

Математическая статистика

Матвеев Сергей М3338

5 семестр

1 Предмет математической статистики. Задачи, решаемые методами математической статистики.

Пусть у нас есть генеральная совокупность, но мы хотим ее как-то изучать, тогда мы можем взять ее часть - выборку.

Мы хотим по выборке сделать содержательные вероятностные выводы о генеральной совокупности.

Примеры задач, которые могут быть решены таким способом:

1. Бросок монеты (оценить вероятность орла, честно|нечестно)
2. Замеры показателя: какие типичные значения для показателя
3. Как учатся мальчики и девочки (одинаково или по разному)
4. Цена на недвижимость, расстояние до метро (оценка зависимости)

2 Модель простейшей выборки. Эмпирическая функция распределения и её свойства. Способы визуализации выборки.

Простейшая модель выборки - X_1, X_2, \dots, X_n - *i.i.d.*, F - функция распределения (теоретическая функция).

$X_1, \dots, X_n \sim F$ (F мы не знаем априори)

x_1, \dots, x_n - реализация выборки

Цель: оценить из реализации x_1, \dots, x_n теор F .

Эмпирическая функция распределения:

$$\mu_n(t) = \sum_{i=1}^n \mathbb{1}(X_i \leq t)$$

$$F_n(t) = \frac{\mu_n(t)}{n} - \text{эмпирическая функция распределения.}$$

Свойства

$$\mathbb{1}(X_i \leq t) \sim \text{Bern}(F(t))$$

$$\mathbb{E}(F_n(t)) = F(t) \text{ (это называется несмещенность)}$$

$$\text{Var}(F_n(t)) = \frac{F(t)(1-F(t))}{n}$$

$$\text{ЗБЧ: } F_n(t) \xrightarrow{P} F(t) \text{ - это называется состоятельность}$$

$$\text{ЦПТ: } \frac{\mu_n(t) - nF(t)}{\sqrt{F(t)(1-F(t))n}} \xrightarrow{d} U \sim N(0, 1) =$$

$$= \sqrt{n} \frac{F_n(t) - F(t)}{\sqrt{F(t)(1-F(t))}} \text{ (асимптотическая нормальность)}$$

Теорема Гливенко-Кантелли

$$\sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow[n \rightarrow \infty]{a.s.} 0$$

Теорема Колмогорова

$$D_n = \sup_x |F_n(x) - F(x)| \Rightarrow P(\sqrt{n}D_n \leq t) \xrightarrow[n \rightarrow \infty]{} K(t) = \sum_{j=-\infty}^{+\infty} (-1)^j e^{-2j^2 t^2}$$

$$F \in C(\mathbb{R})$$

Такая функция называется функцией Колмогорова

Теорема Смирнова

$X_1, \dots, X_n, Y_1, \dots, Y_n$ - независимы

Обе распределены по $F \in C(\mathbb{R})$

$$D_{n,m} = \sup_x |F_n(x) - F_m(x)| \Rightarrow P\left(\sqrt{\frac{mn}{m+n}} D_{n,m} \leq t\right) \xrightarrow[n \rightarrow \infty, m \rightarrow \infty]{} K(t)$$

Стоит отметить, что обе теоремы имеют смысл при $t \geq 0$

Выборку можно визуализировать используя гистограмму или box plot

3 Начальные выборочные моменты и их свойства, в том числе выборочное среднее.

$\alpha_k = EX_1^k$ - k -ый теоретический момент.

$$\overline{g(X)} = \frac{1}{n} \sum_{k=1}^n g(X_k), g: \mathbb{R} \rightarrow \mathbb{R}$$

$$\widehat{\alpha}_k = \overline{X^k} = \frac{1}{n} \sum_{j=1}^n X_j^k \text{ - } k\text{-ый выборочный момент.}$$

$E\widehat{\alpha}_k = \alpha_k$ (несмещенность, мы просто воспользовались линейностью математического ожидания)

$$\text{Var } \widehat{\alpha}_k = \frac{1}{n} \text{Var}(X_1^k) = \frac{1}{n} (EX_1^{2k} - (EX_1^k)^2)$$

По ЦПТ получаем:

$$\sqrt{n} \frac{\widehat{\alpha}_k - \alpha_k}{\sqrt{\alpha_{2k} - \alpha_k^2}} \approx N(0, 1)$$

$\sqrt{n} \frac{\widehat{\alpha}_k - \alpha_k}{\sqrt{\widehat{\alpha}_{2k} - \widehat{\alpha}_k^2}} = \sqrt{n} \frac{\widehat{\alpha}_k - \alpha_k}{\sqrt{\alpha_{2k} - \alpha_k^2}} \cdot \sqrt{\frac{\alpha_{2k} - \alpha_k^2}{\widehat{\alpha}_{2k} - \widehat{\alpha}_k^2}}$ - первый множитель по ЦПТ сходится к $N(0, 1)$ по распределению

Давайте посмотрим что будет со вторым множителем. Он будет сходиться к 1 по вероятности.

Таким образом:

$$\sqrt{n} \frac{\widehat{\alpha}_k - \alpha_k}{\sqrt{\widehat{\alpha}_{2k} - \widehat{\alpha}_k^2}} \cdot \sqrt{\frac{\alpha_{2k} - \alpha_k^2}{\widehat{\alpha}_{2k} - \widehat{\alpha}_k^2}} \xrightarrow{d} N(0, 1)$$

А почему вторая дробь сходится к единице?

$\widehat{\alpha}_k \xrightarrow{P} \alpha_k$ (по ЗБЦ, это называется состоятельность)

$$\widehat{\alpha}_{2k} - \widehat{\alpha}_k^2 \xrightarrow{P} \alpha_{2k} - \alpha_k^2$$

\bar{X} - выборочное среднее.

4 Выборочные центральные моменты и их свойства, в том числе выборочная дисперсия. Дельта-метод.

$\beta_k = E(X_1 - EX_1)^k$ - k -ый центральный момент.

$$\widehat{\beta}_k = \overline{(X - \bar{X})^k} = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X})^k$$
 - k -ый выборочный момент.

$\widehat{\beta}_2 = S_*^2$ - выборочная дисперсия.

S_* - выборочное стандартное отклонение (выборочное среднеквадратичное отклонение).

Note : выборочные моменты есть ничто иное как моменты посчитанные относительно эмпирического распределения.

$$S_* = \overline{X^2} - (\bar{X})^2$$

$$\widehat{\beta}_k = Poly(\widehat{\alpha}_k, \dots, \widehat{\alpha}_1)$$

$\widehat{\alpha}_1, \dots, \widehat{\alpha}_k$ - состоятельные оценки (имеет место сходимость по вероятности)

Так как полином это непрерывная функция, то $\widehat{\beta}_k \xrightarrow{P} \beta_k$

$$ES_*^2 = \frac{1}{n} \sum_{i=1}^n EX_i^2 - E(\bar{X})^2 = \alpha_2 - \alpha_1^2 - \frac{\beta_2}{n} = \beta_2 - \frac{\beta_2}{n} = \frac{n-1}{n} \beta_2$$

Т.к. $ES_*^2 \neq \beta_2$ такую дисперсию называют несмещенной

$$S^2 = \frac{n}{n-1} S_*^2$$
 - несмещенная дисперсия

5 Выборочные квантили, в том числе выборочная медиана. Распределение k-ой порядковой статистики.

Определение. Вариационный ряд

$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ - вариационный ряд

Определение. Порядковая статистика

$X_{(k)}$ - k-я порядковая статистика.

Квантиль порядка α

$$P(X \geq q_\alpha) \geq 1 - \alpha$$

$$P(X \leq \alpha) \geq \alpha$$

Это общее определение

Если F строго возрастает:

$$F(q_\alpha) = \alpha \Leftrightarrow q_\alpha = F^{-1}(\alpha)$$

$$F^{-1}(\alpha) : \sup\{x : F(x) \leq \alpha\}, \inf\{x : F(x) \geq \alpha\}$$

Определение. Выборочный квантиль порядка α

$\alpha \in (0, 1)$

$$\exists 0 \leq k \leq n - 1 : \frac{k}{n} \leq \alpha < \frac{k+1}{n}$$

$X_{(k+1)}$ - выборочный квантиль порядка α

$\alpha = 0$ - нулевой квантиль

$$F^{-1} = \sum \{x \in \mathbb{R} F_n(x) \leq \alpha\}$$

$\alpha = \frac{1}{4}$ - первый квартиль (нижний квартиль)

$\alpha = \frac{1}{2}$ - второй квартиль (выборочная медиана)

$\alpha = \frac{3}{4}$ - третий квартиль (верхний квартиль)

$\alpha = 1$ - max(X) (четвертый квартиль)

$$n = 2m \Rightarrow med(X) = \frac{X_{(m)} + X_{(m+1)}}{2}$$

$$n = 2m + 1 \Rightarrow med(X) = X_{(m+1)}$$

$IQR = \Delta$ между верхним и нижним квартилем.

$$P(X_{(k)} \leq t) = P(\mu_n(t) \geq k) = \sum_{j=k}^n C_n^j F^j(t) (1 - F(t))^{n-j}$$

$$B(z, a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^z t^{a-1} (1-t)^{b-1} dt = B(F(t), k, n-k+1)$$

$$0 \leq z \leq 1$$

Пусть $p()$ - теоретическая плотность, то есть $p = F'$

$$(P(X_{(k)} \leq t))'_t = \frac{\Gamma(n+1)}{\Gamma(k)\Gamma(n-k+1)} \cdot F^{k-1}(t)(1-F(t))^{n-k} \cdot p(t) - \text{плотность}$$

k-й порядковой статистики

6 Теоремы об асимптотиках среднего и крайнего членов вариационного ряда (идеи доказательств).

Средний член вариационного ряда: $\frac{K(n)}{n} \rightarrow const \in (0, 1)$

Крайний член вариационного ряда:

$X_{(r)}$, r - огр.

$X_{(n+1-s)}$, s - огр.

