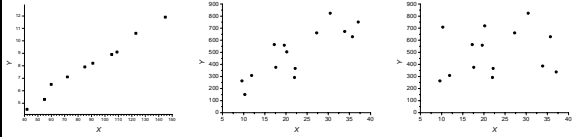


## Korelacije

Razni stupnjevi rasipanja:



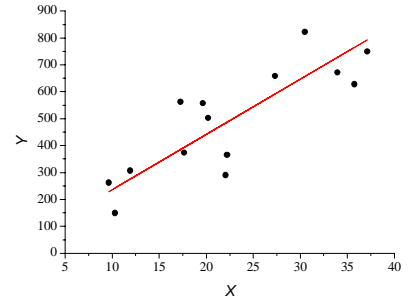
Proučavamo srednju sliku:

ovisnost Y o X:

$$y = ax + b$$

$$a = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2}$$

$$b = \frac{\overline{x^2}\bar{y} - \bar{x}\overline{xy}}{\overline{x^2} - \bar{x}^2}$$



$$a = \frac{m_{11} - m_{10}m_{01}}{m_{20} - m_{10}^2} = \frac{M_{11}}{M_{20}} = \frac{\sigma_{XY}}{\sigma_X^2}$$

kovarijanca :

$$M_{11} = \sigma_{XY} = \frac{1}{n} \sum_i (x_i - \bar{x})(y_i - \bar{y})$$

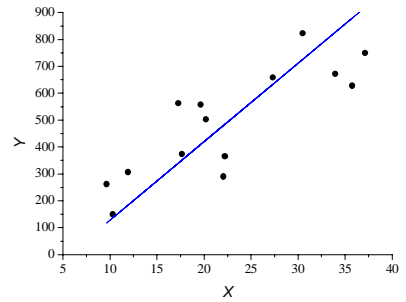
$$M_{11} = \frac{1}{n} \sum_i x_i y_i - \bar{x}\bar{y} = \frac{1}{n} \sum_i x_i y_i - \frac{1}{n^2} \sum_i x_i \sum_i y_i$$

Ovisnost X o Y:

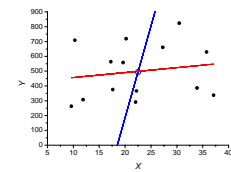
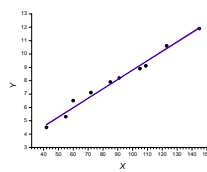
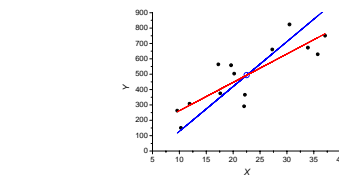
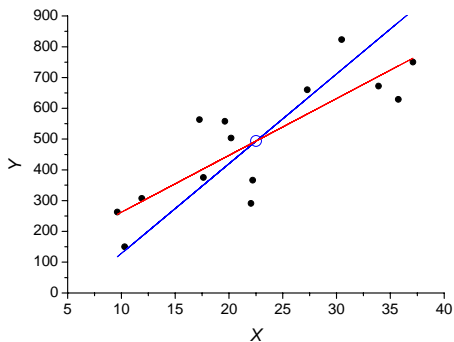
$$x = cy + d$$

$$c = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{y^2} - \bar{y}^2}$$

$$d = \frac{\overline{y^2}\bar{x} - \bar{y}\overline{xy}}{\overline{y^2} - \bar{y}^2}$$



$$c = \frac{m_{11} - m_{10}m_{01}}{m_{02} - m_{01}^2} = \frac{M_{11}}{M_{02}} = \frac{\sigma_{XY}}{\sigma_Y^2}$$



**Jednadžbe pravaca regresije pomoću kovarijance**

Ovisnost Y o X:

jednadžba pravca:  $(y - \bar{y}) = \frac{\sigma_{XY}}{\sigma_X^2} (x - \bar{x})$

