

Poglavje 9

Osnove matematične statistike

9.0 Uvod

V dosedanjih poglavjih smo izdelali način opisa naključnih dogodkov, spremenljivk in funkcij naključnih spremenljivk s pomočjo verjetnosti.

Verjetnost za nastop naključnega dogodka smo opisali s porazdelitveno funkcijo verjetnosti.

Porazdelitveno funkcijo lahko določimo:

- 1) analitično na osnovi predpostavk o lastnostih dogodkov
- 2) empirično na podlagi poskusov

Problem na katerega naletimo pri opisu naključnega pojava na osnovi poskusov je, omejenost števila poskusov.

Osnovna naloga statistike je sklepanje o statističnih lastnosti (o porazdelitveni funkciji, momentih) naključne spremenljivke X na osnovi omejenega števila poskusov.

9.1 Osnovni pojmi

Množico vseh možnih izidov poskusa, ki ustreza celotnemu vzorčnemu prostoru S imenujemo populacija.

⇒

Populacijo predstavimo z naključno spremenljivko X .

V večini praktičnih primerov je nemogoče izvesti meritev na celotni populaciji X .

Zadovoljimo se z izbrano podmnožico n vrednosti meritev x_i iz populacije X oziroma vzorcem v .

$$v = (x_1, x_2, \dots, x_n)$$

Pri tem n imenujemo velikost

ali razsežnost vzorca v .

Primer: kvaliteta izdelka v serijski proizvodnji.

1- dober, 0-slab:

$$v = (x_1, x_2, \dots, x_n) = (1, 1, 0, \dots, 1)$$

Ob ponoviti izbire vzorca v v splošnem dobimo različne realizacije:

$$v_1 = (x_{11}, x_{12}, \dots, x_{1n})$$

$$v_k = (x_{k1}, x_{k2}, \dots, x_{kn})$$

.....

$$v_m = (x_{m1}, x_{m2}, \dots, x_{mn})$$

Zato vzorec v v splošnem predstavlja naključno spremenljivko:

$$V = (X_1, X_2, \dots, X_n)$$

V primeru, ko so posamezne naključne spremenljivke X_i , ki nastopajo v vzorcu \mathbf{V} medsebojno **statistično neodvisne** in imajo **isto porazdelitev** $f_X(x)$, predstavlja vzorec

$$\mathbf{V} = (X_1, X_2, \dots, X_n)$$

naključni vzorec.

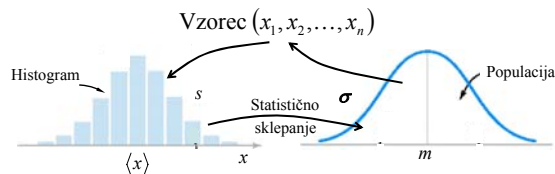
⇒

V praktičnih primerih moramo **zagotoviti naključnost vzorca!!**

Osnovni naloga statistike je na osnovi **izbranega naključnega vzorca**

$$\mathbf{V} = (X_1, X_2, \dots, X_n)$$

sklepati na statistične lastnosti obravnavane populacije X .



V ta namen vpeljemo **vzorčne karakteristike**, ki jih imenujemo **statistike**.

Primeri statistik:

Vzorčno povprečje:

$$\langle X \rangle_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Vzorčna varianca:

$$s_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \langle X \rangle_n)^2$$

Popravljen vzorčna varianca:

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \langle X \rangle_n)^2$$

Razpon vzorca:

$$\Delta X_n = X_{\max, n} - X_{\min, n}$$

Vzorčna relativna frekvenca dogodka A :

$$p_n(A) = \frac{n(X_i \in A)}{n}$$

Vzorčna porazdelitvena funkcija:

$$F_n(x) = \frac{n(X_i \leq x)}{n}$$

V splošnem predstavlja statistiko \mathbf{Z} poljubna merljiva funkcija \mathbf{Z} opredeljena na naključnem vzorcu \mathbf{V} .

$$\mathbf{Z} = Z(\mathbf{V}) = Z(X_1, X_2, \dots, X_n)$$

⇒

Statistika \mathbf{Z} je naključna spremenljivka.

