

χ^2 тест сапастности са расушеном

У овом тестирању $H_0: F=F_0$ vs $H_1: F \neq F_0$, при чему F не мора бити аус. нпр.

Идеја: 1. пожељно домени одређења у k дисјунктних категорија B_1, \dots, B_k
 \downarrow
скуп свих могућих вредности које одређење може узети ако знамо или расушено F_0 (шј. ако боде H_0)

2. одредимо колико је путаке из узорка ушло у j -ту категорију и тај број означимо са m_j

(4) 3. Ако је H_0 тачна сл. вел M_j која представља број X_i -ева који су узети вредности из j -те категорије има $\text{Bin}(n, p_j)$, $p_j = P\{X \in B_j \mid H_0 \text{ тачна}\}$

$M_j = \sum_{i=1}^n I_k$
 $I_k = I\{X_k \in B_j\}$
 I_1, \dots, I_n су независни!

Зашто ћемо посматрати колико се свако m_j разликује од np_j , што је очекивање од M_j ако важи H_0
Ако је H_0 тачна, m_j је блиско $\mathbb{E} M_j = np_j$

Сада смо на корак од формирања тест статистике.
Поменимо прво једну теорему:

Т Турсовска теорема

Нека је $p_j = P\{X \in B_j\}$, $j \in \{1, \dots, k\}$, при чему
 $B_i \cap B_j = \emptyset$ за $i \neq j$, $\sum_{j=1}^k p_j = 1$

Нека је (X_1, \dots, X_n) нсу, где $X_i \stackrel{d}{=} X$, n велико



Нека је M_j др. X_i -ева коју су узели вредности из B_j .
Тада $M_j \sim \text{Bin}(n, p_j)$, $E M_j = n p_j$ и:

$$\sum_{j=1}^k \frac{(M_j - n p_j)^2}{n p_j} \sim \chi^2_{k-1}$$

Дакле, за тестирање

$$H_0: F = F_0 \quad \text{vs} \quad H_1: F \neq F_0$$

када је одим узорка велики можемо користити исти стат.

$$T = \sum_{j=1}^k \frac{(M_j - n p_j)^2}{n p_j}, \quad T \sim \chi^2_{k-1} \quad \text{по } H_0$$

Где су M_j и p_j дефинисани у кораку 3. (ск)

У случају H_1 угу бете вредности $(M_j - n p_j)^2$, то
и бете вредности T , то је критичка област:

$$W = \{T \geq c\}$$

- категорије можемо груписати произвољно с тим
што се саветује да j -ту категорију спојимо
са неком другом ако је $n p_j < 5$

- Ако F_0 зависи од непознатих параметара, обично*
их оценимо методом максималне вероватноће.

Оцене тих параметара искористићемо за рачунање
истог стат. T , која је сада таква да:

$$T \underset{\text{по } H_0}{\sim} \chi^2_{k-1} - \text{др. оценок параметара}$$

* није теоријски исправно, али се ради у пракси

51.

Видимо је $P\{X=k\} = 2^{-k} = \frac{1}{2^{k-1}} \cdot \frac{1}{2}$, $k \in \mathbb{N}$, одговара за-

коњу геометријске расподеле са параметром $\frac{1}{2}$

Дакле, пошредно је да изјавимо следеће хипотезе:

$$H_0: F = F_0 \quad \text{vs} \quad H_1: F \neq F_0$$

где је F_0 Φ -ја $G(\frac{1}{2})$ расподеле. (у овом задатку
нам F_0 не треба, наводимо само како гласе H_0 и H_1)

Шта ћемо из даље падење у задатку?

Видимо да се јединица у узорку појавила 45 пута,
двојка 30 пута, тројка 15 пута и четворка 10 пута,

Ово падење нам може послужити да одредимо ка-
тегорије у које ћемо разврстати паде из
узорка.

Природно би било да за категорије узмемо ску-
пове $\{1\}$, $\{2\}$, $\{3\}$ и $\{4\}$, али као што смо рекли,
укупна свих категорија мора покривати цео домен
однежане X ако она зависи има $G(\frac{1}{2})$ расподелу.