Ovisnost X o Y:

jednadžba pravca:  $(x - \bar{x}) = \frac{\sigma_{XY}}{\sigma_Y^2} (y - \bar{y})$

standardizirane jednadžbe pravaca:

$$\frac{(y - \bar{y})}{\sigma_Y} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y} \frac{(x - \bar{x})}{\sigma_X}$$

$$\frac{(x - \bar{x})}{\sigma_X} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y} \frac{(y - \bar{y})}{\sigma_Y}$$

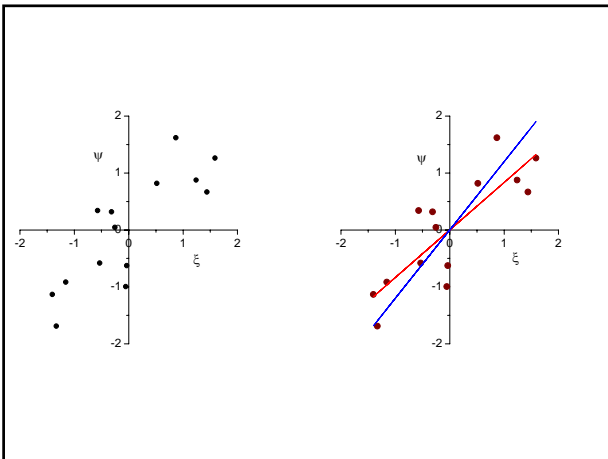
bezdimenzionalne varijable:  $\xi = \frac{(x - \bar{x})}{\sigma_X}$  i  $\psi = \frac{(y - \bar{y})}{\sigma_Y}$

novi pravci regresije:  $\psi = \rho \xi$        $\xi = \rho \psi$

Def.: **koeficijent korelacije:**

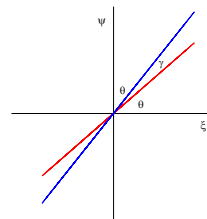
$$\rho = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

bezdimenzionalno  
ne ovisi o jedinicama



**Značenje koeficijenta korelacije**

poprima vrijednosti:  $-1 \leq \rho \leq 1$



$$\rho = \text{tg } \theta$$

$$\gamma = \frac{\pi}{2} - 2\theta$$

$$\text{tg } \gamma = \frac{1 - \rho^2}{2\rho}$$

X i Y nezavisni  $\Rightarrow \gamma = \frac{\pi}{2}$  ,  $\rho = 0$

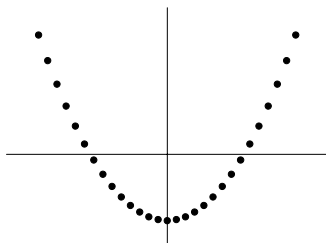
X i Y linearno deterministički vezani  $\Rightarrow \gamma = 0$  ,  $\rho = -1$  ili  $\rho = 1$

Dogovor: dobra korelacija  $\rightarrow |\rho| > 0,5$

**Napomena:**

Kad su X i Y nezavisne varijable, koeficijent korelacije težiti će nuli ( $\rho \approx 0$ ). Međutim, činjenica da je  $\rho = 0$  ne znači nužno da su varijable nezavisne.

Primjer:



**Provjeravanje (testiranje) hipoteza**

Provjera hipoteze vrlo je bitan dio statističkog zaključivanja.

Da bi se takva provjera formulirala, potrebno je postaviti neku teoriju koja se želi dokazati.

Npr.:

- Novi lijek bolji je za liječenje određenih simptoma od starog.
- Igrača kocka ima pomaknuto težište
- U svijetu se rađa više muškaraca nego žena

U svakom takvom problemu postavljamo dvije tvrdnje (hipoteze) od kojih je točno jedna istinita:

$H_0$  = nul-hipoteza

$H_1$  = alternativna hipoteza

One se **ne** tretiraju **ravnopravno**:

U sudskom procesu:  $H_0$ : optuženi je nevin

$H_1$ : optuženi je kriv

Hipoteza  $H_0$  smatra se ispravnom dok se ne dokaže  $H_1$ .