Kar pomeni, da ima \mathbf{Z} zalogo vrednosti in neko porazdelitev verjetnosti $F_Z(z)$.

Statistike, ki se uporabljajo za oceno parametrov θ populacije X na osnovi naključnega vzorca V imenujemo tudi **cenilke**.

Primeri parametrov, ki jih najpogosteje ocenjujemo so $E(X)$, $\text{Var}(X)$,

V splošnem cenilko parametra θ populacije X označimo z $\hat{\theta}$ vrednost cenilke parametra populacije X pa z $\hat{\theta}$.

Da neko statistiko $Z = \hat{\theta}$ uporabimo kot cenilko za parametra θ morajo vrednosti statistike $\hat{\theta} = \hat{\theta}(V)$ imeti lastnost:

$$\lim_{n \rightarrow \infty} P\left[|\hat{\theta} - \theta| < C\right] = 1$$

kjer je C izbrana pozitivna konstanta.

Pri tem sta vrednost cenilke $\hat{\theta}$ naključna in parameter θ deterministična količina.

Splošne lastnosti cenilk

1) Doslednost cenilke

Statistika $\hat{\theta}$ je dosledna cenilka parametra θ , če za poljubno majhen ε velja:

$$\lim_{n \rightarrow \infty} P\left[|\hat{\theta} - \theta| < \varepsilon\right] = 1$$

Kar pomeni, da vrednost $\hat{\theta}$ cenilke $\hat{\theta}$ v smislu verjetnosti konvergira k parametru θ .

2) Nepristranost cenilke

Statistika $\hat{\theta}$ je nepristranska ali centrirana cenilka parametra θ , če je njeno statistično povprečje:

$$E[\hat{\theta}] = \theta$$

3) Asimptotska nepristranost cenilke

Statistika $\hat{\theta}$ je asimptotsko nepristranska cenilka parametra θ če velja:

$$\lim_{n \rightarrow \infty} E[\hat{\theta}] = \lim_{n \rightarrow \infty} (\hat{\theta} + O(n)) = \theta$$

9.2 Lastnosti pomembnejših statistik

9.2.1 Vzorcna relativna frekvenca

Z A označimo dogodek, da smo pri merjenju naključen spremenljivke X izmerili vrednost, ki leži v intervalu S_A .

$$A = (x_i \in S_A)$$

Izide n ponovitev poskusa zabeležimo z vzorcem

$$V = (X_1, X_2, \dots, X_n)$$

z n_A označimo število izidov dogodka A .

Vzorčno relativno frekvenca dogodka A opredelimo s statistiko:

$$p_n(A) = \frac{n(X_i \in A)}{n}$$

Nastop dogodka A v vzorcu V obravnavamo kot izid Bernoullijevega poskusa pri katerem se dogodek A zgodi z verjetnostjo p ali pa se z verjetnostjo $1-p$ ne zgodi.

⇒

Relativno frekvenco $p_n(A)$ lahko obravnavamo kot skupino n Bernoullijevih poskusov oziroma kot binomsko spremenljivko.

Število nastopov n_A dogodka A ima Bernoullijevo ali binomsko porazdelitev verjetnosti:

$$P(n_A; n) = \binom{n}{n_A} p^{n_A} (1-p)^{n-n_A}$$

z $E[n_A] = np$ in $\text{Var}(n_A) = np(1-p)$.

⇒

$$E[p_n(A)] = E\left[\frac{n_A}{n}\right] = \frac{1}{n} E[n_A] = p$$

⇒

Relativna frekvenca dogodka A je nepristranska cenilka verjetnosti dogodka A .

