

Оценка распределения и статистические функционалы

а) Эмпирическая функция распределения

Пусть задана i.i.d. выборка $X_1, \dots, X_n \sim F$

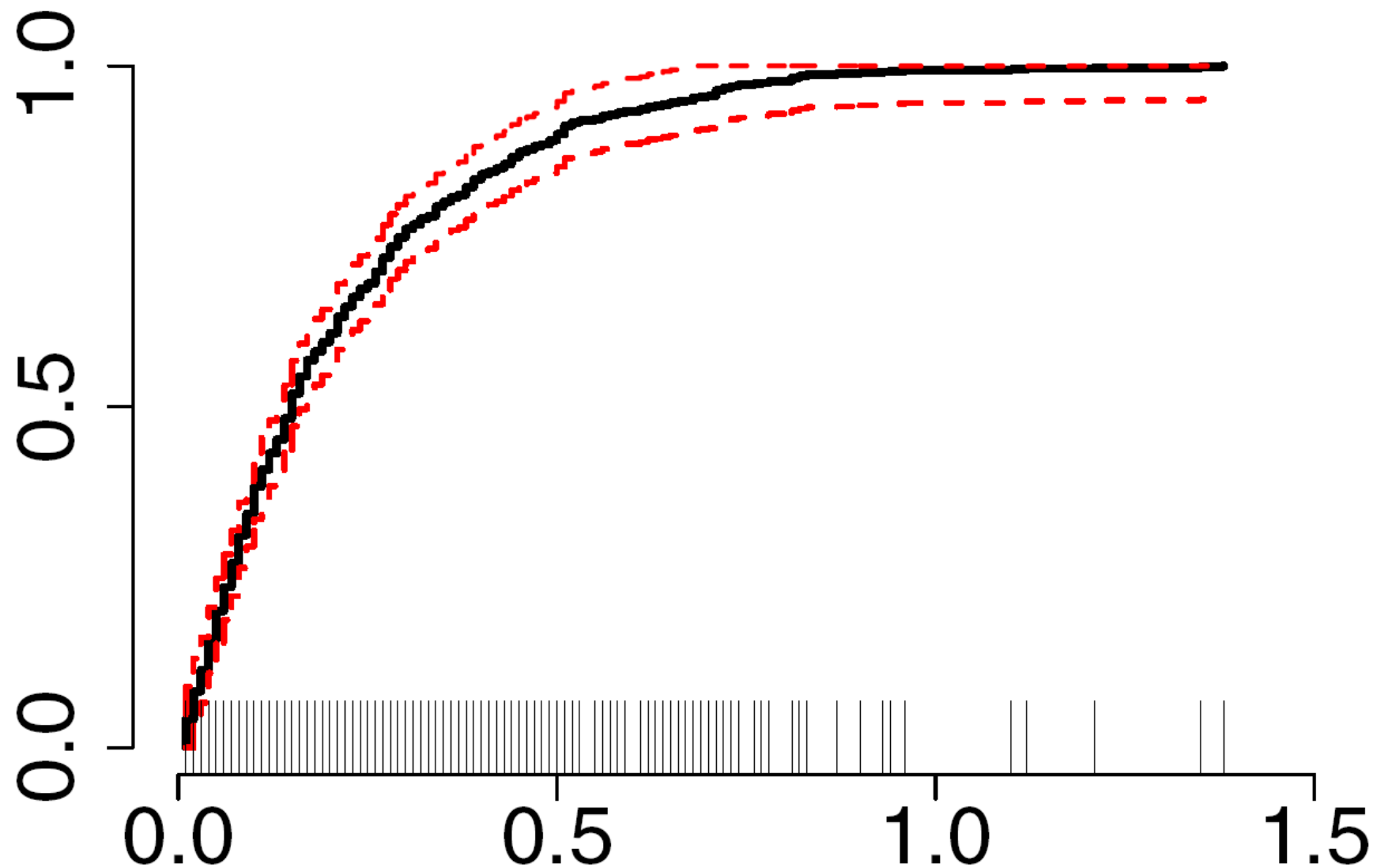
Необходимо оценить распределение F

1 Эмпирическая функция распределения \hat{F}_n имеет вид

$$\hat{F}_n(x) = \frac{\sum_{i=1}^n I(X_i \leq x)}{n}, \text{ где}$$

$$I(X_i \leq x) = \begin{cases} 1, & X_i \leq x \\ 0, & X_i > x \end{cases}$$

2



Эмпирическая функция распределения. Моменты между последовательными импульсами вдоль нервного волокна

Утверждение. Для любого фиксированного x

$$\mathbb{E} \left(\hat{F}_n(x) \right) = F(x),$$

$$\mathbb{V} \left(\hat{F}_n(x) \right) = \frac{F(x)(1 - F(x))}{n},$$

3

$$\text{MSE} = \frac{F(x)(1 - F(x))}{n} \rightarrow 0,$$

$$\hat{F}_n(x) \xrightarrow{\text{P}} F(x).$$

Утверждение (Th. Гливенко-Кантелли). Пусть $X_1, \dots, X_n \sim F$ - i.i.d. выборка. Тогда

$$\sup_x |\hat{F}_n(x) - F(x)| \xrightarrow{\text{P}} 0$$

Утверждение (Неравенство Dvoretzky, Kiefer, Wolfowitz). Пусть $X_1, \dots, X_n \sim F$ - i.i.d. выборка. Тогда для любого $\epsilon > 0$

$$\mathbb{P}\left(\sup_x |F(x) - \hat{F}_n(x)| > \epsilon\right) \leq 2e^{-2n\epsilon^2}$$

Пусть

4

$$L(x) = \max\{\hat{F}_n(x) - \epsilon_n, 0\}$$

$$U(x) = \min\{\hat{F}_n(x) + \epsilon_n, 1\}$$

$$\epsilon_n = \sqrt{\frac{1}{2n} \log\left(\frac{2}{\alpha}\right)}$$

