

Математическая статистика. ДЗ 2.

ПРОХОРОВ ЮРИЙ, 776

Задача 1

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — простая выборка из неизвестного непрерывного распределения $F(x)$. Доказать, что распределение случайной величины

$$Z_n = \int_{-\infty}^{+\infty} (\hat{F}_n(x) - F(x))^2 dF(x),$$

где $\hat{F}_n(x) = \frac{1}{n} \sum_{k=1}^n \mathbb{I}\{X_k < x\}$ — эмпирическая функция распределения, не зависит от функции $F(x)$.

Так как противного не сказано, будем понимать данный интеграл как интеграл Римана-Стилтьеса по траекториям подынтегральной случайной функции.

Решение:

Введем случайные величины $Y_k = F(X_k)$. Найдем их распределение:

$$\begin{aligned} x \leq 0 & \implies F_{Y_k}(x) = \mathbb{P}\{Y_k < x\} = 0 \\ 0 < x < 1 & \implies F_{Y_k}(x) = \mathbb{P}\{F^{-1}(Y_k) < F^{-1}(x)\} = \mathbb{P}\{X_k < F^{-1}(x)\} = F(F^{-1}(x)) = x \\ x \geq 1 & \implies F_{Y_k}(x) = \mathbb{P}\{Y_k < x\} = 1 \end{aligned}$$

Если функция $F(x)$ не строго монотонна, то под обратной функцией $F^{-1}(y)$ подразумевается

$$F^{-1}(y) = \inf\{x \mid F(x) = y\}$$

Таким образом, $Y_k \sim \mathcal{U}[0, 1]$. С учетом этого перепишем случайную величину Z :

$$\begin{aligned} Z_n &= \int_{-\infty}^{+\infty} \left(\frac{1}{n} \sum_{k=1}^n \mathbb{I}\{X_k < x\} - F(x) \right)^2 dF(x) = \int_{-\infty}^{+\infty} \left(\frac{1}{n} \sum_{k=1}^n \mathbb{I}\{Y_k < F(x)\} - F(x) \right)^2 dF(x) = \\ &= \int_{-\infty}^{+\infty} \left(\frac{1}{n} \sum_{k=1}^n \mathbb{I}\{Y_k < y\} - y \right)^2 dy \end{aligned}$$

Видно, что это выражение уже не зависит от F , что и требовалось доказать.

Оказывается, что для случайной величины Z_n выполнено

$$nZ_n \xrightarrow[n \rightarrow \infty]{d} \int_0^1 B^2(t) dt, \quad B(t) = W(t) - t \cdot W(1),$$

где $B(t)$ — броуновский мост, $W(t)$ — винеровский процесс, а интеграл понимается как стохастический интеграл Римана.

Распределение квадратичного интеграла броуновского моста табулировано, про него можно прочитать в статье “On the Distribution of the Square Integral of the Brownian Bridge”, Leonid Tolmatz, 2002.

Задача 2

(а) Пусть $\xi_1, \dots, \xi_n \sim \text{Be}(p) - \text{i.i.d.}$, $S_n = \xi_1 + \dots + \xi_n$. Доказать, что при $z \geq 0$:

$$\mathbb{P}\{S_n - np \geq z\} \leq \exp\left\{-nH\left(p + \frac{z}{n}\right)\right\} \leq \exp\left\{-\frac{2z^2}{n}\right\},$$

$$\mathbb{P}\{S_n - np \leq -z\} \leq \exp\left\{-nH\left(p - \frac{z}{n}\right)\right\} \leq \exp\left\{-\frac{2z^2}{n}\right\},$$

где функция $H(x) = x \ln \frac{x}{p} + (1-x) \ln \frac{1-x}{1-p}$.

(б) Пусть \hat{F}_n — эмпирическая функция распределения произвольной случайной величины с теоретической функцией распределения F .

Оценить вероятность

$$\mathbb{P}\left\{\sqrt{n}|\hat{F}_n(x) - F(x)| \geq \Delta\right\}$$

Решение:

(а) Докажем первое неравенство.