Statistični raztros relativne frekvenca je podan z (glej primer $Y = aX + b$):

$$\text{Var}(p_n(A)) = \text{Var}\left(\frac{n_A}{n}\right) = \frac{1}{n^2} \text{Var}(n_A) = \frac{p(1-p)}{n}$$

in z $n \rightarrow \infty$ pada proti 0.

⇒

Kakovost ocene verjetnosti p_A nastopa dogodka A s pomočjo relativne frekvenca $p_n(A)$ narašča z $n \rightarrow \infty$.

Verjetnost, da se ocena verjetnosti nastopa dogodka A (s pomočjo statistike relativne frekvenca) loči od verjetnosti nastopa dogodka A za več kot ε podamo z neenačbo Čebiševa:

$$P(|X - E[X]| \geq \varepsilon) \leq \frac{\text{Var}(X)}{\varepsilon^2}$$

$$P(|p_n(A) - p| \geq \varepsilon) \leq \frac{\text{Var}(p_n(A))}{\varepsilon^2} = \frac{p(1-p)}{n\varepsilon^2}$$

⇒

$$\lim_{n \rightarrow \infty} P[|p_n(A) - p| \geq \varepsilon] = 0$$

Relativna frekvenca je dosledna cenilka verjetnosti.

9.2.2 Vzorčna porazdelitvena funkcija

Na razpolago imamo vzorec meritev:

$$\mathbf{V} = (X_1, X_2, \dots, X_n)$$

Z A označimo meritev pri kateri dobimo rezultat $X_i \leq x$.

Z $n_A(x) = n(X_i < x)$ označimo število izmerjenih vrednosti, ki zadoščajo pogoju za nastop dogodka A na izbranem vzorcu \mathbf{V} :

Vzorčno porazdelitveno funkcijo definiramo z:

$$F_n(X) = \frac{n(X_i \leq x)}{n}$$

Vzorčna porazdelitveno funkcijo uporabimo jo kot cenilko zbirne porazdelitve verjetnosti:

$$F(X) = P(X \leq x)$$

Iz definicije $F_n(X)$ in ugotovljenih lastnosti relativne frekvenca lahko pokažemo, da je:

$$F_n(X) = \frac{n(X_i \leq x)}{n}$$

nepristranska in dosledna cenilka $F(X)$.

V praktičnih primerih nas pogosto zanima vrednost naključne spremenljivke x_p pri kateri velja:

$$F_n(X) = \frac{n(X_i \leq x_p)}{n} = p$$

Vrednost x_p imenujemo **fraktil**.

Najpogosteje uporabljamo naslednje fraktile:

mediana: $x_{1/2}$

kvartili: $x_{1/4}, x_{1/2}, x_{3/4}$

decili: $x_{1/10}, \dots, x_{9/10}$

centili: $x_{1/100}, \dots, x_{99/100}$

Vzorčna porazdelitev:

zalogo vrednosti spremenljivke X razdelimo na k intervalov.

Posamezni interval opišemo s *centralno vrednostjo* x_j in širino intervala Δx_j :

$$x_j - \Delta x_j / 2 \leq X \leq x_j + \Delta x_j / 2 \quad j = 1, \dots, k$$

Z n_j označimo število vzorčnih vrednosti, ki ležijo v j -tem intervalu.

Z $p_{nj} = n_j / n$ označimo *intervalne relativne frekvence*.

Z množico *intervalnih relativnih frekvenc*:

$$\{p_{nj}; j=1, \dots, k\}$$

je podana *vzorčna porazdelitev*.