Теорема об асимптотике среднего члена вариационного ряда

$0 < \alpha < 1$ - теоретическая плотность.

q_α - теоретический квантиль порядка α

$p \in C^1$ (окр-сть q_α)

$p(q_\alpha) > 0$

Тогда:

$$\sqrt{n} \cdot f(q_\alpha) \frac{X_{(\lfloor n\alpha \rfloor)} - q_\alpha}{\sqrt{\alpha(1-\alpha)}} \xrightarrow{d} N(0, 1)$$

Идея доказательства

Пусть $\lfloor n\alpha \rfloor = k$

Мы умеем писать плотность для $X_{(k)}$

Затем у нас идет преобразование:

$$g(x) = \sqrt{np(q_\alpha)} \frac{x - q_\alpha}{\sqrt{\alpha(1-\alpha)}} \rightsquigarrow p_{g(X_{(k)})}(t) = p_{X_{(k)}}(g^{-1}(t)) |g^{-1}(t)'| \text{ (теорема из}$$

прошлого семестра)

Там вылезут факториалы, от них мы умеем избавляться по Стирлингу

Затем надо будет воспользоваться непрерывной дифференцируемостью:

$p \in C^1$ (окр-сть q_α)

$p(q_\alpha) > 0$

Тогда в пределе наша новая плотность будет стремиться к плотности нормального стандартного закона.

Теорема об асимптотике крайних членов вариационного ряда

$r, s, F, x, p()$ - плотность

Тогда:

$$nF(X_{(r)}) \xrightarrow{d} \Gamma(r, 1)$$

$$n(1 - F(X_{(n+1-s)})) \xrightarrow{d} \Gamma(s, 1)$$

И оба распределения независимы.

Идея доказательства

У нас есть совместная плотность и какое-то преобразование, тогда мы можем написать плотность после преобразования

Затем берем предел и мы получим плотность равная произведению двух этих двух законов.

7 Постановка задачи точечного оценивания параметров. Свойства оценок.

$X_1, \dots, X_n \sim F_\theta \in \Theta \subset \mathbb{R}^d$

θ - некий фиксированный неизвестный вектор.

Наша цель оценить θ в виде $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$

Определение. Состоятельность

$\hat{\theta}$ - состоятельная оценка $\theta \Leftrightarrow \hat{\theta} \xrightarrow{P} \theta$

Определение. Несмещенность

$b.as(\hat{\theta}) \stackrel{def}{=} E\hat{\theta} - \theta$ - смещение

$b.as(\hat{\theta}) = 0 \Leftrightarrow$ несмещенная

Определение. Асимптотическая нормальность

$\sqrt{n}(\hat{\theta} - \theta) \rightarrow N(0, \Sigma_\theta)$

Определение. Эффективность (оптимальность)

$\hat{\theta}_1$ эффективнее $\hat{\theta}_2 \Leftrightarrow MSE\hat{\theta}_1 < MSE\hat{\theta}_2$

$MSE\hat{\theta} = E \left\| \hat{\theta} - \theta \right\|^2 = E(\hat{\theta} - \theta)^T (\hat{\theta} - \theta)$

Утверждение

$MSE\hat{\theta} = tr(\text{Var } \hat{\theta}) + \left\| b.as\hat{\theta} \right\|^2$

Доказательство

$MSE = E(\hat{\theta} - \theta)^T (\hat{\theta} - \theta) = E(\hat{\theta} - E\hat{\theta} + E\hat{\theta} - \theta)^T (\hat{\theta} - E\hat{\theta} + E\hat{\theta} - \theta) = E(\hat{\theta} - E\hat{\theta})^T (\hat{\theta} - E\hat{\theta}) + \sum \text{Var } \hat{\theta}_i + \left\| b.as\hat{\theta} \right\|^2$

1. Асимптотическая нормальность \Rightarrow состоятельность

$$\hat{\theta} - \theta = \frac{1}{\sqrt{n}} \sqrt{n}(\hat{\theta} - \theta) \xrightarrow{P} 0$$

2. Асимптотическая нормальность $\Rightarrow b.as\hat{\theta} \rightarrow 0$

Пусть $d = 1$

$$P(|\hat{\theta} - E\hat{\theta}| > \varepsilon) = P\left(\frac{\sqrt{n}|\hat{\theta} - E\hat{\theta}|}{\sigma} > \frac{\varepsilon\sqrt{n}}{\sigma}\right) = 1 - P(\dots < \frac{\varepsilon\sqrt{n}}{\sigma}) \approx$$

$$1 - (2\Phi(\frac{\varepsilon\sqrt{n}}{\sigma}) - 1) = 2(1 - \Phi(\frac{\varepsilon\sqrt{n}}{\sigma})) \rightarrow 0$$

3. Состоятельность $\Rightarrow b.as\hat{\theta} \rightarrow 0$

Следует из усиленного закона больших чисел

$\bar{X} \xrightarrow{a.s} \mu \Rightarrow E\bar{X} \rightarrow \mu$ (По теореме Лебега о мажорируемой сходимости)

4. Пусть $d = 1$, $b.as\hat{\theta} \rightarrow 0$, $\text{Var } \hat{\theta} \rightarrow 0 \Rightarrow \hat{\theta}$ - сост.

8 Метод моментов.

Рассмотрим g_1, \dots, g_d

$\exists Eg_1(X_1) = m_1(\theta_1, \dots, \theta_d)$

$$\exists E g_2(X_2) = m_2(\theta_1, \dots, \theta_d)$$

...

$$\exists E g_d(X_d) = m_d(\theta_1, \dots, \theta_d)$$

$$\begin{cases} \overline{g_1(X)} = m_1(\hat{\theta}_1, \dots, \hat{\theta}_d) \\ \dots \\ \overline{g_d(X)} = m_d(\hat{\theta}_1, \dots, \hat{\theta}_d) \end{cases}$$

Пусть $\exists!$ решение:

$$\begin{cases} \hat{\theta}_1 = \alpha_1(\overline{g_1(X)}, \dots, \overline{g_d(X)}) \\ \dots \\ \hat{\theta}_d = \alpha_d(\overline{g_1(X)}, \dots, \overline{g_d(X)}) \end{cases}$$

Тогда это будет оценка методов моментов.

9 Метод максимального правдоподобия.

probability must function: $p(x, \theta) = p(x|\theta)$

probability identity function: $p(x, \theta) = p(x|\theta)$

Будем называть оба случая плотностью.

Пусть у нас есть выборка $X_1, \dots, X_n \sim p(x|\theta)$

$L(x|\theta) = \prod p(x_i|\theta)$ - функция правдоподобия

$\theta^* = \underset{\hat{\theta}}{\operatorname{argmax}} (L(x, \theta))$ - оценка максимума правдоподобия

$\theta \in \Theta$ - открыто

$\theta_1 \neq \theta_2 \Rightarrow L(x, \theta_1) \neq L(x, \theta_2)$

Доказательство

1. Посмотреть и подумать
2. Рассмотреть $\ln L(x|\theta)$; $\frac{\partial \ln L(x, \theta)}{\partial \theta}$
3. $\frac{\partial \ln L(x, \theta)}{\partial \theta} = 0$
4. Проверить достаточные условия максимума

10 Оценки максимального правдоподобия для нормальной и полиномиальной моделей.

1) $N(\theta_1, \theta_2)$

$$L(x, \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\theta_2}} \exp\left(-\frac{(x_i - \theta_1)^2}{2\theta_2}\right)$$

$$\ln L(x, \theta) = \sum_{i=1}^n \left[-\frac{1}{2} \ln 2\pi - \frac{1}{2} \theta_2 - \frac{(x_i - \theta_1)^2}{2\theta_2} \right]$$

$$\frac{\partial \ln L(x, \theta_1)}{\partial \theta_1} = \sum_{i=1}^n \frac{2(x_i - \theta_1)}{2\theta_2} = \sum_{i=1}^n \frac{x_i - \theta_1}{\theta_2}$$

$$\frac{\partial \ln L(x, \theta)}{\partial \theta_2} = \sum_{i=1}^n \left[-\frac{1}{2\theta_2} + \frac{(x_i - \theta_1)^2}{2\theta_2^2} \right]$$

$$\sum_{i=1}^n \frac{(x_i - \hat{\theta}_1)}{\hat{\theta}_2} = 0 \Rightarrow \hat{\theta}_1 = \bar{X}$$

$$\sum_{i=1}^n \left[-\frac{1}{2\hat{\theta}_2} + \frac{(x_i - \hat{\theta}_1)^2}{2\hat{\theta}_2^2} \right] = 0 \Rightarrow \hat{\theta}_2 = S_*^2$$

2) $Poly(1, p) : p = (p_1, \dots, p_m)$

Рассмотрим частоты:

ν_1 - кол-во наблюдений типа 1

...

ν_m - кол-во наблюдений типа m

Суммируем и смотрим на функцию правдоподобия

$$L(X, p) = p_1 \dots p_m$$

$$\ln L(X, p) = \sum_{j=1}^{m-1} \nu_j \ln p_j + \nu_m \ln (1 - p_1 - \dots - p_{m-1})$$

$$\frac{\partial \ln L \dots}{\partial p_j} = \frac{\nu_j}{p_j} - \frac{\nu_m}{1 - p_1 - \dots - p_{m-1}} = 0$$

$$\sum \text{уравнения: } \nu_j (1 - \hat{p}_1 - \dots - \hat{p}_{m-1}) = \hat{p}_j \cdot \nu_m$$

$$\hat{p}_m (n - \nu_m) = \nu_m (1 - \hat{p}_m)$$

$$\hat{p}_m n - \hat{p}_m \nu_m = \nu_m$$

$$\hat{p}_m = \frac{\nu_m}{n}$$

$$\hat{p}_j = \frac{\nu_j \hat{p}_m}{\nu_m} = \frac{\nu_j}{n}$$

11 Информация Фишера.

$$d = 1 : L(X, \theta) = \prod p(X_j, \theta)$$

$$\ln L(X, \theta) = \sum \ln p(x_j, \theta)$$

$$V(X, \theta) = \frac{\partial \ln L \dots}{\partial \theta} = \sum \frac{\partial \ln p \dots}{\partial \theta} - \text{вклад выборки}$$