Za  $H_0$  obično se uzima stara, postojeća teorija:

Npr., za lijek:  $H_0$  = stari lijek jednako je dobar kao i novi

za kocku:  $H_0$  = kocka je poštena

**Ako ne odbacimo  $H_0$ , to ne znači da je ona ispravna, nego samo da nemamo dovoljno dokaza da ju odbacimo.**

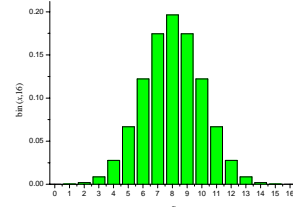
Primjer: ispitujemo simetričnost Galtonove daske sa 16 redova.

### 1. Postavljanje problema, izricanje hipoteza

$H_0$ : daska je simetrična, tj.  $X \sim \text{Bin}(16, 1/2)$

$H_1$ : daska je nagnuta udesno

Raspodjela dana nul-hipotezom:

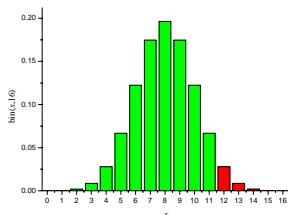


### 2. Postavljanje problema, kritično područje

Prije uzimanja uzorka podijelimo skup mogućih ishoda u dva područja:

A = područje prihvatanja  $H_0$

B = područje odbacivanja  $H_0$  = **kritično područje**



Odlučimo se: A = {0, 1, 2, ..., 11}      B = {12, 13, ..., 16}

### Vrste pogrešaka

		odluka	
		odbaci $H_0$	prihvati $H_0$
istina	$H_0$	pogreška I. vrste	ispravan zaključak
	$H_1$	ispravan zaključak	pogreška II. vrste

Smatra se da je pogreška I. vrste mnogo ozbiljnija od pogreške II. vrste.

Vjerojatnost pogreške I. vrste:

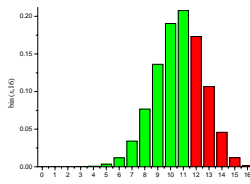
$\alpha$  = vjerojatnost odbacivanja  $H_0$  kada je istinita =  $P(B/H_0)$

Vjerojatnost pogreške II. vrste:

$\beta$  = vjerojatnost prihvatanja  $H_0$  kada nije istinita =  $P(A/H_1)$

U našem primjeru je  $\alpha = 0,038 \approx 4\%$ .

Međutim, čak i da je daska toliko nagnuta da je  $p = 2/3$ , raspodjela varijable  $X$  bi bila



pa bi vjerojatnost pogreške II. vrste uz naš odabir kritičnog područja bila  $\beta = 0,6488$ .

Veliki  $\beta$  se događa kada je uzorak premalen ( $n = 16$ ).

### Signifikantnost testa (važnost, značajnost)

Def:

Za postupak provjere kažemo da ima **razinu signifikantnosti** (važnosti)  $\alpha$  ako je

$$P(\text{pogreška I. vrste}) \leq \alpha.$$

Kažemo da je to **test razine  $\alpha$** .

Naš primjer je test razine 0,04.

Tradicionalno se kao razine signifikantnosti uzimaju vrijednosti 0,01; 0,05 ili 0,10.

Za  $\alpha = 0,05$  kažemo da je test **signifikantan**, a za  $\alpha = 0,01$  kažemo da je test **vrlo signifikantan**.

### Moć testa

Moć statističkog testa mjeri sposobnost testa da odbaci nul-hipotezu kad je uistinu pogrešna, tj. da učini ispravnu odluku.

$$P = 1 - \beta$$

U našem primjeru je  $P = 0,34$ . Idealno je  $P = 1$ . Za snažniji test morali bismo imati veći uzorak.

### Opća pravila odabira testa

Prije uzimanja uzorka:

- 1) Izreci nul-hipotezu i alternativnu hipotezu.
- 2) Razmotri odgovarajuću raspodjelu danu nul-hipotezom.
- 3) Odluči o razini signifikantnosti testa.
- 4) Odredi kritično područje (odluči o kriteriju odbacivanja nul-hipoteze)

Sada uzmi uzorak!