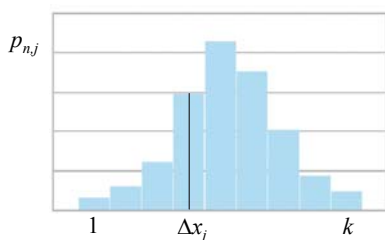
Intervalne relativne frekvence zadoščajo pogoju:

$$\sum_{j=1}^k p_{nj} = 1$$

\Rightarrow Intervalne relativne frekvence so *nepristranske in dosledne cenilke verjetnosti*, ki pripada izbranemu intervalu:

$$p_j = \int_{x_j - \Delta x_j}^{x_j + \Delta x_j} dP(x)$$

Intervalne frekvence $\{p_{nj}; j=1, \dots, k\}$ grafično predstavimo s histogramom:



Razdelitev na intervale je poljubna vendar po priporočilih je število intervalov k med 10 in 100.

9.2.3 Vzorčno povprečje

Iz populacije X izberemo vzorec razsežnosti n :

$$V = (X_1, X_2, \dots, X_n)$$

in opredelimo vzorčno povprečje z izrazom:

$$\langle X \rangle_n = \frac{1}{n} \sum_{i=1}^n X_i$$

$\langle X \rangle_n$ se v splošnem spreminja od vzorca do vzorca

$\Rightarrow \langle X \rangle_n$ je *naključna spremenljivka*

Opredelimo lahko statistično povprečje vzorčnega povprečja:

$$E[\langle X \rangle_n] = E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n E[X_i]$$

Komponente vzorca so med seboj neodvisne in imajo enako porazdelitev verjetnosti $E[X_i] = m$.

\Rightarrow

$$E[\langle X \rangle_n] = \frac{1}{n} \sum_{i=1}^n m = m$$

\Rightarrow *Vzorčno povprečje je nepristrana cenilka statističnega povprečja oz. povprečne vrednosti populacije X*

Raztros vzorčnega povprečja je podan z:

$$\text{Var}(\langle X \rangle_n) = E[(\langle X \rangle_n - E[\langle X \rangle_n])^2] = E[(\langle X \rangle_n - m)^2]$$

Komponente vzorca so med seboj neodvisne in imajo enako porazdelitev verjetnosti: (glej primer $Y = aX + b$)

\Rightarrow

$$\begin{aligned} \text{Var}[\langle X \rangle_n] &= \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X) = \frac{\text{Var}(X)}{n} = \frac{\sigma_X^2}{n} \end{aligned}$$

Iz neenačbe Čebiševa

$$P(|X - E[X]| \geq \varepsilon) \leq \frac{\text{Var}(X)}{\varepsilon^2}$$

$$P(|\langle X \rangle_n - m| \geq \varepsilon) \leq \frac{\sigma_X^2}{n\varepsilon^2}$$

sledi:

$$\lim_{n \rightarrow \infty} P(|\langle X \rangle_n - m| \geq \varepsilon) \leq \frac{\sigma_X^2}{n\varepsilon^2} = 0$$

$$\lim_{n \rightarrow \infty} P(\langle X \rangle_n = m) = 1$$

⇒

Vzorčno povprečje je **dosledna cenilka** statističnega povprečja m populacije X .

9.2.4 Vzorčni momenti

Iz populacije X izberemo vzorec razsežnosti n :

$$V = (X_1, X_2, \dots, X_n)$$

in opredelimo vzorčno varianco:

$$s^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \langle X \rangle_n)^2$$

in popravljeno vzorčno varianco:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \langle X \rangle_n)^2$$

Statistični povprečji opredeljenih varianc sta podani z:

$$E[s^2] = \frac{n-1}{n} \sigma^2 = \sigma^2 + O\left(\frac{1}{n}\right)$$

$$E[S^2] = \sigma^2$$

⇒

s^2 je **asimptotsko nepristranska** med tem ko je S^2 **nepristranska cenilka** variance populacije X .

Varianca nepristranske cenilke S^2 je:

$$\text{Var}(S^2) = \frac{1}{n} \left(\mu^4 - \frac{n-3}{n-1} \sigma^4 \right)$$

$$\lim_{n \rightarrow \infty} \text{Var}(S^2) = 0$$

⇒

S^2 je **dosledna cenilka** variance populacije X .