$\theta \in \Theta$ - открыто

$$\theta_1 \neq \theta_2 \Rightarrow p(X, \theta_1) \neq p(X, \theta_2)$$

Регулярность:

$$1. \frac{\partial}{\partial \theta} \int_X T(X) L(X, \theta) dX = \int \frac{\partial}{\partial \theta} L(X, \theta) \cdot T(X) dX$$

Необходимое условие $\sup p_x$ не зависит от θ

$$U[0, \theta] \int_0^\theta \frac{1}{\theta} dt = 1$$

$$\left(\int_0^\theta \frac{1}{\theta} dt\right)'_\theta = \left(\frac{1}{\theta} \int_0^\theta dt\right)'_\theta = -\frac{1}{\theta^2} \int_0^\theta dt + \frac{1}{\theta} = 0 \neq \int_0^\theta \left(\frac{1}{\theta}\right)'_\theta dt$$

$$2. EV^2(X, \theta) < \infty$$

$$\int_X L(X, \theta) dX = 1 \xrightarrow{\frac{\partial}{\partial \theta}} \int_X \frac{\partial L(\cdot)}{\partial \theta} dX = \int_X \frac{\frac{\partial L(\dots)}{\partial \theta}}{L(\dots)} \cdot L(\dots) dX = \int_X V(X, \theta) L(X, \theta) dX = EV(X, \theta) = 0$$

$I(\theta) = \text{Var}(V(X_i, \theta)) = E(V^2(X_i, \theta))$ - информация Фишера для всей выборки

$$V(X, \theta) = \sum_j \frac{\partial \ln p(x, \theta)}{\partial \theta} \Rightarrow \text{Var}(V(X, \theta)) = n \cdot \text{Var} \frac{\partial \ln p(x, \theta)}{\partial \theta}$$

$i(\theta)$ - информация Фишера для 1 наблюдения

$$i(\theta) = E\left(\frac{\partial \ln p(x_j, \theta)}{\partial \theta}\right)^2$$

$$\frac{\partial}{\partial \theta} \int_{\mathbb{R}} \frac{\partial \ln p(x, \theta)}{\partial \theta} \cdot p(x, \theta) dx = \int_{\mathbb{R}} \frac{\partial^2 \ln p(x, \theta)}{\partial^2 \theta} dx + \int_{\mathbb{R}} \frac{\partial \ln \dots}{\partial \theta} \frac{\partial p \dots}{\partial \theta} \text{ вот тут домножаем и делим на плотность}$$

$$E \frac{\partial^2 \ln p(x, \theta)}{\partial \theta^2} + E\left(\frac{\partial \ln p(x, \theta)}{\partial \theta}\right)^2 = 0$$

$$i(\theta) = E\left(\frac{\partial \ln p(x_j, \theta)}{\partial \theta}\right)^2 = -E \frac{\partial^2 \ln p(x, \theta)}{\partial \theta^2}$$

Произвольное d :

$$i(\theta) = -\left(E \frac{\partial^2 \ln p(X, \theta)}{\partial \theta_i \partial \theta_j}\right)_{1 \leq i, j \leq d}$$

$$I(\theta) = ni(\theta)$$

12 Неравенство Рао-Крамера.

Модель регулярная, $d = 1$

$\tau(\theta)$ - оцениваемая функция

$\tau \in C^1$ (как правило $\tau(\theta) = \theta$)

$$E\widehat{\tau(\theta)} = \theta$$

Тогда:

$$\text{Var} \widehat{\tau(\theta)} \geq \frac{[\tau'(\theta)]^2}{ni(\theta)}$$

$$\tau'(\theta) = \int \widehat{\tau(\theta)} \frac{\partial L(X, \theta)}{\partial \theta} dX = \int \widehat{\tau(\theta)} V(X, \theta) L(X, \theta) dX - EV(X, \theta) \cdot E\widehat{\tau(\theta)} =$$

$$\text{Cov}(V(X, \theta), \widehat{\tau(\theta)})$$

$$\text{Cov}^2(V(X, \theta), \widehat{\tau(\theta)}) \leq \text{Var}(V(X, \theta)) \cdot \text{Var}(\widehat{\tau(\theta)})$$

Многомерный случай

$\tau(\theta) : \mathbb{R}^d \rightarrow \mathbb{R}$

$\tau \in C^1$

$$E\widehat{\tau\theta} = \tau\theta \Rightarrow \text{Var} \widehat{\tau(\theta)} \geq \frac{\nabla \tau(\theta) i^{-1}(\theta) \nabla^T \tau(\theta)}{n}$$

13 Свойства оценок максимального правдоподобия.

Состоятельность

Пусть θ_0 - реальный параметр $\Rightarrow p_{\theta_0}(L(X, \theta_0) > L(X, \theta)) \rightarrow 1$

$$\frac{L(X, \theta)}{L(X, \theta_0)} < 1$$

$$\frac{1}{n} \sum \ln \frac{p(X_j, \theta)}{p(X_j, \theta_0)} < 0$$

$$\text{По ЗБЧ} \Rightarrow E_{\theta_0} \ln \frac{p(x_j, \theta)}{p(x_j, \theta_0)} \leq E_{\theta_0} \left[\frac{p(X_j, \theta)}{p(X_j, \theta_0)} - 1 \right] = \int_X p(X, \theta) dX - \int_X p(X, \theta_0) dX = 0$$

Давайте введем события:

$$S_n = \{X : \ln L(X, \theta_0) > \ln L(X, \theta_0 - a)\} \cap \{X : \ln L(X, \theta_0) > \ln L(X, \theta_0 + a)\}$$

$$P_{\theta_0}(S_n) \rightarrow 1$$

$$A_n = \{X : |\hat{\theta} - \theta_0| < a\}$$

$$B_n = \{X : \frac{\partial \ln L(X, \theta)}{\partial \theta} \big|_{\theta=\hat{\theta}} = 0\}$$

$$S_n \subset A_n B_n \subset A_n \Rightarrow P(A_n) \rightarrow 1$$

Принцип инвариантности

$$\theta \in \Theta \xrightarrow{\text{bijection}} \gamma \in \Gamma$$

$$\theta = \varphi^{-1}(\gamma) \Leftrightarrow \gamma = \varphi(\theta)$$

$$\sup_{\theta} L(X, \varphi(\gamma)) = \sup_{\gamma} L(x, \gamma)$$

$$\gamma^* = \varphi(\theta^*)$$

Теорема Асимптотическая нормальность ОМП

Пусть наша модель регулярная, так же пусть:

$$\left| \frac{\partial^3 \ln f(x, \theta)}{\partial \theta_i \partial \theta_j \partial \theta_k} \right| \leq M$$

θ_* - ОМП для θ

Уравнение $\nabla \ln L(X, \theta) = 0$ имеет единственное решение. Тогда:

$$1. \sqrt{n}(\theta_* - \theta) \rightarrow N(0, i^{-1}(\theta))$$

$$2. \tau(\theta) - \text{оцениваемая функция от } \theta$$

$$\tau \in C^1$$

$$\sqrt{n}(\tau(\theta_*) - \tau(\theta)) \rightarrow N(0, \sigma^2)$$

$$\sigma^2 = \nabla \tau(\theta) i^{-1}(\theta) \nabla^T \tau(\theta)$$

$$3. \sigma^2 - \text{непрерывная функция от } \theta \Rightarrow \sqrt{\frac{\tau(\theta_*) - \tau(\theta)}{\sigma(\theta_*)}} \rightarrow N(0, 1)$$

14 Экспоненциальное семейство распределений.

Пусть наше распределение относится к экспоненциальному семейству распределений если:

$$p(x, \theta) = \exp(A(\theta)B(x) + C(\theta) + D(x))$$

К таким распределениям относятся: $N(), \Gamma(), Pois(), Bin, NB$

$$\ln p(x, \theta) = A(\theta)B(x) + C(\theta) + D(x)$$

$$\frac{\partial \ln p(x, \theta)}{\partial \theta} = A'(\theta)B(x) + C'(\theta)$$

$$V(X, \theta) = A'(\theta) \sum B(X_i) + nC'(\theta)$$

$$V(X, \theta) = n(A'(\theta)\overline{B(X)} + C''(\theta))$$

$$\frac{V(X, \theta)}{n} - C'(\theta) = A'(\theta)\overline{B(X)}$$

$$\frac{n}{\overline{B(X)}} = \frac{V(X, \theta)}{nA'(\theta)} - \frac{C'(\theta)}{A'(\theta)}$$

$$\overline{B(X)} - \text{оптимальная оценка для } \left(-\frac{C'(\theta)}{A'(\theta)}\right)$$

15 Байесовские оценки.

$$X_1, \dots, X_n \sim F_\theta$$

$$\theta \sim \pi(\theta) - \text{prior}$$

$$l(\hat{\theta}, \theta) - \text{функция потерь}$$

$$l(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2 \text{ (default)}$$

$$R(\hat{\theta}, \theta) = El(\hat{\theta}, \theta) - \text{риск}$$

$$r(\hat{\theta}) = E_{\pi(\theta)} R(\hat{\theta}, \theta) - \text{байесовский риск}$$

$$\hat{\theta}_B = \underset{\theta}{\operatorname{argmin}} r(\hat{\theta})$$

$$r(\hat{\theta}) = El(\hat{\theta}, \theta)$$

Давайте вспомним теорему Байеса:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$P(\theta|X) = \frac{L(X|\theta)\pi(\theta)}{\int L(X|\theta)\pi(\theta)d\theta}$$

$$\text{posterior} = \text{likelihood} \times \text{prior}$$

$$\hat{\theta}_B = \underset{\hat{\theta}}{\operatorname{argmin}} E[l(\hat{\theta})|X]$$

$$r(\theta_*) \leq r(\hat{\theta})$$

16 Минимаксные оценки.

$$m(\hat{\theta}) = \sup_{\theta} R(\hat{\theta}, \theta)$$

$$\hat{\theta}_{WC} = \operatorname{argmin}_{\theta} m(\hat{\theta}) - \text{минимаксная оценка}$$

$$r(\hat{\theta}) \leq m(\hat{\theta})$$

Утверждение

$$\exists \pi(\theta) - \text{prior} : R(\hat{\theta}_B, \theta) = \text{const} \Rightarrow \hat{\theta}_{WC} = \hat{\theta}_B$$

17 Доверительные интервалы. Общая схема построения доверительного интервала. «Универсальный» рецепт.