- 5) Očitaj rezultat testa (izračunaj vrijednost statistike testa)
- 6) Učini odluku:

Ako je vrijednost statistike testa u kritičnom području, odbaci  $H_0$ !

Ako vrijednost statistike testa nije u kritičnom području, nemoj odbaciti  $H_0$ !

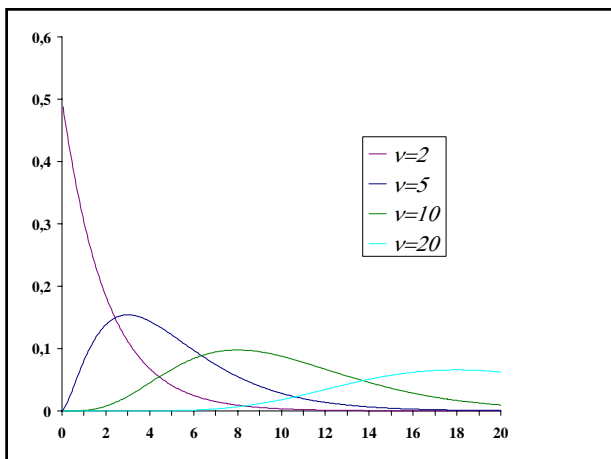
## Testiranje dobrote prilagodbe

### $\chi^2$ raspodjela

Neka je  $\nu \in \mathbb{N}$ . Kontinuirana slučajna varijabla  $X$  ima  $\chi^2$  raspodjelu s parametrom  $\nu$  ako je njezina funkcija gustoće vjerojatnosti dana s

$$f(x; \nu) = \begin{cases} \frac{x^{\nu/2-1} e^{-x/2}}{2^{\nu/2} \Gamma(\nu/2)} & , x \geq 0 \\ 0 & , x < 0 \end{cases}$$

Parametar  $\nu$  zove se "broj stupnjeva slobode" varijable  $X$ .

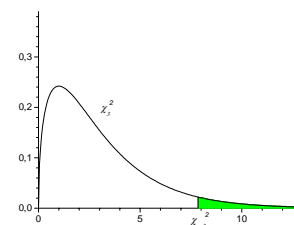


Vjerojatnost da varijabla  $X$  poprimi vrijednost veću od neke određene vrijednosti

$$\chi_{\alpha, \nu}^2$$

dana je površinom ispod repa krivulje:

$$P(X > \chi_{\alpha, \nu}^2) = \int_{\chi_{\alpha, \nu}^2}^{\infty} \chi_{\nu}^2(x) dx$$



Indeks  $\alpha$  označava iznos te vjerojatnosti. Ako  $X$  poprimi vrijednost veću od  $\chi_{\alpha, \nu}^2$ , kažemo da je u **kritičnom području**.

Obično unaprijed odredimo koliku vjerojatnost  $\alpha$  želimo pa na osnovu toga određujemo kritičnu vrijednost  $\chi^2_{\alpha, \nu}$ .

Te vjerojatnosti je teško izračunati pa se služimo tablicama.

U testiranju hipoteza najčešće nas zanimaju vjerojatnosti:

$$P(X > \chi^2_{0.05, \nu}) = 0.05 \quad P(X > \chi^2_{0.01, \nu}) = 0.01$$

Npr., za  $\nu = 3$ , vrijednosti su:

$$\chi^2_{0.05, \nu} = 7.82$$

$$\chi^2_{0.01, \nu} = 11.34$$

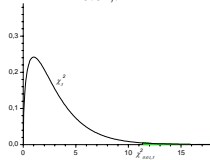
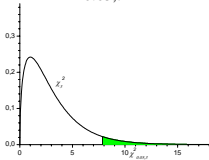
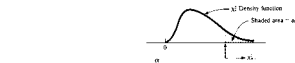


