

## Домашнее задание № 3

Тема: Доверительные интервалы. Общая теория построения статистических тестов.

Крайний срок сдачи: 15 ноября 2020 г. (до конца дня).

### 1

- T1 (i) (1.5 балла) Пусть даны i.i.d. случайные величины  $X_1, \dots, X_n$  с непрерывной функцией распределения  $F(x)$  (возможно, зависящей от параметров). Докажите, что

$$-\sum_{i=1}^n \log(1 - F(X_i)) \sim \Gamma(n, 1).$$

Пользуясь этим результатом, постройте точный доверительный интервал методом центральной функции для параметра  $\lambda$  из распределения Вейбулла с функцией распределения

$$F(x) = 1 - e^{-(x/\lambda)^\tau}, \quad x > 0.$$

Параметр  $\tau$  предполагается известным.

- (ii) (1.5 балла) Для случая  $\tau = 1$  (экспоненциальное распределение) постройте асимптотические доверительные интервалы для параметра  $\lambda$ , используя асимптотическую нормальность статистик

$$(iii a) \quad \frac{1}{n} \sum_{i=1}^n X_i; \quad (iii b) \quad \sqrt{\frac{1}{2n} \sum_{i=1}^n X_i^2}.$$

Выясните, какой из этих двух асимптотических интервалов имеет меньшую длину при больших  $n$  (для ответа на этот вопрос можно использовать компьютер).

(iii) (1 балл) Для этого же интервала постройте непараметрическую оценку параметра  $\lambda$ , используя зависимость медианы распределения от параметра  $\lambda$ .

N1 (1 балл) Промоделируйте выборку длины  $n = 1000$  из экспоненциального распределения с фиксированным параметром  $\lambda$ . Расширяя подвыборку от 50 до 1000 наблюдений, вычислите длины доверительных интервалов, построенных в п. (iii) и (iv) (при этом зафиксируйте  $\alpha = 0.05$ ) и постройте графики, выражающие зависимость длины интервала от размера выборки (отдельно для каждого метода). Отобразите зависимость длины интервала от  $\alpha$  (для этого зафиксируйте параметр  $n = 1000$ ).

## 2

T2 (1 балл) Проводится опрос общественного мнения с целью понять различия в предпочтениях между жителями городской и сельской местности. Задавался единственный вопрос - "Поддерживаете ли Вы действия правительства", возможные ответы - да или нет. Было решено опросить равное количество респондентов -  $n$  человек в городе и  $n$  человек в сельской местности. Нужно проводить опрос до тех пор, пока разность между вероятностями поддержки в городе и селе не будет оценена с точностью 0.05. Решение принимается на основе асимптотического доверительного интервала с уровнем доверия 0.95. Определите, какое минимальное количество респондентов  $2n$  нужно опросить.

T3 (1 балл) Количество респондентов для опроса из [T2] вычисляется исходя тестирования гипотезы "уровни поддержки в городе и селе совпадают" против альтернативы "уровень поддержки в городе превышает уровень поддержки в сельской местности на 3%". По-прежнему в городе и селе опрашивается равное количество респондентов. Исходя из проведённых ранее опросов предполагается, что уровень поддержки в городе примерно равен 14%, а в сельской местности - 9%.

Определите, какое минимальное количество респондентов нужно опросить, чтобы допустимая вероятность ошибки первого рода составила 0.05, а допустимая вероятность ошибки второго рода - 0.04.

*Указание.* Условие "Исходя из проведённых ранее опросов предполагается, что уровень поддержки в городе примерно равен 14%, а в сельской местности - 9%" следует использовать для оценивания асимптотической дисперсии (метод подстановки параметра).

### 3

Т4 (1.5 балла) Дана выборка длины  $n$  из распределения с функцией плотности

$$p(x, \theta) = \theta(1 - x)^{\theta-1}, \quad x \in [0, 1]$$

с неизвестным параметром  $\theta > 0$ . Опишите равномерно наиболее мощный тест для проверки гипотезы  $\theta = \theta_0$  против альтернативы  $\theta > \theta_0$ . Докажите, что мощность этого теста равна

$$W(\theta) = F_{(n, \theta_0)}\left(\frac{\theta}{\theta_0} \cdot q_{(n, \theta_0)}(\alpha)\right),$$

где  $F_{(n, \theta_0)}, q_{(n, \theta_0)}(\alpha)$  - это функция распределения и  $\alpha$ -квантиль гамма-распределения, имеющего плотность

$$p(x) = \frac{\theta_0^n}{\Gamma(n)} x^{n-1} e^{-\theta_0 x}, \quad x > 0.$$

Т5 (1.5 балла) Наблюдаемые величины  $X_1, \dots, X_n$  имеют нормальное распределение с неизвестным средним  $\mu$  и известной дисперсией  $\sigma^2$ . Для тестирования гипотезы  $H_0 : \mu = \mu_0$  против альтернативы  $H_1 : \mu = \mu_1 > \mu_0$  используются тесты вида  $\{\bar{X} \geq c\}$ , где  $\bar{X} = (X_1 + \dots + X_n)/n$ . Докажите, что при выборе  $c = (\mu_1 + \mu_0)/2$  получается тест, у которого сумма ошибок первого и второго рода минимальна среди всех возможных тестов такого вида.

## 4

Т6\* (2 балла) Пусть задана выборка  $x_1, \dots, x_n$  из экспоненциального семейства распределений с плотностью

$$p(x, \theta) = g(x)e^{x\theta - d(\theta)},$$

где  $g : \mathbb{R} \rightarrow \mathbb{R}_+$  и  $d : \mathbb{R} \rightarrow \mathbb{R}$ . Для тестирования гипотезы  $H_0 : \theta = \theta_0$  против альтернативы  $\theta > \theta_0$  используется LR (likelihood ratio) тест вида

$$\left\{ \Lambda(\vec{x}) := \frac{\max_{\theta > \theta_0} L(\vec{x}, \theta)}{L(\vec{x}, \theta_0)} > c_\alpha \right\}, \quad (1)$$

где  $L(\vec{x}, \theta) = \prod_{i=1}^n p(x_i, \theta)$  - функция правдоподобия, и пороговый уровень  $c_\alpha$  подбирается из условия, что уровень значимости теста равен  $\alpha \in (0, 1)$ .

(i) Докажите, что

$$\log(\Lambda(\vec{x})) = \begin{cases} nK(\hat{\theta}, \theta_0), & \text{если } \hat{\theta} > \theta_0, \\ 0, & \text{иначе,} \end{cases}$$

где  $\hat{\theta}$  - оценка максимального правдоподобия (вычисленная по выборке  $x_1, \dots, x_n$ ), а  $K$  - расстояние Кульбака-Ляйблера, определяемое как

$$K(\hat{\theta}, \theta_0) = \int \log\left(\frac{p(x, \hat{\theta})}{p(x, \theta_0)}\right) p(x, \theta_0) dx.$$

(ii) Докажите, что LR - тест (1) в данной модели может быть записан в виде

$$\left\{ \hat{\theta} > \theta_0 + \tilde{c}_\alpha \right\},$$

где  $\tilde{c}_\alpha > 0$ .