Како сп. величина из $G(\frac{1}{2})$ расподеле узима вред-
ности из скупа \mathbb{N} , можемо поделу на категорије
извршити на следећи начин:

$$B_1 = \{1\}, B_2 = \{2\}, B_3 = \{3\}, B_4 = \{4, 5, \dots\}$$

(није грешка погледати B_i као интервале у збирци,
али ово је мало природније..)

- Дакле наше категорије - B_j и бројеви појављивања из узорка које су и наше y класе - m_j су:

B_j	$\{1\}$	$\{2\}$	$\{3\}$	$\{4, 5, \dots\}$
m_j	45	30	15	10

- Да бисмо одредили шестиперцентну, потврду је да израчунамо p_1, p_2, p_3 и p_4 .

$$p_1 = P\{X \in B_1 | H_0 \text{ важи}\} = P\{X=1 | X \sim G(\frac{1}{2})\} = \frac{1}{2^1}$$

$$\text{слично, } p_2 = \frac{1}{2^2}, p_3 = \frac{1}{2^3}.$$

$$p_4 \text{ можемо добити као } 1 - \frac{1}{2} - \frac{1}{2^2} - \frac{1}{2^3} = 1 - \frac{7}{8}$$

- Требао би да се уверимо да ли је закључак $n p_j \geq 5$ (ако није, спајемо неке категорије):

$$n = m_1 + m_2 + m_3 + m_4 = 100$$

$$\Rightarrow n p_1 = 50, n p_2 = 25, n p_3 = 12,5, n p_4 = 12,5$$

- Дакле, категорије не треба спајати.

- Израчунајмо сада шестиперцентну:

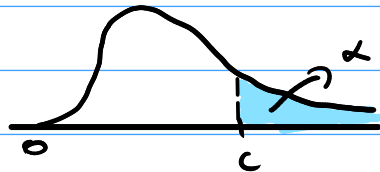
$$T = \sum_{j=1}^4 \frac{(m_j - n p_j)^2}{n p_j} \Rightarrow \text{реализоване вредности } T \text{ је:}$$

$$\frac{(45 - 50)^2}{50} + \frac{(30 - 25)^2}{25} + \frac{(15 - 12,5)^2}{12,5} + \frac{(10 - 12,5)^2}{12,5}$$

$$= 2.5$$

Одређујемо сада критичну област. Рекли смо да је она
облика $W = \{T \geq c\}$.

c је и.г. $P\{T \geq c \mid H_0 \text{ важи}\} = \alpha$, а како знамо да
 $T \sim \chi^2_{k-1} = \chi^2_{n-1}$, за c важи:



$$F_{\chi^2_3}(c) = 1 - \alpha$$

$$\Rightarrow c = F_{\chi^2_3}^{-1}(1 - \alpha)$$

$$\text{у R-у: } c = qchisq(1 - 0,05, 3)$$

↓
α код
H0

52. Ако са X означимо оделенке које представља
одржавање пречника осовине и са F означимо
 ϕ -ју раширене оделенке, пошредно је интерпретаци:

$$H_0: F = F_0 \quad \text{vs} \quad H_1: F \neq F_0$$

где је F_0 ϕ -ја $N(m, \sigma^2)$ раширене

Видимо да у овом задатку F_0 зависи од непознатих
параметара, па је пошредно да их одређимо.

Рекли смо да ћемо непознате параметаре одређивати
методом макс. веродостојносног, па су изражене оне:

$$\hat{m}_{mnv} = \bar{X}_n \quad \text{и} \quad \hat{\sigma}_{mnv}^2 = \bar{S}_n^2$$

Међутим, ове ситуације не можемо изразити на основу података које имамо, па ћемо користити модификоване оцене \bar{X}_n^* и \bar{S}_n^{2*} , које се добијају као узорак средина и узорак дисперзија узорака сачињених од средина датих интервала.