V splošnem definiramo cenilke vzorčne momentov z izrazom:

$$m_{k,n} = \langle X^k \rangle_n = \frac{1}{n} \sum_{i=1}^n X_i^k, \quad k = 3, 4, \dots$$

Cenilke centralnih momentov pa z:

$$\mu_{k,n} = \langle (X - \langle X \rangle_n)^k \rangle_n = \frac{1}{n} \sum_{i=1}^n (X_i - \langle X \rangle_n)^k$$

Z upoštevanjem neodvisnosti in enakih porazdelitev spremenljivk x_i dobimo:

$$E[m_{k,n}] = m_k$$

$$\text{Var}(m_{k,n}) = \frac{m_{2r} - m_r^2}{n}$$

$$E[\mu_{k,n}] = \mu_k + O(1/n)$$

⇒

$m_{k,n}$ so **nepristranske in dosledne cenilke** momentov m_k
 $\mu_{k,n}$ so **asimptotično nepristranske in dosledne cenilke** centralnih momentov μ_k .

9.3 Verjetnostne porazdelitve statistik

Statistike so funkcije naključnega vzorca \mathbf{V} :

$$Z = Z(\mathbf{V}) = Z(X_1, X_2, \dots, X_n)$$

⇒

Statistika je **naključna spremenljivka**, ki ima ustrezno **zalogo vrednosti** in porazdelitev verjetnosti oziroma **gostoto verjetnosti**:

$$f_Z(z)$$

Določitev $f_Z(z)$ ni splošno rešljiv problem ker je rešitev odvisna od funkcije Z .

V primeru, ko je poznana $f_X(x)$ in je statistika opredeljena z **funkcijo vsote** (glej primer $Z=X+Y$):

$$Z = Z(X) = \sum_{i=1}^n Z(X_i)$$

statistično neodvisnih komponent vzorca:

$$\mathbf{V} = (X_1, X_2, \dots, X_n)$$

je porazdelitev $f_Z(z)$ podana z $n-1$ **kratno konvolucijo** gostot porazdelitev verjetnosti posameznih komponent.

Pri tem je potrebno predhodno določiti gostoto verjetnosti za funkcijo $Z(X)$, ki je v primeru monotone funkcije podana z:

$$f_Z(z) = f_X(x(z)) \left| \frac{dx(z)}{dz} \right|$$

9.3.1 Porazdelitev statistike vzorčnega povprečja

Vzorčno povprečje naključnega vzorca smo opredelili z:

$$Z_n = \langle X \rangle_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Pri tem so X_i statistično neodvisne in imajo enako gostoto porazdelitve:

$$f_{X_i}(x) = f_{X_j}(x)$$

V primeru, ko je X normalno porazdeljena z:

$$f_{X_i}(x) = f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}, \quad \text{in } m_X = 0, \sigma_X^2 = 1$$

Je za $n = 2$, gostota verjetnosti Z podana s konvolucijo:

$$\begin{aligned} f_{Z_2}(z) &= \int_{-\infty}^{\infty} f_{X_1}(x_1) f_{X_2}(nz - x_1) dx_1 \\ &= \frac{1}{\sqrt{2\pi}\sqrt{2}} e^{-\frac{z^2}{4}} \end{aligned}$$

$f_{Z_2}(z)$ je normalna z $m_{Z_2} = m_X = 0$ in $\sigma_{Z_2} = \sigma_X / \sqrt{2}$

V splošnem lahko pokažemo, če ima **populacija X normalno porazdelitev** z:

$$E[X] = m \quad \text{in} \quad \text{Var}(X) = \sigma_X^2$$

Ima statistika **vzorčnega povprečja**:

$$Z_n = \langle X \rangle_n = \frac{1}{n} \sum_{i=1}^n X_i$$

normalno porazdelitev z:

$$E[\langle X \rangle_n] = \frac{1}{n} \sum_{i=1}^n m = m \quad \text{in} \quad \text{Var}(\langle X \rangle_n) = \frac{1}{n^2} \sum_{i=1}^n \sigma_X^2 = \frac{\sigma_X^2}{n}$$