Table A.6 Critical Values  $\chi^2_{\alpha, \nu}$  for the Chi-Squared Distribution



$\alpha$	99	95	90	85	80	75	70	65	60
1	0.004	0.004	0.004	0.004	0.004	0.004	0.004	0.004	0.004
2	0.010	0.010	0.010	0.010	0.010	0.010	0.010	0.010	0.010
3	0.016	0.016	0.016	0.016	0.016	0.016	0.016	0.016	0.016
4	0.020	0.020	0.020	0.020	0.020	0.020	0.020	0.020	0.020
5	0.024	0.024	0.024	0.024	0.024	0.024	0.024	0.024	0.024
6	0.027	0.027	0.027	0.027	0.027	0.027	0.027	0.027	0.027
7	0.289	0.289	0.289	0.289	0.289	0.289	0.289	0.289	0.289
8	1.344	1.344	1.344	1.344	1.344	1.344	1.344	1.344	1.344
9	1.735	1.735	1.735	1.735	1.735	1.735	1.735	1.735	1.735
10	2.158	2.158	2.158	2.158	2.158	2.158	2.158	2.158	2.158
11	2.603	2.603	2.603	2.603	2.603	2.603	2.603	2.603	2.603
12	3.078	3.078	3.078	3.078	3.078	3.078	3.078	3.078	3.078
13	3.581	3.581	3.581	3.581	3.581	3.581	3.581	3.581	3.581
14	4.075	4.075	4.075	4.075	4.075	4.075	4.075	4.075	4.075
15	4.600	4.600	4.600	4.600	4.600	4.600	4.600	4.600	4.600
16	5.142	5.142	5.142	5.142	5.142	5.142	5.142	5.142	5.142
17	5.705	5.705	5.705	5.705	5.705	5.705	5.705	5.705	5.705
18	6.268	6.268	6.268	6.268	6.268	6.268	6.268	6.268	6.268
19	6.841	6.841	6.841	6.841	6.841	6.841	6.841	6.841	6.841
20	7.434	7.434	7.434	7.434	7.434	7.434	7.434	7.434	7.434
21	8.034	8.034	8.034	8.034	8.034	8.034	8.034	8.034	8.034
22	8.641	8.641	8.641	8.641	8.641	8.641	8.641	8.641	8.641
23	9.264	9.264	9.264	9.264	9.264	9.264	9.264	9.264	9.264
24	9.891	9.891	9.891	9.891	9.891	9.891	9.891	9.891	9.891
25	10.521	10.521	10.521	10.521	10.521	10.521	10.521	10.521	10.521
26	11.164	11.164	11.164	11.164	11.164	11.164	11.164	11.164	11.164
27	11.821	11.821	11.821	11.821	11.821	11.821	11.821	11.821	11.821
28	12.461	12.461	12.461	12.461	12.461	12.461	12.461	12.461	12.461
29	13.121	13.121	13.121	13.121	13.121	13.121	13.121	13.121	13.121
30	13.787	13.787	13.787	13.787	13.787	13.787	13.787	13.787	13.787
31	14.451	14.451	14.451	14.451	14.451	14.451	14.451	14.451	14.451
32	15.114	15.114	15.114	15.114	15.114	15.114	15.114	15.114	15.114
33	15.775	15.775	15.775	15.775	15.775	15.775	15.775	15.775	15.775
34	16.434	16.434	16.434	16.434	16.434	16.434	16.434	16.434	16.434
35	17.091	17.091	17.091	17.091	17.091	17.091	17.091	17.091	17.091
36	17.745	17.745	17.745	17.745	17.745	17.745	17.745	17.745	17.745
37	18.397	18.397	18.397	18.397	18.397	18.397	18.397	18.397	18.397
38	19.046	19.046	19.046	19.046	19.046	19.046	19.046	19.046	19.046
39	19.692	19.692	19.692	19.692	19.692	19.692	19.692	19.692	19.692
40	20.337	20.337	20.337	20.337	20.337	20.337	20.337	20.337	20.337

