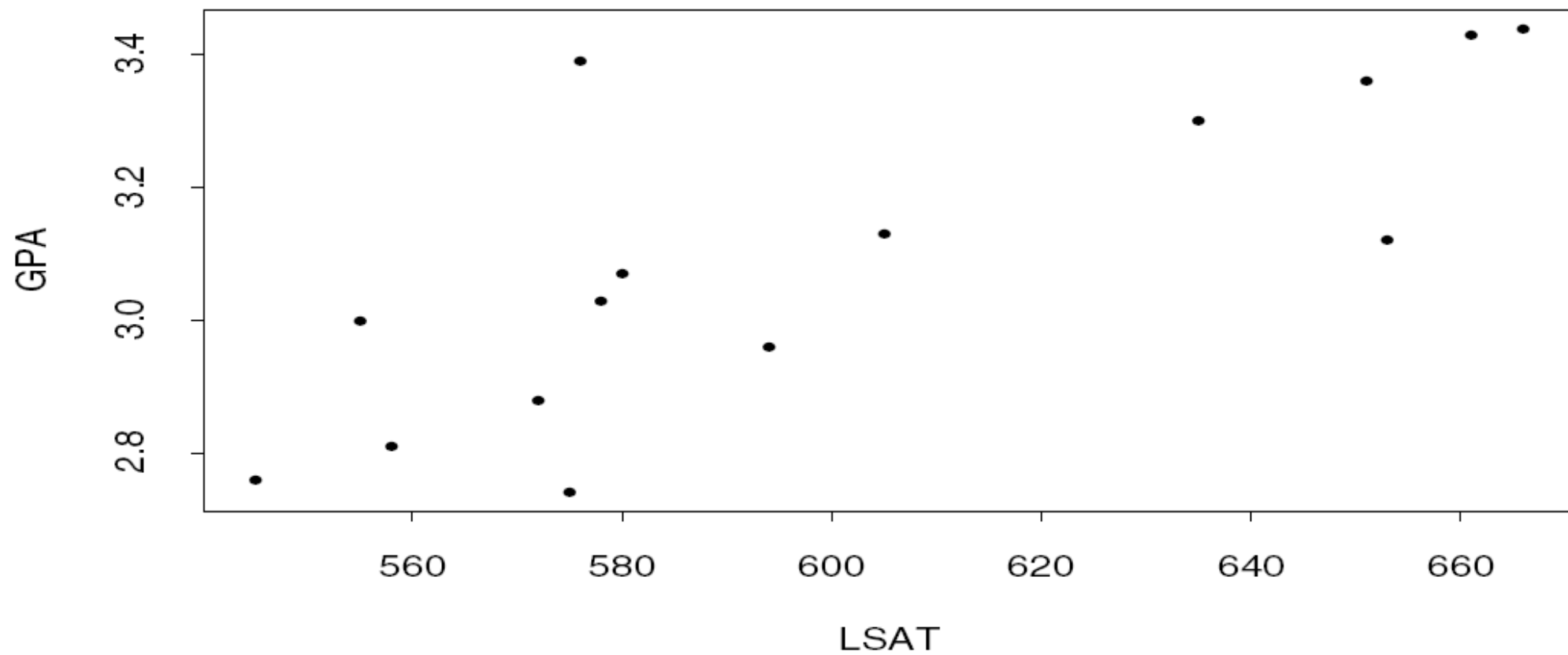
19

LSAT	576	635	558	578	666	580	555	661
	651	605	653	575	545	572	594	

GPA	3.39	3.30	2.81	3.03	3.44	3.07	3.00	3.43
	3.36	3.13	3.12	2.74	2.76	2.88	3.96	

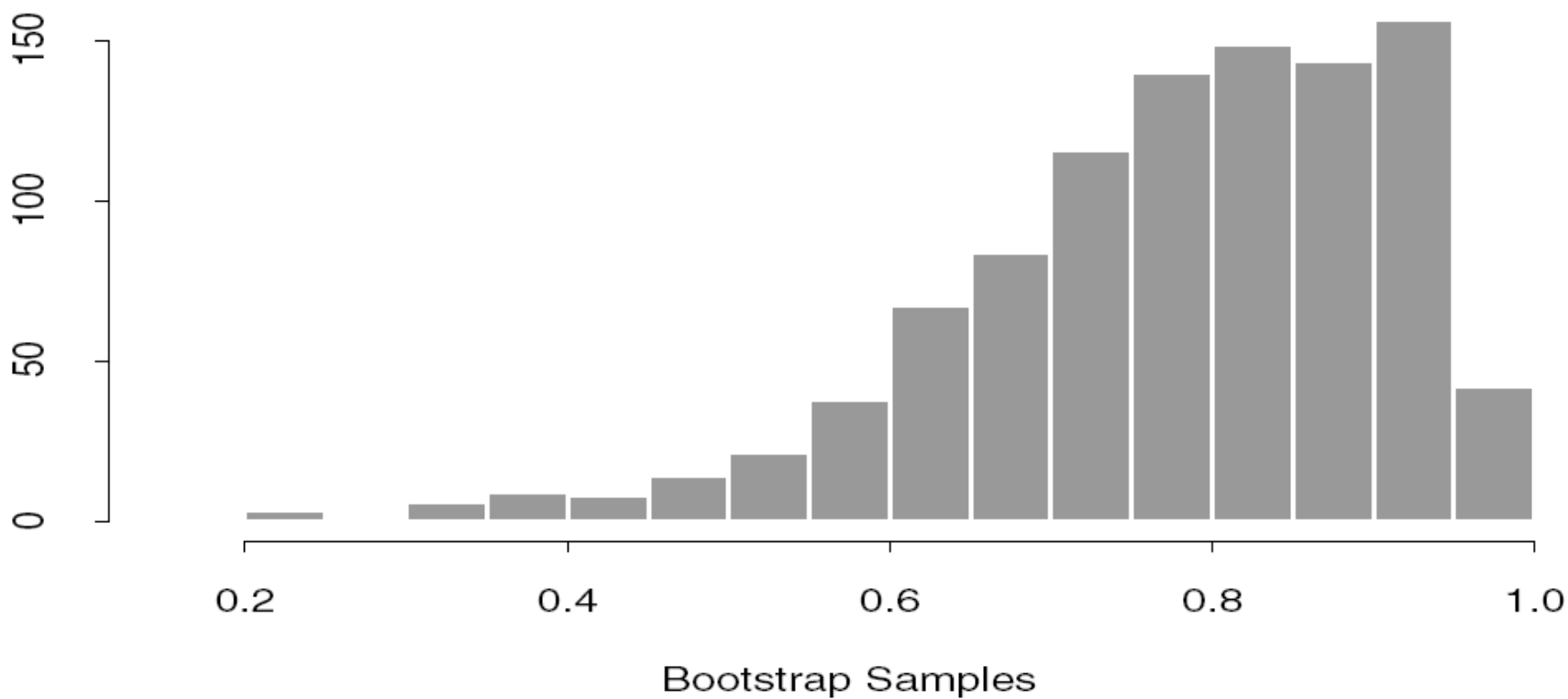
$$X_i = (Y_i, Z_i), \text{ где } Y_i = \text{LSAT}_i \text{ и } Z_i = \text{GPA}_i$$

20



$$\theta = \frac{\int \int (y - \mu_Y)(z - \mu_Z) dF(y, z)}{\sqrt{\int (y - \mu_Y)^2 dF(y) \int (z - \mu_Z)^2 dF(z)}}$$

$$\hat{\theta} = \frac{\sum_i (Y_i - \bar{Y})(Z_i - \bar{Z})}{\sqrt{\sum_i (Y_i - \bar{Y})^2 \sum_i (Z_i - \bar{Z})^2}}, \hat{\theta} = .776$$



Гистограмма значений $\hat{\theta}_1^*, \dots, \hat{\theta}_B^*$ корреляций, подсчитанных по бутстреп-выборкам

Гистограмма = аппроксимация распределения выборочных значений $\hat{\theta}$

На основе бутстрепа с $B = 1000$ была получена оценка стандартной ошибки оценки $\hat{se} = .137$

95% интервал для коэффициента корреляции имеет вид

Нормальный интервал: $.78 \pm 2\hat{se} = (.51, 1.00)$

Интервал на основе процентилей: $(.46, .96)$

Причина сильного отличия – малый объем выборки

Пример: доказательство эквивалентности эффекта от медицинских препаратов

23

номер	плацебо	старое	новое	старое-плацебо	новое-старое
1	9243	17649	16449	8406	-1200
2	9671	12013	14614	2342	2601
3	11792	19979	17274	8187	-2705
4	13357	21816	23798	8459	1982
5	9055	13850	12560	4795	-1290
6	6290	9806	10157	3516	351
7	12412	17208	16570	4796	-638
8	18806	29044	26325	10238	-2719

Обозначим $Z =$ старое-плацебо и $Y =$ старое-новое. Считается, что эффект от препаратов совпадает, если $|\theta| \leq .20$, где