Nadalje, za vzorčno povprečje poljubne populacije velja **centralni limitni teorem**, ki pravi:

Če vzorčimo iz populacije X z neznano porazdelitvijo verjetnosti bo porazdelitev vzorčnega povprečja $\langle X \rangle_n$ približno **normalna** z srednjo vrednostjo in varianco:

$$E[\langle X \rangle_n] = m \quad \text{in} \quad \text{Var}(\langle X \rangle_n) = \sigma_X^2 / n$$

Oziroma če je naključni vzorec

$$V = (X_1, X_2, \dots, X_n)$$

velikosti n izbran iz poljubne populacije X z:

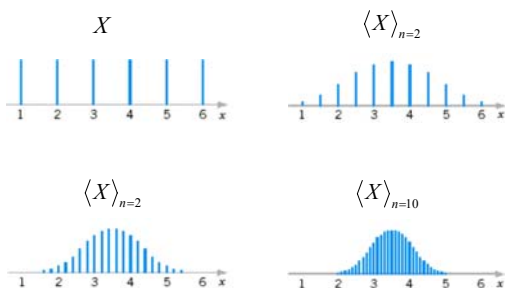
$$E[X] = m \quad \text{in} \quad \text{Var}(X) = \sigma_X^2$$

in če je $\langle X \rangle_n$ **vzorčno povprečje** potem porazdelitev statistike:

$$Z = \frac{\langle X \rangle_n - m}{\sigma / \sqrt{n}}$$

z $n \rightarrow \infty$ limitira k **standardni normalni porazdelitvi**.

Primer: Povprečje meta n kock.



Primer: Populacija X ima enakomerno zvezno porazdelitev:

$$f_X(x) = \begin{cases} 1/2, & x \in [4,6] \\ 0, & \text{izven} \end{cases}$$

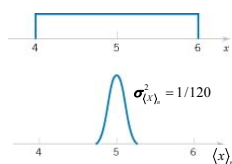
Sprašujemo po porazdelitvi vzorčnega povprečja velikosti vzorca $n=40$:

$$E[\langle X \rangle_n] = m = 5 \quad \text{in} \quad \text{Var}(\langle X \rangle_n) = \frac{(b-a)^2}{12} = 1/3$$

Po centralnem limitnem teoremu je porazdelitev $\langle X \rangle_n$ aproksimirana z normalno porazdelitvijo z:

$$E[\langle X \rangle_n] = E[X] = m = 5$$

$$\text{Var}(\langle X \rangle_n) = \sigma_{\langle X \rangle_n}^2 = \frac{\sigma_X^2}{n} = \frac{1}{3 \cdot 40} = \frac{1}{120}$$



9.3.2 χ^2 Hi-kvadrat porazdelitev

Imamo populacijo X z normalno porazdelitvijo, ki ima:

$$E[X] = m = 0 \quad \text{in} \quad \text{Var}(X) = \sigma_X^2 = 1$$

Zanima nas gostota porazdelitve naključne spremenljivke oziroma statistike Z_n opredeljena z:

$$Z_n = \sum_{i=1}^n X_i^2$$

Ker so X_i in s tem X_i^2 statistično neodvisne je porazdelitev vsote $\sum_{i=1}^n X_i^2$ podana s konvolucijo posameznih gostot porazdelitve $f_{X_i^2}(x)$:

Porazdelitve so enake:

$$f_{X_i^2}(x) = f_{X_j^2}(x) = f_{X_k^2}(x) = \dots = f_{X_n^2}(x) \quad i = 1, 2, \dots, n$$