For  $\nu > 40$ ,  $\chi^2_{\alpha, \nu} \approx \nu \left( 1 - \frac{2}{3\nu} + z_{\alpha} \sqrt{\frac{2}{3\nu}} \right)$

### $\chi^2$ -test

Razmotrimo pokus koji daje  $n$  opaženih frekvencija  $f_i$ .  
 Želimo li postaviti hipotezu da ta opažanja slijede neku teorijsku raspodjelu izračunamo očekivane teorijske vrijednosti  $f_{ii}$ .  
 Opažene frekvencije odstupaju od teorijskih, a mi želimo donijeti odluku mogu li se odstupanja pripisati slučaju ili ne.  
 Nul-hipoteza je: "Opažanja slijede teorijsku raspodjelu".

$$\text{Definiramo statistiku: } \chi^2_{op} = \sum_{i=1}^n \frac{(f_i - f_{ii})^2}{f_{ii}}$$

**Teorem (K. Pearson):**

Statistika  $\chi^2_{op}$  približno je raspodijeljena prema  $\chi^2$  raspodjeli sa stupnjem slobode  $\nu$  koji ovisi o broju opažanja (ili razreda)  $i$  o broju postavljenih ograničenja.

Napomene:

- Uzimanjem  $\chi^2$  raspodjele, diskretnu raspodjelu aproksimiramo kontinuiranom. Ta aproksimacija ne vrijedi ako je očekivana frekvencija manja od 5. Taj problem prevladavamo grupiranjem nekoliko razreda male frekvencije u jedan razred dovoljno velike frekvencije.
- Ako je vrijednost  $\chi^2_{op}$  vrlo mala, mudro je upitati se nisu li opaženi rezultati namješteni.

### Stupnjevi slobode

Parametar  $\nu$  definiran je kao:  
 $\nu = \text{broj nezavisnih varijabli uključenih u izračun}$

Nalazimo ga na sljedeći način:  
 $\nu = \text{broj razreda} - \text{broj ograničenja}$

### Normalna raspodjela, $\mu$ i $\sigma$ poznati

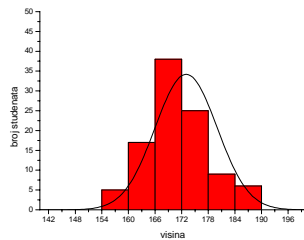
**Primjer:**  
 Dugogodišnje statistike pokazuju da je visina studenata na nekom sveučilištu normalno raspodijeljena s očekivanjem  $\mu = 173$  cm i standardnom devijacijom  $\sigma = 7$  cm.  
 Iz jedne generacije studenata na tom sveučilištu izdvojeno je 100 studenata, mjerene su njihove visine i svrstane u razrede.  
 Opažene su frekvencije:

$x$	154-160	160-166	166-172	172-178	178-184	184-190	Total
$f(x)$	5	17	38	25	9	6	100

Provjeri s 5% signifikantnosti je li raspodjela normalna s s očekivanjem  $\mu = 173$  cm i standardnom devijacijom  $\sigma = 7$  cm !

Rješenje:

$H_0$ : Raspodjela je Gaussova s očekivanjem  $\mu = 173$  cm i standardnom devijacijom  $\sigma = 7$  cm.



Imamo jedno ograničenje:  $\sum f_i = \sum f_{ii} = N = 100$

Uvodimo standardnu normalnu slučajnu varijablu  $z_i = \frac{x - \mu}{\sigma}$ .