для заданного значения α . Тогда для любого распределения F

$$\mathbb{P}\left(L(x) \leq F(x) \leq U(x) \quad \text{для } \forall x\right) \geq 1 - \alpha$$

б) Статистические функционалы

Пусть задана i.i.d. выборка $X_1, \dots, X_n \sim F$

Статистическим функционалом $T(F)$ называется любая функция от F ,
например,

$$\mu = \int x dF(x)$$

$$\sigma^2 = \int (x - \mu)^2 dF(x)$$

$$m = F^{-1}(1/2)$$

Оценка $T(F)$ получается с помощью величины $\hat{\theta}_n = T(\hat{F}_n)$

Функционалы вида $T(F) = \int r(x) dF(x)$ для некоторой функции $r(x)$
называются линейными статистическими функционалами, поскольку

$$T(aF + bG) = aT(F) + bT(G)$$

Для линейных статистических функционалов

$$T(\hat{F}_n) = \int r(x) d\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n r(X_i)$$

Пусть se - стандартное отклонение для $T(\hat{F}_n)$, а \hat{se} - его оценка.

Тогда во многих случаях

6

$$T(\hat{F}_n) \approx N(T(F), \hat{se}^2)$$

Приближенный доверительный интервал с доверительной вероятностью

$1 - \alpha$ для $T(F)$ будет иметь вид

$$T(\hat{F}_n) \pm z_{\alpha/2} \hat{se}$$

Пример (среднее значение). Пусть $\mu = T(F) = \int x dF(x)$. Тогда

$$\hat{\mu} = \int x d\hat{F}_n(x) = \bar{X}_n$$

$$\text{se} = \sqrt{\mathbb{V}(\bar{X}_n)} = \sigma / \sqrt{n}$$

$$\bar{X}_n \pm z_{\alpha/2} \hat{\text{se}}$$

Пример (дисперсия).