Определение. Доверительный интервал

$$X_1, \dots, X_n \sim F_{\theta}, \theta \in \Theta \subset \mathbb{R}$$

$1 - \alpha = \gamma \in (0, 1)$ - уровень доверия

default: 0.9, 0.95, 0.99

$(T_l(X), T_r(X))$ - доверительный интервал уровня $\gamma = 1 - \alpha$ если $p(\theta \in$

$$(T_l(X), T_r(X)) \geq \gamma)$$

Пусть $T(X, \theta) \sim G$ - не зависит от θ

Рассмотрим $p(q_1 < T(X, \theta) < q_2) = 1 - \alpha$

$$q_1 = q_{\frac{\alpha}{2}}$$

$$q_2 = q_{1 - \frac{\alpha}{2}}$$

Из данной вероятности можно выразить θ

Универсальный рецепт (нет)

$$\text{а) } F_{\theta}(X_k)$$

$$P(F_{\theta}(X_k) \leq t) = P(X_k \leq F_{\theta}^{-1}(t)) = F_{\theta}(F_{\theta}^{-1}(t)) = t$$

$$\text{б) } -\ln F_{\theta}(X_k) \sim \text{Exp}(1)$$

$$P(-\ln U \leq t) = P(U \geq e^{-t}) = 1 - e^{-t}$$

$$\text{в) } -\sum \ln F_{\theta}(X_k) \sim \Gamma(n, 1)$$

18 Теорема Фишера и примыкающие к ней леммы. Распределение хи-квадрат.

Лемма о независимости линейной и квадратичной статистик

$$X_1, \dots, X_n \sim N(\mu, \sigma^2)$$

$$T = AX, X = (X_1, \dots, X_n)^T, A \in M_{m \times n}(\mathbb{R})$$

$$Q = X^T B X, B \in M_n(\mathbb{R}), B = B^T$$

$$AB = 0$$

Тогда T, Q - независимы

Доказательство

$$\Lambda = U^T B U$$

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m, 0, 0)$$

λ_k - собственное число не 0

$U = (u_1, \dots, u_n)$ - собственные векторы ортонормированного базиса $\Leftrightarrow B =$

$$U \Lambda U^T = \sum_{j=1}^m \lambda_j u_j u_j^T \Rightarrow Q = \sum_{j=1}^M \lambda_j (X^T U_j)(U_j^T X) = \sum_j \lambda_j (U_j^T X)^2$$

$$A \left(\sum_{j=1}^m \lambda_j u_j u_j^T \right) = 0$$

$$\sum_{j=1}^m \lambda_j A U_j u_j^T = 0$$

Зафиксируем $1 \leq k \leq m$ домножим справа на u_k

$$A u_k = 0 \Rightarrow \forall i A[i, *] u_k = 0$$

Нам надо доказать, что $\forall i, k A[i, *] X$ и $u_k^T X$ - нез

$$\text{Cov}(A[i, *] X, u_k^T X) = \text{Cov}(A[0, *] X, X^T u_k) = A[i, *] \text{Var } X u_k = \sigma^2 A[i, *] u_k = 0$$

Лемма о независимости двух квадратичных статистик

$$Q_1 = X^T B_1 X$$

$$Q_2 = X^T B_2 X$$

$$B_1 B_2 = B_2 B_1 = 0$$

Тогда Q_1, Q_2 - нез

Определение X_n - квадратичная

$$X_1, \dots, X_n \sim N(0, 1)$$

$$\sum_{k=1}^n X_k^2 \sim \Xi^2(n) \text{ (распределение)}$$

хи-квадрат с n степенями свободы

Лемма о распределении квадратичной статистики

$$X_1, \dots, X_n \sim N(0, 1)$$

$$Q = X^T B X$$

$$B = B^2$$

$$\text{Тогда } Q \sim \Xi^2(r), r = \text{rank}(B) = \text{tr}(B)$$

$$Q = \sum_{k=1}^n (u_k^T X)^2 \sim \Xi^2(r)$$

Доказательство

$$u_k^T \sim N(u_k^T E X, u_k^T I_n u_k) = N(0, 1)$$

$$\text{Cov}(u_k^T X, u_j^T X) = 0$$

$$B = U \Lambda U^T$$

$$\text{rank } B = \text{rank } \Lambda = \text{tr } \Lambda$$

$$\text{tr } B = \text{tr}(U \Lambda U^T)$$

$$\text{Заметим что } B_{j,j} = \lambda_j u_j u_j^T = \lambda_j$$

Теорема Фишера

$$X_1, \dots, X_n \sim N(\mu, \sigma^2) \Rightarrow$$

1. $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$
2. $\frac{nS_*^2}{\sigma^2} = \frac{(n-1)S^2}{\sigma^2} \sim \Xi^2(n-1)$
3. S^2, \bar{X} - нез
4. S_*^2, \bar{X} - нез

$$Y_j = \frac{X_j - \mu}{\sigma}$$

$$\bar{Y} = \frac{1}{\sigma}(\bar{X} - \mu)$$

$$S_*^2(Y) = \frac{1}{n} \sum_{j=1}^n (Y_j - \bar{Y})^2 = \frac{S_*^2(X)}{\sigma^2}$$

$$\bar{Y} = \frac{\sum Y_j}{n} = (\frac{1}{n}, \dots, \frac{1}{n})(Y_1, \dots, Y_n)^T = bY$$

$$nS_*^2(Y) = (Y - BY)^T(Y - BY) = Y^T(I - B)^T(I - B)Y = \sum \Xi^2(tr(I - B)) -$$

по предыдущей лемме

Для того чтобы доказать третье утверждение

$$b(I - B) = b - b = 0$$

Тогда мы пользуемся первой леммой

Таким образом теорема Фишера доказана.

19 Распределения Фишера, Стюдента. Построение доверительных интервалов для параметров нормального закона.

Определение. Распределение Стюдента

X_0, X_1, \dots, X_n - нез, $N(0, 1)$

$$\frac{X_0}{\sqrt{\frac{1}{n} \sum_{k=1}^n X_k^2}} \sim T(n)$$

$$\sqrt{\frac{1}{n} \sum_{k=1}^n X_k^2}$$

n - степени свободы (deg of freedom)

Давайте выведем статистику:

$$\frac{\sqrt{n} \frac{\bar{X} - \mu}{\sigma}}{\sqrt{\frac{1}{n-1} \frac{nS_*^2}{\sigma^2}}} = \sqrt{n-1} \frac{\bar{X} - \mu}{S_*} \sim T(n-1)$$

$$\frac{\sqrt{n} \frac{\bar{X} - \mu}{\sigma}}{\sqrt{\frac{1}{n-1} \frac{(n-1)S^2}{\sigma^2}}} = \sqrt{n} \frac{\bar{X} - \mu}{S}$$

Доверительный интервал:

$$\bar{X} - \frac{q_{1-\frac{\alpha}{2}}S}{\sqrt{n}}, \bar{X} + \frac{q_{1-\frac{\alpha}{2}}S}{\sqrt{n}}$$

Определение. Распределение Фишера

$$\Xi_n^2 \sim \Xi^2(n)$$

$$\Xi_m^2 \sim \Xi^2(m)$$

Они независимы

$$\frac{\Xi_n^2(n)}{n}$$

$$\frac{\Xi_m^2(m)}{m} \sim F(n, m)$$

Мы доказали теорему Фишера давайте теперь с помощью теоремы мы рассмотрим задачу построения доверительных интервалов нормального закона

$$P(-q_{1-\frac{\alpha}{2}} \leq \frac{\bar{X} - \mu}{\sigma^2} \leq q_{1-\frac{\alpha}{2}}) - \text{ну и потом просто выражаем}$$

- σ^2 - известно, $\mu = ?$

$$\text{Рассмотрим два варианта: } \frac{X_1 - \mu}{\sigma} \sim N(0, 1)$$

$$\sqrt{n} \frac{\bar{X} - \mu}{\sigma} \sim N(0, 1)$$

$$\text{Доверительный интервал уровня } 1 - \alpha \left[\bar{X} - \frac{q_{1-\frac{\alpha}{2}}\sigma}{\sqrt{n}}, \bar{X} + \frac{q_{1-\frac{\alpha}{2}}\sigma}{\sqrt{n}} \right]$$

- μ - известно, $\sigma^2 = ?$

$$-q \leq \sqrt{n} \frac{(\bar{X} - \mu)}{\sigma} \leq q$$

$$-q\sigma \leq \sqrt{n}(\bar{X} - \mu) \leq q\sigma$$

$$\frac{\sqrt{n}(\bar{X} - \mu)}{q} \leq \sigma$$

$$-\frac{q}{\sqrt{n}(\bar{X} - \mu)} \leq \sigma$$

Рассмотрим следующую статистику:

$$\sum \frac{(X_i - \mu)^2}{\sigma^2} \sim \Xi^2(n)$$

$$P(q_{\frac{\alpha}{2}} \leq \sum \frac{(X_i - \mu)^2}{\sigma^2} \leq q_{1-\frac{\alpha}{2}}) = 1 - \alpha$$

Доверительный интервал:

$$\frac{\sum (X_i - \mu)^2}{q_{1-\frac{\alpha}{2}}} \leq \sigma^2 \leq \frac{\sum (X_i - \mu)^2}{q_{\frac{\alpha}{2}}}$$

Давайте теперь рассмотрим задачу построения доверительного интервала $\mu = ?$, $\sigma^2 = ?$

Воспользуемся теоремой Фишера:

$$\frac{nS_*^2}{\sigma^2} \sim \Xi^2(n-1)$$

$$\text{Доверительный интервал: } \frac{nS_*^2}{q_{1-\frac{\alpha}{2}}} \leq \sigma^2 \leq \frac{nS_*^2}{q_{\frac{\alpha}{2}}}$$

20 Асимптотические доверительные интервалы. Доверительный интервал для математического ожидания, дисперсии, медианы.