Ker poznamo gostoto porazdelitve X je porazdelitev funkcije naključne spremenljivke $Z=X^2$ podana z:

$$P(Z < z) = P(-\sqrt{z} \leq X \leq \sqrt{z}) = \int_{-\sqrt{z}}^{\sqrt{z}} f_X(x) dx$$

Porazdelitvena funkcije naključne spremenljivke $Z=X^2$ podana z:

$$P(Z < z) = P(-\sqrt{z} \leq X \leq \sqrt{z}) = \frac{1}{\sqrt{2\pi}} \int_{-\sqrt{z}}^{\sqrt{z}} e^{-x^2/2} dx$$

Gostoto porazdelitve $f_Z(z)$ dobimo z odvajanjem porazdelitvene funkcije $P(Z < z)$ po spremenljivki z :

$$f_Z(z) = \begin{cases} \frac{1}{\sqrt{2\pi}} \frac{e^{-z/2}}{\sqrt{z}}, & \text{za } z > 0 \\ 0, & \text{za } z \leq 0 \end{cases}$$

Z poznavanjem gostote verjetnosti $f_Z(z)$ in uporabe **konvolucije** lahko določimo gostoto verjetnosti vsote dveh kvadratov Z_2

$$\begin{aligned} f_{Z_2}(z) &= \int f_Z(z) f_Z(z-x) dx \\ &= \frac{1}{(2\pi)^{2/2}} \int_0^z \frac{1}{\sqrt{x}} \frac{1}{\sqrt{z-x}} e^{-x/2} e^{-(z-x)/2} dx \\ &= \begin{cases} \frac{1}{2\pi} e^{-z/2} \int_0^z \frac{dx}{\sqrt{x(z-x)}}, & \text{za } z > 0 \\ 0, & \text{za } z \leq 0 \end{cases} \end{aligned}$$

S substitucijo spremenljivke $x=zv$ in $dx=zdv$ v zadnjem integralu dobimo:

$$f_{Z_2}(z) = \begin{cases} \frac{1}{2\pi} e^{-z/2} z^0 \int_0^1 v^{-1/2} (1-v)^{-1/2} dv, & \text{za } z > 0 \\ 0, & \text{za } z \leq 0 \end{cases}$$

Gornji integral lahko izrazimo z beta funkcijo :

$$B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} = \int_0^1 z^{a-1} (1-z)^{b-1} dz$$

Kjer je Γ gama funkcija.

Z uporabo **beta** in **gama** funkcije lahko gostoto verjetnosti $f_{Z_2}(z)$ za $z > 0$ zapišemo v obliki:

$$\begin{aligned} f_{Z_2}(z) &= \frac{1}{(2\pi)^{2/2}} z^0 e^{-z/2} B(1/2, 1/2) \\ &= \frac{1}{(2\pi)^{2/2}} z^0 e^{-z/2} \frac{\Gamma(1/2)\Gamma(1/2)}{\Gamma(1/2+1/2)} \end{aligned}$$

Z upoštevanjem, da je $\Gamma(1/2) = \sqrt{\pi}$ dobimo:

$$f_{Z_2}(z) = \begin{cases} \frac{1}{2^{2/2}\Gamma(2/2)} z^0 e^{-z/2}, & z > 0 \\ 0, & z \leq 0 \end{cases}$$

S pomočjo popolne matematične indukcije lahko pokažemo da velja:

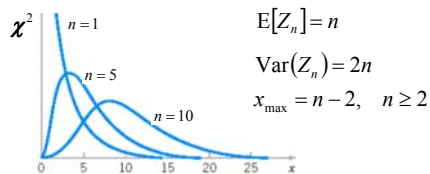
$$f_{Z_n}(z) = \chi^2 = \begin{cases} \frac{1}{2^{n/2}\Gamma(n/2)} z^{(n/2)-1} e^{-z/2}, & z > 0 \\ 0, & z \leq 0 \end{cases}$$

Dobljena gostota verjetnosti za χ^2 se imenuje **"hi-kvadrat"** z n prostostnimi stopnjami.