7 Пусть $\sigma^2 = T(F) = \mathbb{V}(X) = \int x^2 dF(x) - \left(\int x dF(x) \right)^2$. Тогда

$$\hat{\sigma}^2 = \int x^2 d\hat{F}_n(x) - \left(\int x d\hat{F}_n(x) \right)^2$$

$$= \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2$$

$$= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

Однако полученная оценка является смещенной. Несмещенная оценка имеет вид

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

Пример (коэффициент асимметрии).

8

$$\kappa = \frac{\mathbb{E}(X - \mu)^3}{\sigma^3} = \frac{\int (x - \mu)^3 dF(x)}{\left\{ \int (x - \mu)^2 dF(x) \right\}^{3/2}}$$
$$\hat{\kappa} = \frac{\int (x - \mu)^3 d\hat{F}_n(x)}{\left\{ \int (x - \mu)^2 d\hat{F}_n(x) \right\}^{3/2}} = \frac{\frac{1}{n} \sum_i (X_i - \hat{\mu})^3}{\hat{\sigma}^3}$$

Пример (корреляция). Пусть $Z = (X, Y)$ - двумерная случайная величина. Определим корреляцию по формуле

$$\rho = T(F) = \mathbb{E}(X - \mu_X)(Y - \mu_Y) / (\sigma_x \sigma_y),$$

где $F(x, y)$ - двумерная функция распределения

$$T_1(F) = \int x dF(z), \quad T_2(F) = \int y dF(z), \quad T_3(F) = \int xy dF(z)$$

9 $T_4(F) = \int x^2 dF(z), \quad T_5(F) = \int y^2 dF(z),$

$$a(t_1, \dots, t_5) = \frac{t_3 - t_1 t_2}{\sqrt{(t_4 - t_1^2)(t_5 - t_2^2)}}$$

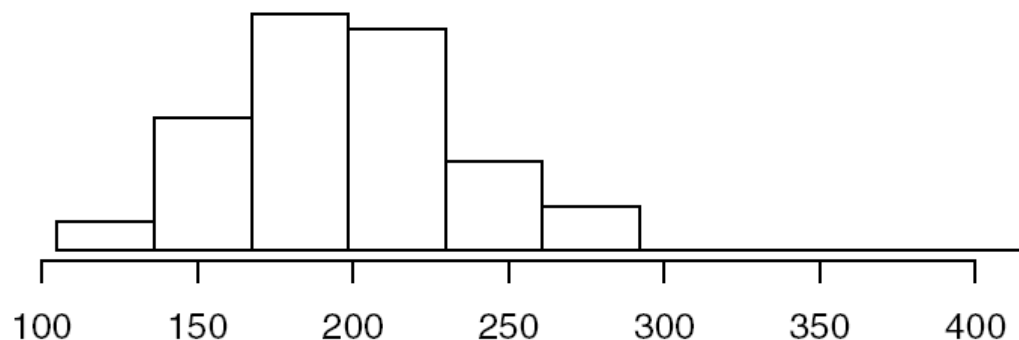
$$\hat{\rho} = a(T_1(\hat{F}_n), T_2(\hat{F}_n), T_3(\hat{F}_n), T_4(\hat{F}_n), T_5(\hat{F}_n))$$

$$\hat{\rho} = \frac{\sum_i (X_i - \bar{X}_n)(Y_i - \bar{Y}_n)}{\sqrt{\sum_i (X_i - \bar{X}_n)^2} \sqrt{\sum_i (Y_i - \bar{Y}_n)^2}}$$

Пример (квантиль). Пусть распределение F строго возрастает с плотностью f . p -квантилью при $0 < p < 1$ распределения F называется величина $T(F) = F^{-1}(p)$. Поскольку \hat{F}_n - необратимая функция, то