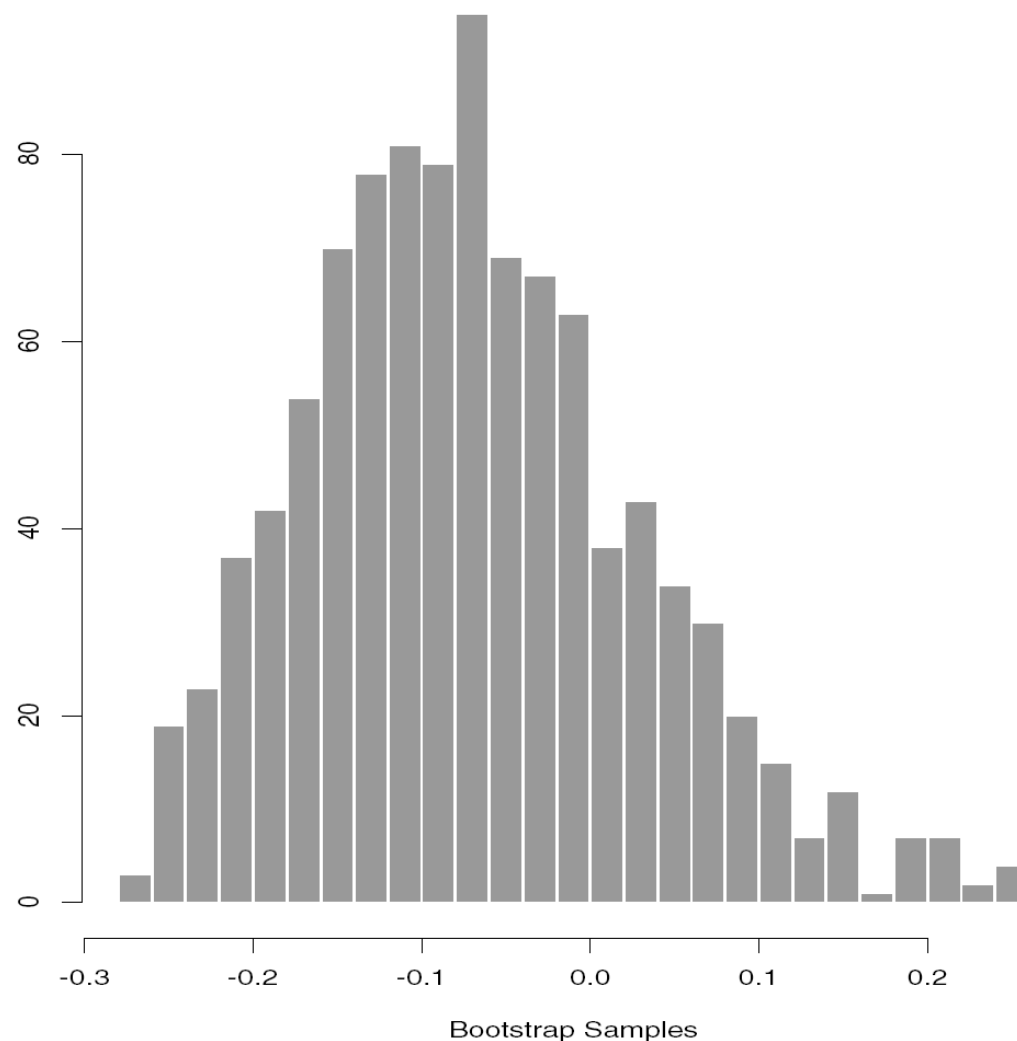
$$\theta = \frac{\mathbb{E}_F(Y)}{\mathbb{E}_F(Z)}$$

$$\hat{\theta} = \frac{\bar{Y}}{\bar{Z}} = \frac{-452.3}{6342} = -0.0713$$

Оценка стандартной ошибки оценки параметра с помощью бутстрепа равна $\hat{se} = 0.105$

При $B = 1000$ 95% центральный интервал равен $(-0.24, 0.15) \notin (-0.20, 0.20) \Rightarrow$ эквивалентность показана не была

На рис. изображена гистограмма значений θ , полученных на основе бутстреп-выборки



d) Метод складного ножа

Пусть $T_n = T(X_1, \dots, X_n)$ - значение статистики, подсчитанное на основе простой выборки X_1, \dots, X_n

26 Обозначим через $T_{(-i)}$ значение статистики, подсчитанное на основе выборки $X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n$

Пусть $\bar{T}_n = n^{-1} \sum_{i=1}^n T_{(-i)}$. Оценка дисперсии $\text{var}(T_n)$ по методу складного ножа равна

$$v_{\text{jack}} = \frac{n-1}{n} \sum_{i=1}^n (T_{(-i)} - \bar{T}_n)^2$$

Оценка стандартной ошибки по методу складного ножа равна

$$\widehat{se}_{jack} = \sqrt{v_{jack}}.$$

Можно показать, что $v_{jack}/\text{var}(T_n) \xrightarrow{P} 1$ при $n \rightarrow \infty$. Однако в отличие от бутстрепа оценка стандартной ошибки оценки квантили на основе метода складного ножа является несостоятельной