Для произвольного $\lambda > 0$:

$$\begin{aligned} \mathbb{P}\{S_n - np \geq z\} &= \mathbb{P}\{e^{S_n} \geq e^{np+z}\} = \mathbb{P}\{e^{\lambda S_n} \geq e^{\lambda(np+z)}\} = \left/ \begin{array}{c} \text{неравенство} \\ \text{Маркова} \end{array} \right/ \leq \\ &\leq \frac{\mathbb{E}e^{\lambda S_n}}{e^{\lambda(np+z)}} = \left/ \begin{array}{c} S_n - \text{сумма} \\ \text{независимых с.в.} \end{array} \right/ = \frac{\prod_{k=1}^n \mathbb{E}e^{\lambda \xi_k}}{e^{\lambda(np+z)}} = \frac{[1 + p(e^\lambda - 1)]^n}{e^{\lambda(np+z)}} = \\ &= \left[\frac{1 + p(e^\lambda - 1)}{e^{\lambda(p + \frac{z}{n})}} \right]^n = f(\lambda) \end{aligned}$$

Отсюда следует, что

$$\mathbb{P}\{S_n - np \geq z\} \leq \inf_{\lambda > 0} \left[\frac{1 + p(e^\lambda - 1)}{e^{\lambda(p + \frac{z}{n})}} \right]^n = \inf_{\lambda > 0} [qe^{-\lambda(p+\alpha)} + pe^{\lambda(q-\alpha)}]^n = \inf_{\lambda > 0} f(\lambda),$$

где $q = 1 - p$, $\alpha = \frac{z}{n}$.

По сути мы применили сейчас неравенство Чернова.

Из условия $f'(\lambda) = 0$ находим минимум

$$\lambda^* = \ln \frac{q(p+\alpha)}{p(q-\alpha)}, \quad f(\lambda^*) = \left(q \left[\frac{q(p+\alpha)}{p(q-\alpha)} \right]^{-(p+\alpha)} + p \left[\frac{q(p+\alpha)}{p(q-\alpha)} \right]^{q-\alpha} \right)^n = \left[\frac{p^{p+\alpha} q^{q-\alpha}}{(p+\alpha)^{p+\alpha} (q-\alpha)^{q-\alpha}} \right]^n$$

После преобразований получаем

$$\mathbb{P}\{S_n - np \geq z\} \leq e^{-nH(p+\alpha)}$$

Для дальнейшей оценки, заметим, что

$$H(p) = H'(p) = 0, \quad H''(x) = \frac{1}{x(1-x)} \leq 4, \quad H(p+\alpha) \geq \frac{1}{2} \cdot \max |H''| \cdot \alpha^2 = 2\alpha^2$$

Строгое доказательство:

$$\frac{d}{d\alpha} [H(p+\alpha) - 2\alpha^2] = 0 \quad \implies \quad \varphi(\alpha) \equiv H'(p+\alpha) - 4\alpha = \ln \frac{p+\alpha}{p} - \ln \frac{1-p-\alpha}{1-p} - 4\alpha = 0$$

$$\left. \begin{array}{l} \varphi'(\alpha) = H''(p+\alpha) - 4 \geq 0 \\ \varphi(-p) = -\infty \\ \varphi(1-p) = +\infty \end{array} \right\} \implies \begin{array}{l} \text{по теореме о промежуточном значении:} \\ \exists! \text{ решение уравнения } \varphi(\alpha) = 0 \end{array}$$

Это решение можно угадать: $\alpha^* = 0$, поэтому

$$H(p + \alpha) - 2\alpha^2 \geq H(p + \alpha^*) - 2(\alpha^*)^2 = H(p) = 0$$

Наконец, получаем требуемое неравенство

$$\mathbb{P}\{S_n - np \geq z\} \leq e^{-nH(p+\alpha)} \leq e^{-2\alpha^2 n}$$

Второе неравенство получается аналогично.

(b) Пусть \hat{F}_n строится по выборке X_1, \dots, X_n . Обозначим случайные величины

$$\xi_k = \mathbb{I}\{X_k < x\} \sim \text{Be}(F(x)) = \text{Be}(p) - \text{i.i.d.}$$

$$S_n = \xi_1 + \dots + \xi_n$$

Тогда можно воспользоваться пунктом (a):

$$\mathbb{P}\left\{\sqrt{n}|\hat{F}_n(x) - F(x)| \geq \Delta\right\} = \mathbb{P}\{|S_n - np| \geq \sqrt{n}\Delta\} \leq 2e^{-2\Delta^2}$$