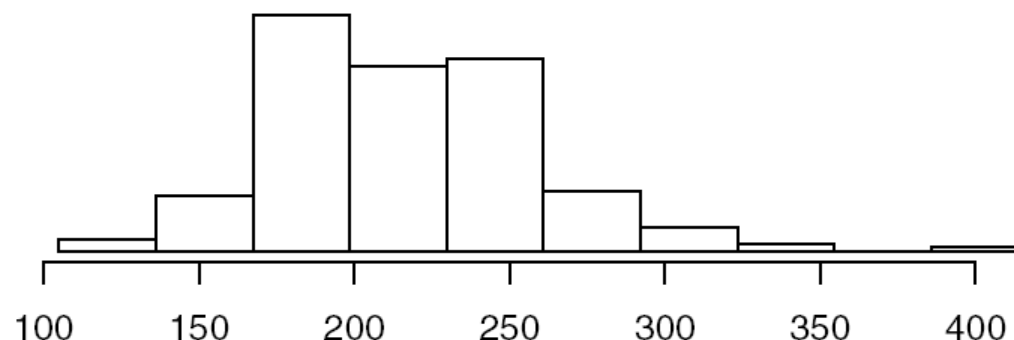
$$T(\hat{F}_n) = \hat{F}_n^{-1}(p) = \inf\{x : \hat{F}_n(x) \geq p\}$$

Пример (проверка гипотезы равенства средних значений). 371 пациент с болью в груди. У 51 нет признаков болезни сердца, у 320 – наблюдается сужение артерий.



11

Распределение в % пациентов в зависимости от величины холестерина в плазме (здоровые)



Распределение в % пациентов в зависимости от величины холестерина в плазме (сердечники)

$$\hat{\mu}_1 = \int x d\hat{F}_{n,1}(x) = \overline{X}_{n,1} = 195.27$$

$$\hat{\mu}_2 = \int x d\hat{F}_{n,2}(x) = \overline{X}_{n,2} = 216.19$$

$$\text{se}(\hat{\mu}) = \sqrt{\mathbb{V}\left(\frac{1}{n} \sum_{i=1}^n X_i\right)} = \sqrt{\frac{1}{n^2} \sum_{i=1}^n \mathbb{V}(X_i)} = \sqrt{\frac{n\sigma^2}{n^2}} = \frac{\sigma}{\sqrt{n}}$$

$$\widehat{\text{se}}(\hat{\mu}) = \frac{\hat{\sigma}}{\sqrt{n}}$$

$$\hat{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \overline{X})^2}$$

$$\widehat{\text{se}}(\hat{\mu}_1) = 5.0, \widehat{\text{se}}(\hat{\mu}_2) = 2.4$$

95% доверительны интервалы для μ_1 и μ_2

$$\hat{\mu}_1 \pm 2\widehat{\text{se}}(\hat{\mu}_1) = (185, 205)$$

$$\hat{\mu}_2 \pm 2\widehat{\text{se}}(\hat{\mu}_2) = (211, 221)$$

13 Рассмотрим функционал

$$\theta = T(F_2) - T(F_1)$$

$$\hat{\theta} = \hat{\mu}_2 - \hat{\mu}_1 = 216.19 - 195.27 = 20.92$$

$$\text{se} = \sqrt{\mathbb{V}(\hat{\mu}_2 - \hat{\mu}_1)} = \sqrt{\mathbb{V}(\hat{\mu}_2) + \mathbb{V}(\hat{\mu}_1)} = \sqrt{(\text{se}(\hat{\mu}_1))^2 + (\text{se}(\hat{\mu}_2))^2}$$

$$\widehat{se} = \sqrt{(\widehat{se}(\widehat{\mu}_1))^2 + (\widehat{se}(\widehat{\mu}_2))^2} = 5.55$$

95% доверительны интервалы для θ

$$\widehat{\theta} \pm 2 \widehat{se}(\widehat{\theta}_n) = (9.8, 32.0)$$

Таким образом, у пациентов с суженными артериями повышенное содержание холестерина. Однако, это совершенно не означает, что холестерин – причина сужения. Может существовать какой-то третий фактор, который влияет и на сужение артерий и на повышение уровня холестерина