Раньше мы говорили $P(\theta \in (l_n, r_n)) \geq 1 - \alpha$

Теперь же мы будем говорить $\lim_{n \rightarrow \infty} P(\theta \in (l_n, r_n)) \geq 1 - \alpha$

$T(X, \theta) \xrightarrow{d} G$ не зависит от θ

ЦПТ и ее следствия

$$\sqrt{n} \frac{\bar{X} - \mu}{\sigma} \rightarrow N(0, 1)$$

$$\sqrt{n} \frac{\bar{X} - \mu}{S} \rightarrow N(0, 1)$$

$$\text{Доверительный интервал: } \bar{X} \pm \frac{q_{1-\frac{\alpha}{2}} S}{\sqrt{n}}$$

$$\sqrt{n} \frac{S_*^2 - \sigma^2}{\sqrt{\beta_4 - S_*^4}} \rightarrow N(0, 1)$$

$$\text{Доверительный интервал: } S_*^2 \pm \frac{q_{1-\frac{\alpha}{2}} \sqrt{\beta_4 - S_*^4}}{\sqrt{n}}$$

Теорема об асимптотике среднего члена вариационного ряда

$$\sqrt{n} \frac{X_{(\lfloor np \rfloor)} - q_p}{\sqrt{p(1-p)}} f(q_p) \rightarrow N(0, 1)$$

$$\text{Доверительный интервал для медианы: } p = \frac{1}{2}$$

$$\sqrt{n} f(q_p) \frac{X_{(\lfloor \frac{n}{2} \rfloor)} - q_p}{\frac{1}{2}} \quad (\text{зачастую } f \text{ это константа})$$

$$\text{Доверительный интервал: } X_{(\lfloor \frac{n}{2} \rfloor)} \pm \frac{q_{1-\frac{\alpha}{2}}}{\sqrt{n} \cdot \text{const}}$$

21 Постановка задачи проверки статистических гипотез: выбор нулевой и альтернативной гипотез, общий принцип работы статистического теста, p-value, ошибки I и II рода (false positive and false negative).

Нам надо будет выделить основное предположение (по умолчанию) и альтернативное предположение (наше подозрение или то, что мы хотим доказать)

Давайте рассуждать:

рациональное, с точки зрения инопланетянина (не опираться на жизненный опыт)

X_1, \dots, X_n - выборка в широком смысле.

$(X_1, \dots, X_n) \sim F$

H_0 - нулевая гипотеза.

H_1 - альтернатива.

Так же пусть нам дали уровень значимости

$\alpha \in (0, 1)$ (по умолчанию 0.1, 0.05, 0.01, 0.001)

Статистический тест (критерий)

$$\delta(X, \alpha, H_0, H_1) = \begin{cases} \text{accept } H_0 \\ \text{reject } H_0 (w. \text{ respect to } H_1) \end{cases} \quad \text{То есть в первом случае}$$

данные не противоречат H_0 , а втором противоречат.

Но это не значит, что мы доказали утверждение.

Пусть у нас есть функция $T(X)$ - статистика критерия

$T(X)$ либо в точности, либо в пределе стремится к G при условии $H_0 (\sim$
or $\rightarrow)$

$P(T(X) \in T_0(\alpha) | H_0) = 1 - \alpha$

if $T(x) \in T_1(\alpha)$: reject H_0

else: accept H_0

$T_0(\alpha)$ - область принятия

$T_1(\alpha)$ - область опровержения

1. left: $T_0(\alpha) = [q_\alpha, +\infty)$

2. right: $T_0(\alpha) = (-\infty, q_\alpha]$

3. two: $T_0(\alpha) = [q_{\frac{\alpha}{2}}, q_{1-\frac{\alpha}{2}}], T_1(\alpha) = \overline{T_0(\alpha)}$

$p_l = P(U \leq T(x) | H_0)$

$p_r = P(U > T(x) | H_0)$

$p = 2\min(p_l, p_r)$

if $p < \alpha$: reject H_0

else: accept H_1

Если мы опровергли нулевую гипотезу, но она верна, то это ошибка первого рода (false positive), ее максимальная вероятность это в точности α

Если мы не опровергли нулевую гипотезу, но она была не верна, то это ошибка второго рода (false negative), ее вероятность это β

$\beta = P(T(X) \in T_0(\alpha) | H_1)$

22 Статистические тесты, основанные на доверительных интервалах (z-тест для одной/двух выборок, t-тест для одной/двух выборок, F-тест).

Когда мы строили доверительные интервалы то мы зажимали статистику между квантилями. Это похоже на двухсторонний тест.

$$X_1, \dots, X_n \sim F_\theta$$

$$T(X, \theta) \rightarrow U \sim G$$

$$P(q_{\frac{\alpha}{2}} \leq T(X, \theta) \leq q_{1-\frac{\alpha}{2}}) = 1 - \alpha$$

$$H_0 : \theta = \theta_0$$

$$P(T(X, \theta_0) \in T_0(\alpha) | \theta = \theta_0) = 1 - \alpha$$

$$H_1 = \theta \neq \theta_0, \theta > \theta_0, \theta < \theta_0$$

Пример:

$$1) X_1, \dots, X_n \sim F, \mu = EX_1, \exists \text{Var } X$$

$$H_0 : \mu = \mu_0$$

$$T(X) = \sqrt{n} \frac{\bar{X} - \mu_0}{S} \rightarrow N(0, 1) \text{ если } \mu = \mu_0$$

При $H_1 : \mu \neq \mu_0$ у нас двухсторонняя критическая область

При $H_1 : \mu > \mu_0$ у нас правосторонняя критическая область

При $H_1 : \mu < \mu_0$ у нас левосторонняя критическая область

Дальше для этого приводится пример с больницей (нам либо надо просто проверить, что температура не стандартная, либо нам важно знать, что она больше нормы, либо нам важно знать, что она ниже нормы)

$$P(\sqrt{n} \frac{\bar{X} - \mu_0}{S} \in T_0(\alpha) | \mu \neq \mu_0) = P(\sqrt{n} \frac{\bar{X} - \mu}{S} + \frac{\mu - \mu_0}{S} \sqrt{n} | \mu \neq \mu_0)$$

Это будет стремиться либо к $\Phi(-\infty)$ либо $1 - \Phi(+\infty)$

23 Критерии Колмогорова-Смирнова.

Критерий Колмогорова

$$X_1, \dots, X_n \sim F$$

$$H_0 : F = F_0 \text{ (} F_0 \text{ - непр)}$$

$$H_1 : F \neq F_0$$

Идея основана на теореме Колмогорова (было в начале семестра)

$$D_n = \sqrt{n} \sup_{x \in \mathbb{R}} |F_n(x) - F_0(x)|, F_n - \text{эмпирическая функция распределения}$$

if $D_n > q_{1-\alpha}$ then reject H_0 else accept H_0

1) $n \geq 20$ работает хорошо, при маленьких n есть спец таблицы

2) Так же есть приближенные формулы для D_n

$$3) H_0 : F = F(\theta), H_1 : \neg H_0 \Rightarrow D_n = \sqrt{n} \sup_x |F_n(x) - F_0(x, \theta)|$$

$$\theta \rightarrow \hat{\theta}$$

В пределе будет более сложная формула

Критерий Смирнова

X_1, \dots, X_n
 Y_1, \dots, Y_m
 Они независимы
 $H_0 : F_X = F_Y (= F_0)$
 $H_1 : \neq H_0$
 Тут идея основана на формуле Смирнова
 $D_{n,m} = \sqrt{\frac{nm}{n+m}} \sup_x |F_n(x) - F_m(x)|$
 $T_1(\alpha) = (q_{1-\alpha}, +\infty)$

24 Критерий согласия Пирсона хи-квадрат для простой и сложных гипотез

$X_1, \dots, X_n \sim F(x)$ - непрерывная
 Давайте дискретизируем данные
 $\Delta_1 : \nu_1$ - количество элементов выборки попадающих Δ_1
 \dots
 Δ_N

$$p_{\Delta_k} = \int_{\delta_k} p(x) dx, p(x) = F'(x)$$

Рассмотрим $\{1, 2, \dots, N\}$, $p = (p_1, \dots, p_N)$ - настоящий вектор вероятностей
 $p_0 = (p_{01}, \dots, p_{0N})$ - ожидаемый фиксированный вектор вероятностей
 ν_k - количество элементов в выборке типа k

$$H_0 : p = p_0$$

$$H_1 : p \neq p_0$$

$$n = \sum_{k=1}^N \nu_k$$

$$\Xi_N^2 = \sum_{k=1}^N \frac{(\nu_k - np_{0k})^2}{np_{0k}}$$

Теорема

$$\Xi_N^2 \xrightarrow{n \rightarrow \infty} \Xi^2(N-1) \text{ при условии } H_0$$

Доказательство

$N = 2 :$

$$\frac{(\nu_1 - np_{01})^2}{np_{01}} + \frac{(\nu_2 - np_{02})^2}{np_{02}} = \frac{(\nu_1 - np_{01})^2}{np_{01}} + \frac{(n - \nu_1 - n(1 - p_{01}))^2}{n(1 - p_{01})} = \frac{(\nu_1 - np_{01})^2}{n} \left(\frac{1}{p_{01}} + \frac{1}{1 - p_{01}} \right) = \frac{(\nu_1 - np_{01})^2}{np_{01}(1 - p_{01})}$$

Без квадрата по ЦПТ это стремится к $N(0, 1)$, но мы можем навесить непрерывную функцию возведения в квадрат и получим то, что нам нужно

Для классического критерий хи-квадрат у нас правосторонняя критическая область (потому что в сумме при H_0 будет маленькая разность в квадрате и тд и тп очев крч)

Критерий состоятельный (то есть вероятность ошибки второго рода стре-

мистия к единице, если погнать объем выборки на бесконечность)

Теорема

$$p_0(\theta) > 0 \forall \theta$$

$$\frac{\partial p_0}{\partial \theta}, \frac{\partial^2 p_0}{(\partial \theta)^2} - \text{непрерывная}$$

$$\text{Тогда } \Xi_N^2 \rightarrow X^2(N-1-d)$$

$$rk\left(\frac{\partial p_{0k}}{\partial \theta_j}\right)_{1 \leq k \leq N, 1 \leq j \leq d} = d$$