Нај узорак код нас изгледа овако:

$$\underbrace{(2.5, 2.5, \dots, 2.5)}_{15 \times}, \underbrace{(7.5, \dots, 7.5)}_{75 \times}, \dots, \underbrace{(22.5, \dots, 22.5)}_{20 \times}$$

$$\Rightarrow \bar{X}_n^* = \frac{15 \cdot 2.5 + 75 \cdot 7.5 + \dots + 20 \cdot 22.5}{15 + 75 + \dots + 20} = 12,2$$

$$\bar{S}_n^{2*} = 25,4$$

Напоменуто да овакав поступак нема никакво теоријско оправдање користимо га јер је једноставан и зато што немамо много других опција. (За радозналике које занима наложба ММВ агеа на основу броја шалака у интервалима: видети ЕМ алгоритам)

Сада када имамо оцене параметара μ и σ^2 , можемо наставити задатак. Следећи корак је одређивање категорије.

Било би природно за категорије узети дате интервале одступања, али претходно продохимо $(0,5)$ на $(-n,5)$ и $(20,25)$ на $(20,n)$, како бисмо покрили цео \mathbb{R} , што је датум $N(\mu, \sigma^2)$ расподеле.

Скенирајмо рогу једноставности исхода са B_j и m_j :

B_j	$(-\infty, 5)$	$[5, 10)$	$[10, 15)$	$[15, 20)$	$[20, +\infty)$
m_j	15	75	100	50	20

Када θ_0 зависи од неопознатих параметара, не можемо га одредити p_j . Ево и зашто:

$$p_1 = P\{X \in B_1 \mid \theta_0 \text{ непознато}\} = P\{X < 5 \mid X \sim N(\mu, \sigma^2)\}$$

$$= P\left\{\frac{X - \mu}{\sigma} < \frac{5 - \mu}{\sigma} \mid X \sim N(\mu, \sigma^2)\right\} = \Phi\left(\frac{5 - \mu}{\sigma}\right)$$

Слику, p_2, p_3, p_4 и p_5 ће зависити од μ и σ . Дакле, сада морамо одређити p_j са \hat{p}_j , користећи годујерне оцене $\hat{\mu}_{mnv}$ и $\hat{\sigma}_{mnv}^2$.

Ипак годујано га је:

$$\hat{p}_1 = \Phi\left(\frac{5 - \hat{\mu}_{mnv}}{\hat{\sigma}_{mnv}}\right) = 0,98$$

$$\hat{p}_2 = \Phi\left(\frac{10 - \hat{\mu}_{mnv}}{\hat{\sigma}_{mnv}}\right) - \Phi\left(\frac{5 - \hat{\mu}_{mnv}}{\hat{\sigma}_{mnv}}\right) = 0,25$$

$$\vdots$$

$$\hat{p}_5 = 1 - \Phi\left(\frac{20 - \hat{\mu}_{mnv}}{\hat{\sigma}_{mnv}}\right) = 0,96$$

Можемо проверити да је $np_j \geq 5$, $j \in \{1, \dots, 5\}$ и да није погрешно приписати категорије

Оцене \hat{p}_j ћемо сад искористити за рачунање статистике:

$$T = \sum_{j=1}^5 \frac{(M_j - n\hat{p}_j)^2}{n\hat{p}_j},$$

која сад под H_0 има χ^2_{5-1-2} расподелу
још имамо 2 параметра

Знамо да је $n = m_1 + \dots + m_5 = 260$, ба знамо све да одређимо реализовану вредност од T :

$$\frac{(15 - 260 \cdot 0,07)^2}{260 \cdot 0,07} + \dots + \frac{(20 - 260 \cdot 0,06)^2}{260 \cdot 0,06} = 6,017$$

Критична област је $W = \{T > c\}$, где је с.в.г.

$$P\{T > c \mid H_0 \text{ вистина}\} = \alpha, \text{ где } P\{T > c \mid T \sim \chi^2_2\} = \alpha$$

$$\Rightarrow F_{\chi^2_2}(c) = 1 - \alpha \Rightarrow c = F_{\chi^2_2}^{-1}(1 - \alpha)$$

$$\text{у R-у: } c = \text{qchis2}(1 - 0,05, 2)$$

$$c = 5,991$$

↓
тамо α