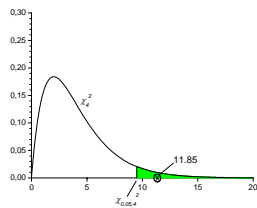
Za te vrijednosti  $z_i$  određujemo iz tablica funkciju raspodjele  $F(z)$  i vjerojatnosti razreda  $f(z)$ :

razred $i$	$f_i$	$z_i$	$F(z)$	$f(z)$	$f_{ii}$	$(f_i - f_{ii})^2 / f_{ii}$
154-160	5	$z < -1,86$	0,0314	0,0314	3	
160-166	17	$-1,86 < z < -1$	0,1587	0,1273	13	2,25
166-172	38	$-1 < z < -0,14$	0,4443	0,2856	28	3,57
172-178	25	$-0,14 < z < 0,71$	0,7611	0,3168	32	1,53
178-184	9	$0,71 < z < 1,57$	0,9418	0,1807	18	4,5
184-190	6	$z > 1,57$	1	0,0582	6	0
zbroj	100			1	100	11,85

Dakle, imamo pet razreda i jedno ograničenje pa je  $\nu = 5 - 1 = 4$

Za 4 stupnja slobode i 5% signifikantnosti, kritična vrijednost (iz tablica) je

$$\chi_{0,05,4}^2 = 9,49$$



Izračunata vrijednost  $\chi_{op}^2 = 11,85$  pada u kritično područje.

Stoga hipotezu  $H_0$  odbacujemo.

### Normalna raspodjela, $\mu$ i $\sigma$ nepoznati

Isti primjer:

Odbacili smo hipotezu o očekivanju i standardnoj devijaciji, ali i dalje vjerujemo da je raspodjela normalna.

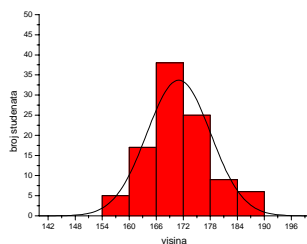
Provjeri s 5% signifikantnosti je li raspodjela normalna!

Rješenje:

$H_0$ : Raspodjela je Gaussova s očekivanjem  $\mu = \bar{x}_{uzorka}$  i standardnom devijacijom  $\sigma = \sigma_{uzorka}$ .

Izračunamo:

$$\bar{x}_{uzorka} = \frac{1}{N} \sum f_i x_i = 171 \text{ cm} \quad \sigma_{uzorka} = \sqrt{\frac{1}{N} \sum f_i (x_i - \bar{x}_{uzorka})^2} = 7,1 \text{ cm}$$



razred $i$	$f_i$	$z_i$	$F(z)$	$f(z)$	$f_{ii}$	$(f_i - f_{ii})^2 / f_{ii}$
154-160	5	$z < -1,55$	0,0606	0,060	6	0,167
160-166	17	$-1,55 < z < -0,70$	0,2420	0,181	18	0,056
166-172	38	$-0,70 < z < 0,14$	0,5557	0,314	32	1,125
172-178	25	$0,14 < z < 0,99$	0,8389	0,283	28	0,321
178-184	9	$0,99 < z < 1,83$	0,9664	0,128	13	0,063
184-190	6	$z > 1,83$	1	0,033	3	
zbroj	100			1	100	1,732

Dakle, imamo pet razreda i tri ograničenja:  $N = 100$

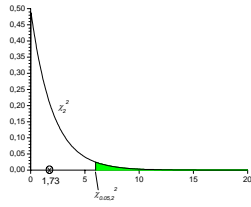
$$\mu = \bar{x}_{uzorka}$$

$$\sigma = \sigma_{uzorka}$$

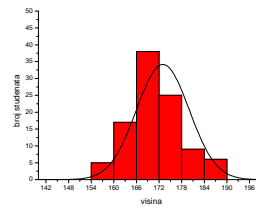
pa je  $\nu = 5 - 3 = 2$

Za 2 stupnja slobode i 5% signifikantnosti, kritična vrijednost je

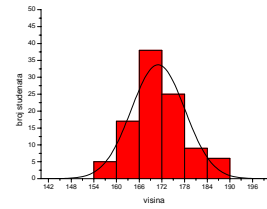
$$\chi^2_{0.05,2} = 5,99$$



Izračunata vrijednost  $\chi^2_{op} = 1,73$  ne pada u kritično područje.  
Stoga hipotezu  $H_0$  **zadržavamo**.



**ODBACUJEMO!**



**PRIHVAĆAMO!**