25 Критерий однородности хи-квадрат, в том числе случай 2×2

Пусть у нас k независимых выборок все они из $\{1, 2, \dots, N\}$

Пусть $p^{(j)}$ - истинный вектор вероятностей для соответствующей выборки

$$H_0 : p^{(1)} = \dots = p^{(k)}$$

$$H_1 : \neq H_0$$

ν_{ij} - количество элементов типа j в i -ой выборке

$$n_i = \sum_j \nu_{ij} = \nu_{i*} - \text{объем } i\text{-ой выборки}$$

$$n = n_1 + \dots + n_k$$

Пусть $p^{(1)}, \dots, p^{(k)}$ - известны

$$\Xi_{n_1}^2 = \sum_{j=1}^N \frac{(\nu_{ij} - n_i p_j^{(i)})^2}{n_i p_j^{(i)}}, df = N - 1$$

$$\Xi_{n_1, n_2, \dots, n_k}^2 = \sum_{i=1}^K \Xi_{n_i}^2, df = k(N - 1)$$

Рассмотрим $p^{(1)}, \dots, p^{(k)}$ - не известны

$$df = k(N - 1) - (N - 1) = (N - 1)(k - 1)$$

$$\hat{p}_j = \frac{\nu_{1j} + \dots + \nu_{kj}}{n} = \frac{\nu_{*j}}{n}$$

$$L(\dots) = p_1^{(\nu_{*1})} \dots p_k^{(\nu_{*k})}, n = \sum_j \nu_{*j}$$

$$Z_{n_1 n_2} = \left(\frac{\nu_{11}}{n_1} - \frac{\nu_{21}}{n_2} \right) \sqrt{\frac{n n_1 n_2}{\nu_{*1} \nu_{*2}}}$$

Утверждается, что $Z_{n_1 n_2}^2 = \Xi_{n_1 n_2}^2 \rightarrow N(0, 1)$

$$H_0 : p_1 = p_2, p_j = P(0 \text{ в } j\text{-ой выборке})$$

H_1 - настраивается

26 Критерий независимости хи-квадрат, в том числе случай 2×2

$$X_1, \dots, X_n : \{1, 2, \dots, N\}$$

$$Y_1, \dots, Y_n : \{1, 2, \dots, M\}$$

ν_{ij} - количество пар, в которых первая компонента равна i , а вторая j

Это можно представить в виде таблицы сопряженности

Просуммируем по каждому столбцу и по каждой строчке

Пусть $p_{ij} = P(X = i, Y = j)$

$p_{xi} = P(X = i)$

$p_{yj} = P(Y = j)$

$H_0 : p_{ij} = p_{xi}p_{yj} \forall i, j$

$H_1 : \neg H_0$

$$\Xi^2 = \sum_{1 \leq i \leq N, 1 \leq j \leq M} \frac{(\nu_{ij} - p_{ij}n)^2}{np_{ij}}, df = MN - 1$$

$$df = MN - 1 - (N - 1) - (M - 1) = MN - N - M + 1 = N(M - 1) - (M - 1) = (M - 1)(N - 1)$$

$$\widehat{p_{Xi}} = \frac{\nu_{i*}}{n}$$

$$\widehat{p_{Yj}} = \frac{\nu_{*j}}{n}$$

$$Z_n = \left(\frac{\nu_{11}}{\nu_{1*}} - \frac{\nu_{21}}{\nu_{2*}} \right) \sqrt{\frac{n\nu_{1*}\nu_{2*}}{\nu_{*1}\nu_{*2}}}$$

Утверждается, что $Z_n^2 = \Xi^2$

$Z_n = \sqrt{n}\rho_n$, ρ_n - выборочный коэффициент корреляции Пирсона (ковариация делить на корень из произведения дисперсий ыыыыыыы)

$$\rho = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E(XY) - EXEY}{\sigma_X \sigma_Y} = \frac{P(X = 1, Y = 1) - P(X = 1)P(Y = 1)}{\sqrt{P(X = 1)P(X = 0)P(Y = 1)P(Y = 0)}} =$$

... Тренер ну тут очев лол кек кд чд бро $= P(Y = 1|X = 1) - P(Y = 1|X = 0)$

27 Критерий квантилей и знаков

27.1 Критерий квантилей

$H_0 :$

$F(q_1) = \alpha_1$

$F(q_2) = \alpha_2$

...

$F(q_N) = \alpha_N$

где $\alpha_0 = 0 < \alpha_1 < \dots < \alpha_N < 1 = \alpha_{N+1}$

$H_1 : \neg H_0$

$q_0 = \inf \text{supp } P < q_1 < \dots < q_N < \sup \text{supp } P = q_{N+1}$

$\Delta_1 = [q_0, q_1)$

...

$\Delta_{N+1} = [q_N, q_{N+1})$

$P(\Delta_1) = \alpha_1 - \alpha_0$

...

$P(\Delta_N) = \alpha_{N+1} - \alpha_N$

27.2 Критерий знаков

Это особый случай $H_0 : F(q) = \frac{1}{2} (X_1, Y_1)^T, \dots, (X_n, Y_n)^T$

Мы хотим проверить что:

а) выборки независимы

б) распределения одинаковы

$$F(x, y) = F_1(x) \cdot F_1(y)$$

$$U = X - Y \Rightarrow \text{med } U = 0$$

ν_1 - количество элементов новой выборки $> \text{med}$

$$Z_n = \frac{2}{\sqrt{n}}(\nu_1 - \frac{n}{2}) \rightarrow N(0, 1)$$

$Z_n = \sqrt{n}\rho_n, \rho_n$ - коэффициент корреляции Пирсона

Теорема

$$(X_1, Y_1)^T, \dots, (X_n, Y_n)^T \sim N(\dots, \dots) \Rightarrow \frac{\sqrt{n-2}\rho_n}{\sqrt{1-\rho_n^2}} \sim T(n-2)$$

$$H_0 : \rho = 0$$

$$H_1 : (\rho \neq 0, \rho > 0, \rho < 0)$$

28 Ранговые критерии. Критерии Манна-Уитни-Уилкоксона

Определение. Ранг

X_1, \dots, X_n - выборка

$r(X_k)$ - номер X_k в вариационном ряде

X_1, \dots, X_n

Y_1, \dots, Y_m

Это две независимых выборки, давайте объединим их в одну

Рассмотрим $(X_1, \dots, X_n, Y_1, \dots, Y_m)$

R_i - ранг X_i , в объединенной выборке

$T = R_1 + \dots, R_n$ - Статистика Вилкоксона

$$Z_{rs} = \mathbb{1}(X_r < Y_s)$$

$$U = \sum_{r=1}^n \sum_{s=1}^m Z_{rs} \text{ - Мант-Уитни}$$

$$T + U = mn + \frac{n(n+1)}{2}$$

Хотим проверить, что распределение X совпадает с Y

$$EU = mnE\mathbb{1}(X < Y) = mnP(X < Y) = \frac{1}{2} \text{ - при условии } H_0$$

$$H_0 : P(X < Y) = \frac{1}{2}$$

$$U \sim N(\frac{mn}{2}, \frac{nm(m+n+1)}{12})$$

Теорема

$$(X_1, Y_1)^T, \dots, (X_n, Y_n)^T \sim N(\dots, \dots) \Rightarrow \frac{\sqrt{n-2}\rho_n}{\sqrt{1-\rho_n^2}} \sim T(n-2)$$

29 Коэффициенты корреляции Пирсона, Спирмена и Кендала. Статистические тесты, основанные на них.

Хотим проверить независимость

R_i - ранг X_i (в своей выборке)

S_i - ранг Y_i (в своей выборке)

ρ - выборочный коэффициент корреляции между R_i и S_i - коэффициент корреляции Спирмена

$$\rho = \frac{12}{n(n^2 - 1)} \sum (R_i - \frac{n+1}{2})(S_i - \frac{n+1}{2})$$

$$H_0 : \rho = 0, \sqrt{n}\rho \rightarrow N(0, 1)$$

$$H_1 : \rho \neq 0, \rho > 0, \rho < 0$$

Коэффициент корреляции Кендала

$$\tau = \frac{2}{n(n-1)} \sum_{i=1}^{n-1} \sum_{j=i+1}^n \text{sign}(T_j - T_i)$$

$$H_0 \text{ верно} \Rightarrow E\tau = 0, \text{Var } \tau = \frac{2(2n+5)}{9n(n-1)}$$

$$\tau \approx N(0, \frac{4}{9n})$$

$$H_0 : \tau = 0$$

$$H_1 : > 0, < 0, \neq 0$$

$$\rho = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E((X - \mu_X)(Y - \mu_Y))}{\sigma_X \sigma_Y}$$

Пирсон проверяет линейную зависимость между двумя случайными величинами

Спирмен проверяет монотонную зависимость между двумя случайными величинами

Кендал делает тоже самое что и Спирмен

30 Критерий инверсий

$X_{(1)} \leq \dots \leq X_{(n)}$ Крайние ситуации: выборка отсортирована, то есть трудно поверить, что у нас все случайно

ν_i - количество инверсий для элемента X_i

$$T = \nu_1 + \dots + \nu_n$$

$$ET = \frac{n(n-1)}{4}$$

$$\text{Var } T = \frac{n(n-1)(2n+5)}{72}$$

Статистика T - асимптотически нормальная

31 Модель линейной регрессии. Минимальные и обычные предположения. Оценка наименьших квадратов.

$$Y = Xb + \varepsilon$$

$X \in M_{n \times m}(\mathbb{R})$ - матрица переменных

x_{ij} - количественная переменная

$Y \in \mathbb{R}^n$ - наблюдения зависимой переменной

$b \in \mathbb{R}^m$ - неизвестный вектор коэффициентов

$\varepsilon \in \mathbb{R}^n$ - случайная ошибка

1) $E\varepsilon = 0$

2) $\text{Var } \varepsilon_i = \sigma^2$ - гомоскедотичность

3) $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$

Наша цель оценить b и σ^2

Определение. Ошибка наименьших квадратов

$$\hat{b} = \text{argmin } S^2(b)$$

$$S^2(b) = (Xb - Y)^T(Xb - Y)$$

$$A = X^T X \in M(m \times n)$$

$$\frac{1}{n} X^T X$$

$$\text{rank}(A) = m$$

Теорема

$$\hat{b} = A^{-1} X^T Y$$

Доказательство

$$S^2(\hat{b} + \delta) = S^2(\hat{b}) + \delta^T A \delta$$

$$t = Tb, T \in M_{k \times m}, k \leq m, \text{rank } T = k$$

$$\hat{t} = T\hat{b}$$

32 Теорема Гаусса-Маркова

1) \hat{t} - несмещенная оценка t

2) \hat{t} - оптимальная оценка в классе линейных несмещенных оценок

Доказательство

$$E\hat{t} = ET\hat{b} = ETA^{-1}X^TY = TA^{-1}X^TEY = TA^{-1}X^TXb = Tb = t$$

$$\text{Var}(\hat{t}) = \text{Var } TA^{-1}X^TY = TA^{-1}X^T \text{Var } YXA^{-1}T^T = \sigma^2 TA^{-1}T^T$$

$$\text{Var } \hat{t} = \sigma^2 TA^{-1}T^T$$

$$\text{Var } \hat{b} = \sigma^2 A^{-1}$$

Пусть LY - несмещенная оценка для t , т.е.

$$ELY = t = Tb \Rightarrow \text{Var } LY = \sigma^2 LL^T$$

$$\text{Заметим, что } LL^T = (TA^{-1}X^T)(TA^{-1}X^T)^T + (L - TA^{-1}X^T)(L - TA^{-1}X^T)^T$$

$$MSE\hat{t} = tr \operatorname{Var} \hat{t}$$

Тогда $L = TA^{-1}X^T$

33 Оценка остаточной дисперсии

$$ES^2(b) = E(Xb - Y)^T(Xb - Y) = E\varepsilon^T\varepsilon = n\sigma^2$$

$$E(\hat{b} - b)A(\hat{b} - b) = \sum_{i,j} a_{ij}E(\hat{b}_i - b_i)(\hat{b}_j - b_j) = \sum_{i,j} a_{ij} \operatorname{Cov}(\hat{b}_i, \hat{b}_j) = \sigma^2 \sum_{i,j} a_{ij}a_{ij}^{-1} =$$

$$\sigma^2 \sum_i \sum_j a_{ij}a_{ji}^{-1} = \sigma^2 \sum_{i=1}^m 1 = \sigma^2 m$$

$$S^2(b) = S^2(\hat{b}) + (\hat{b} - b)^T A(\hat{b} - b) \Rightarrow n\sigma^2 = ES^2(\hat{b}) + m\sigma^2 \Rightarrow \hat{\sigma}^2 = \frac{S^2(\hat{b})}{n - m} -$$

несмещенная оценка для σ^2

34 Условная оценка наименьших квадратов

$$Tb = t_0, T \in M_{k \times m}, \operatorname{rank} T = k$$

$$\hat{b}_{T,t_0} = \operatorname{argmin}_{Tb=t_0} S^2(b) - \text{условная оценка наименьших квадратов}$$

Теорема

$$\hat{b}_{T,t_0} = \hat{b} - A^{-1}T^T D^{-1}(T\hat{b} - t_0), D = TA^{-1}T^T$$

Можно показать: $S^2(b) = S^2(\hat{b}_{T,t_0}) + (\hat{b}_{T,t_0} - b)^T A(\hat{b}_{T,t_0} - b)$

35 Основная теорема о линейной регрессии. Следствия. t-тест.

$$1) S^2(\hat{b}), \hat{b} - \text{незав.}$$

$$S^2(b) - S^2(\hat{b}), S^2(\hat{b}) - \text{незав.}$$

$$2) \hat{b} \sim N(b, \sigma^2 A^{-1})$$

$$3) \frac{S^2(\hat{b})}{\sigma^2} \sim \Xi(n - m)$$

$$\frac{S^2(b) - S^2(\hat{b})}{\sigma^2} \sim \Xi^2(m)$$

Без доказательства

T-test:

$$\frac{\frac{\hat{b}_j - b_j}{\sigma \sqrt{A_{ij}^{-1}}}}{\sqrt{\frac{1}{n - m} \frac{S^2(\hat{b})}{\sigma^2}}} \sim T(n - m)$$

$$\sqrt{\frac{1}{n-m} \frac{S^2(\hat{b})}{\sigma^2}} = \frac{\hat{b}_j - b_j}{S(\hat{b} \sqrt{A_{ij}^{-1}})} \sqrt{n-m}$$

$$H_0 : b_j = b_{0j}$$

$$H_1 : <, \neq, >$$

36 F-тест. Коэффициент детерминации.

$$T \in M_{k \times m}$$

$$H_0 : Tb = t_0 \text{ (default: } T = E_m, t_0 = 0, \text{ то есть все } b_i = 0 \text{ одновременно)}$$

$$H_1 : Tb \neq t_0$$

$$F = \frac{n-m}{k} \cdot \frac{S^2(\hat{b}_{T, t_0}) - S^2(\hat{b})}{S^2(\hat{b})} = \frac{n-m}{k} \cdot \frac{(T\hat{b} - t_0)^T D^{-1} (T\hat{b} - t_0)}{S^2(\hat{b})} \sim F(k, n-m)$$

$$R = \frac{\sum_i (Y_i - \bar{Y})(\hat{Y}_i - \bar{\hat{Y}})}{\sqrt{\sum_i (Y_i - \bar{Y})^2 \sum_i (\hat{Y}_i - \bar{\hat{Y}})^2}}, \hat{Y} = X\hat{b}$$

R - многомерный коэффициент корреляции

R^2 - коэффициент детерминации

37 Однофакторный дисперсионный анализ

$$Y_{ij} = \mu_j + \varepsilon_{ij}, \mu_j - \text{среднее влияние на } j\text{-ом факторе}$$

$$1 \leq j \leq \nu - \text{уровень фактора}$$

$$1 \leq i \leq I_j - i\text{-ое наблюдение для фактора с уровнем } j$$

$$I = \sum_{j=1}^{\nu} I_j$$

$$\varepsilon_{ij} \sim N(0, \sigma^2 E)$$

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_J$$

$$H_1 : \neg H_0$$

$$F = \frac{\frac{S_B^2}{df_B}}{\frac{S_W^2}{df_W}} \sim F(J-1, I-J)$$

$$S_B^2 = \sum_{j=1}^J I_j (\bar{Y}_{*j} - \bar{Y})^2 - \text{межгрупповая дисперсия}$$

$$df_B = J - 1$$

$$S_W^2 = \sum_i \sum_j (Y_{ij} - \bar{Y}_{*j})^2 - \text{внутригрупповая дисперсия}$$

$$df_W = I - J$$

$$S^2 = S_W^2 + S_B^2$$

38 Двухфакторный дисперсионный анализ.

$$Y_{ij} = \mu + a_i + b_j + \varepsilon_{ij}$$

$$1 \leq i \leq I$$

$$1 \leq j \leq J$$

Для каждой пары (i, j) ровно одно наблюдение

M - общее среднее

a_i - среднее влияние на i-ом уровне фактора 1

b_j - среднее влияние на j-ом уровне фактора 2

Факторы независимы

$$\sum a_i = \sum b_j = 0$$

ε_{ij} - незав.

$$\varepsilon_{ij} \sim N(0, \sigma^2)$$

$$H_A : \forall a_i = 0 \rightarrow F_A = \frac{(I-1)(J-1)}{I-1} \frac{S_A^2}{S_{ALL}^2}$$

$$H_B : \forall b_i = 0 \rightarrow F_B = \frac{(I-1)(J-1)}{J-1} \frac{S_B^2}{S_{ALL}^2}$$

$$H : \forall a_i = b_j = 0 \rightarrow F = \frac{(I-1)(J-1)}{I+J-2} \frac{S_A^2 + S_B^2}{S_{ALL}^2}$$

$$S_A^2 = J \sum (\bar{Y}_{i*} - \bar{Y})^2$$

$$S_B^2 = I \sum (\bar{Y}_{*j} - \bar{Y})^2$$

$$S_{ALL}^2 = \sum_{i,j} (Y_{ij} - \bar{Y}_{i*} - \bar{Y}_{*j} + \bar{Y})^2$$

39 Ковариационный анализ.

1 фактор, 1 количественная переменная

$$y_{ij} = \beta_i + \gamma z_{ij} + \varepsilon_{ij}$$

$$1 \leq i \leq I$$

$$1 \leq j \leq J$$

$$n = I \cdot J$$

На каждой паре (i, j) одно наблюдение

β_i - среднее влияние на уровне i

z_{ij} - количественная переменная

ε_{ij} - нез.

$$\varepsilon_{ij} \sim N(0, \sigma^2)$$

$$H_\gamma : \gamma = 0$$

$$\hat{\beta}, \hat{\gamma} - \text{ОНК} \Rightarrow \sqrt{n-I-1} \frac{\hat{\gamma}}{\sqrt{A_{\gamma\gamma}^{-1} S^2(\hat{\beta}, \hat{\gamma})}} \sim T(n-I-1)$$

$$A_{\gamma\gamma} = A_{(I+1)(I+1)}$$

$$S^2(\hat{\beta}, \hat{\gamma}) = \sum_{ij} (y_{ij} - \hat{\beta}_i - \hat{\gamma} z_{ij})^2 - \text{суммарная квадратичная ошибка}$$

$$H_B : \beta_1 = \dots = \beta_I$$

$$F = \frac{n - I - 1}{I - 1} \cdot \frac{\sum_{i,j} y_{ij}^2 - n\bar{y}^2 - p}{p}$$

$$p = \sum_{i,j} y_{ij}^2 - n\bar{y} - \frac{(\sum_{i,j} z_{ij} y_{ij} - n\bar{z})^2}{\sum_{i,j} z_{ij}^2 - n\bar{z}}$$

40 Логистическая регрессия

$$Y_i \approx \frac{1}{1 + e^{-x_i c}}, X = (X_1 \dots X_n)^T, c = (c_1 \dots c_m)^T$$

$$Y_i \sim \text{Bern}\left(\frac{1}{1 + e^{-x_i c}}\right)$$

$$L(X, c) = \prod_{i=1}^n \left(\frac{1}{1 + e^{-x_i c}}\right)^{Y_i} \left(1 - \frac{1}{1 + e^{-x_i c}}\right)^{1-Y_i}$$

$$-\ln L(X, c) = \dots$$

41 Последовательный анализ Вальда. Описание алгоритма. Теорема о конечном числе шагов для простых гипотез.

$$H_0 : F = F_0$$

$$H_1 : F = F_1 \quad f_0, f_1 - \text{соответствующие плотности} \quad A_0 < 1 < A_1$$

$$n = 1$$

$$\text{while true } \{$$

$$\text{if } \frac{L_{0n}}{L_{1n}} \leq A_0 \text{ then return reject } H_0$$

$$\text{if } \frac{L_{0n}}{L_{1n}} \geq A_1 \text{ then return accept } H_0$$

$$++ n \}$$

$$\alpha, \beta, a_0 = \ln A_0, a_1 = \ln A_1$$

$$Z_i = \ln \frac{f_1(X_i)}{f_0(X_i)}$$

$$\sum_{i=1}^n Z_i = \ln \frac{L_{1n}}{L_{0n}}$$

Теорема

Пусть ν - количество итераций

$$\Rightarrow P(\nu \geq n | H_0) \rightarrow 0, P(\nu \geq n | H_1 \rightarrow 0), n \rightarrow \infty$$

Доказательство

$r - fixed$

$$\begin{aligned}
& \eta_1 = Z_1 + \dots + Z_r \\
& \eta_1 = Z_{r+1} + \dots + Z_{2r} \\
& \dots \{ \nu \geq rk \} \Leftrightarrow \{ a_0 < \eta_1 + \dots + \eta_j < a_1 \}_{i \leq j \leq k} = \{ a_0 < \eta_1 < a_1, a_0 < \eta_1 + \eta_2 < a_1, \dots, a_0 < \eta_1 + \dots + \eta_k < a_1 \} \subset \{ |\eta_j| < b = a_i - a_0 \}_{1 \leq j \leq k} \Rightarrow P(\nu \geq rk | H_s) \leq \\
& P_{1 \leq j \leq k}(|\eta_j| < b | H_s) = P_{1 \leq j \leq k}(\nu_j^2 < b^2 | H_s) = P^k(\eta_1^2 < b^2 | H_s) = p_s^k \\
& E_s \eta_1^2 \geq \text{Var}_s \eta_1 = r \text{Var}_s Z_1 > b^2, r > \max\left(\frac{b^2}{\text{Var}(Z_1 | H_0)}, \frac{b^2}{\text{Var} Z_1 | H_1}\right) \stackrel{E_s(\dots) = E(\dots | H_s)}{\text{Var}_s(\dots) = \text{Var}(\dots | H_s)} \Rightarrow \\
& p + s < 1 \\
& P(\nu > n | H_s) \leq P(\nu \geq rk | H_s) \leq p_s^k \rightarrow 0, k = k(n), n \rightarrow \infty
\end{aligned}$$

42 Последовательный анализ Вальда. Теорема об оценке граничных констант и вероятностей ошибок первого и второго рода. Оценка среднего числа итераций.

Теорема

Пусть (α, β) - вероятности $\Rightarrow A_0 \geq A'_0 = \frac{\beta}{1-\alpha}$ и $A_1 \leq A'_1 = \frac{1-\beta}{\alpha}$

if A'_0 и $A_1 = A'_1$ then α', β' - вероятности ошибок

$$\alpha' \leq \frac{\alpha}{1-\beta}, \beta' \leq \frac{\beta}{1-\alpha}$$

$$\alpha' + \beta' \leq \alpha + \beta$$

Доказательство

$$1 = \sum_{n=1}^{\infty} P(\nu = n | H_0) = \sum_{n=1}^{\infty} P(\text{accept at the step } n | H_0) + \sum_{n=1}^{\infty} P(\text{reject at the step } n | H_0) =$$

$$\alpha = \sum_{n=1}^{\infty} P(\text{reject at the step } n | H_0) \leq \frac{1}{A_1} \sum_{n=1}^{\infty} P(\text{reject at the step } n | H_1)$$

$$L_{n1} \geq L_{n0} \cdot A_1$$

$$L_{n1} \geq L_{n0} \cdot A_1$$

Аналогичным образом мы можем получить неравенство:

$$\beta \leq A_0(1-\alpha)$$

$$\alpha \leq \frac{1}{A_1}(1-\beta)$$

Пусть $A'_0 = \frac{\beta}{1-\alpha}$ и $A'_1 = \frac{1-\beta}{\alpha}$ и ошибки $(\alpha', \beta') \Rightarrow \frac{\beta}{1-\alpha} \geq \frac{\beta'}{1-\alpha'}$

$$\frac{1-\beta}{\alpha} \leq \frac{1-\beta'}{\alpha'} \Rightarrow \alpha' \leq \frac{\alpha(1-\beta')}{1-\beta} \leq \frac{\alpha}{1-\beta}$$

$$\beta(1-\alpha) \leq \beta(1-\alpha')$$

$$(1-\beta)\alpha' \leq (1-\beta')\alpha$$

Сложим последние два неравенства и получим требуемое

Давайте вспомним тождество Вальда:

(X_i) - независимые случайные величины с математическим ожиданием a

ν - целочисленная случайная величина не зависящая от X_i

$$E\nu = b$$

$$S_\nu = X_1 + \dots + X_\nu$$

$$\Rightarrow ES_\nu = ab = E(ES_\nu|\nu) = E\nu a = aE\nu = ab$$

$$a_0 = \ln \frac{\beta}{1-\alpha}, a_1 = \ln \frac{1-\beta}{\alpha}$$

$$E_s \nu \cdot E_s Z_i = E_s S_\nu \approx a_0(1-\alpha) + a_1(1-\beta)$$

$$E_1 \nu \cdot E_1 Z_i \approx \beta a_0 + a_1(1-\alpha)$$

43 Бутстреп. Оценка дисперсии, построение д.и., permutation test.

1. $X_1, \dots, X_n \sim F_\theta; \hat{\theta}; \text{Var } \hat{\theta} - ?$ Можно ли ее оценить численно?
2. $X_1, \dots, X_n; \bar{X}; \text{Var } \bar{X} = \frac{\sigma}{\eta}$ можно ли оценить численно?
3. $X_1, \dots, X_n; \text{med } X; \text{Var med}(X) - ?$

$X_1, \dots, X_n; \hat{\theta}$ - оценка чего-то; F_n - Э.Ф.П

for (int i = 0; i <= B - 1; ++i) {

$$X_1^*, \dots, X_n^* \sim F_n$$

$$\theta_i^* = g(X_1^*, \dots, X_n^*)$$

$$\theta^*.append(\theta^*)$$

}

$$\text{Var}^* = \frac{1}{B} \sum_{i=0}^{B-1} (\theta_i^* - \bar{\theta}^*)^2$$

$$\frac{\text{Var}^*}{\widehat{\text{Var } \hat{\theta}}} \xrightarrow[n \rightarrow \infty]{P} 1$$

43.1 CI

$$\widehat{F}(t) = \frac{1}{B} \sum_{j=1}^B \mathbb{1}(\sqrt{n}(\bar{\theta}_j^* - \hat{\theta}) \leq t)$$

$$\hat{\theta} \pm \frac{q_{1-\frac{\alpha}{2}} - \frac{\alpha}{2}}{\sqrt{n}}; q_{1-\frac{\alpha}{2}} - \text{квантиль } \hat{F}$$

$$\hat{\theta} - \frac{q_{\frac{\alpha}{2}}}{\sqrt{n}}; \hat{\theta} + \frac{q_{1-\frac{\alpha}{2}}}{\sqrt{n}}$$

43.2 The permutation test

$$X_1, \dots, X_n \sim F$$

$$Y_1, \dots, Y_m \sim G$$

Both in dependent

$H_0 : F = G$
 $H_1 : F \neq G$
 $W = (X_1, \dots, X_n, Y_1, \dots, Y_m)$
 $X_1 \dots X_n Y_1 \dots Y_m$
 $W_1 \dots W_{n+m}$
 Пусть $T = T(X, Y)$ - некая статистика
 $T = |\bar{X} - \bar{Y}|$
 $T(W, Z) = |\overline{(Z, W)} - \overline{Y(Z, W)}|$
 $X(Z, W) = \{W_i : Z_i = 1\}$
 $Y(Z, W) = \{W_i : Z_i = 2\}$
 Пусть Z^* - перестановка Z
 for i in range(B):
 Z^* - перестановка
 $T^* = T(Z^*, W)$
 $T_{all}.append(T^*)$
 $p = \frac{1}{B} \sum_{j=1}^B \mathbb{1}(T_{all}[j] > t), t = T(W, Z)$
 $p = p_{value}$

44 Введение в байесовскую статистику. Credible intervals. Проверка гипотез.

44.1 Credible interval

$$\begin{aligned}
 P(a_l < \theta < a_r | X) &= 1 - \alpha \\
 \theta &\sim \Gamma(\theta) \\
 p(\theta | X) &= \frac{p(X|\theta) \cdot \pi(\theta)}{\int p(X|\theta) \pi(\theta) d\theta}
 \end{aligned}$$

44.2 Hypothesis. Bayesian

$$\begin{aligned}
 H_0 : \theta &= \theta_0 - p_0 - prior \\
 H_1 : \theta &= \theta_1 - p_1 - prior \\
 L(\theta_0 | X) &= \frac{L(X|\theta_0) \cdot p_0}{const} \\
 const &= P(X) = L(X|\theta_0) \cdot p_0 + L(X|\theta_1) \cdot p_1 \\
 L(\theta_1 | X) &\approx L(X|\theta_1) \cdot p_1 \\
 Y &= Xc + \varepsilon, C \sim N(.,.) \\
 \varepsilon &\sim N(0, \sigma^2 I) \\
 (c, \varepsilon) &\sim (.,.) \\
 (Y|C) &\sim N(Xc, \sigma^2 I) \Rightarrow L(c|Y) \sim N(.,.) \\
 \hat{c} &= E(c|Y)
 \end{aligned}$$