⇒

Naključna spremenljivka opredeljena z vsoto kvadratov standardiziranih normalnih spremenljivk ima **hi-kvadrat** porazdelitev gostote verjetnosti.

Primeri grafov hi kvadrat porazdelitve v odvisnosti od n :



Vrednosti ostalih značilnih parametrov porazdelitve χ^2 :

$$E[Z_n] = m_r = n(n+2) \cdots (n+2r-2)$$

$$g_1 = 2\sqrt{2/n}$$

$$g_2 = 12/n$$

Aditivna lastnost porazdelitve hi-kvadrat

Naj bodo Y_1, Y_2, \dots, Y_n , hi-kvadrat neodvisne naključne spremenljivke s prostostnimi stopnjami k_1, k_2, \dots, k_n .

$$Y = Y_1 + Y_2 + \dots + Y_n$$

je hi-kvadrat naključna spremenljivka z prostostno stopnjo:

$$k = \sum_{i=1}^n k_i$$

Primeri spremenljivk z χ^2 porazdelitvijo

1. Poljubni normalni porazdeljeni spremenljivki X z:

$$E[X] = m_x \quad \text{in} \quad \text{Var}(X) = \sigma_x^2$$

lahko priredimo standardni odmik od srednje vrednosti Z :

$$Z = \frac{X - m_x}{\sigma_x}$$

Z :

$$E[Z] = m_z = 0 \quad \text{in} \quad \text{Var}(Z) = \sigma_z^2 = 1$$

Vzorčni drugi moment spremenljivke Z opredeljen z:

$$m_{2,n} = \frac{1}{n} \sum_{i=1}^n Z_i^2 = \frac{1}{n} \sum_{i=1}^n \left(\frac{X_i - m_x}{\sigma_x} \right)^2 \propto \sum_{i=1}^n X_i^2$$

ima χ^2 porazdelitev z n prostostnimi stopnjami.

2. Za poljubno normalno spremenljivki X z:

$$E[X] = m_x \quad \text{in} \quad \text{Var}(X) = \sigma_x^2$$

je z $X_i - \langle X \rangle_n$ podan odmik od vzor. povprečja $\langle X \rangle_n$

Na osnovi odmika vpeljemo spremenljivko:

$$\chi_{n-1}^2 = \frac{1}{\sigma_x^2} \sum_{i=1}^n (X_i - \langle X \rangle_n)^2$$

ki ima χ^2 porazdelitve z $n-1$ prostostnimi stopnjami.

3. Z uporabo spremenljivke χ^2 lahko izrazimo vzorčni varianci:

$$s^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \langle X \rangle_n)^2 = \frac{\sigma_x^2}{n} \chi_{n-1}^2$$

in

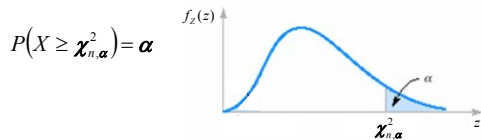
$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \langle X \rangle_n)^2 = \frac{\sigma_x^2}{n-1} \chi_{n-1}^2$$

ki imata χ^2 porazdelitve z $n-1$ prostostnimi stopnjami.

Vrednosti porazdelitvena funkcija χ^2 so podane tabelarično. Tabela podaja verjetnosti:

$$P(X \geq \chi_{n,\alpha}^2) = \int_{\chi_{n,\alpha}^2}^{\infty} f_Z(z) dz = \alpha$$

Kjer je $\chi_{n,\alpha}^2$ označuje vrednost hi-kvadrat spremenljivke X z n prostostnimi stopnjami, pri kateri je:



9.3.3 Studentova porazdelitev t

Vzorčimo iz normalne populacije X z:

$$E[X] = m_X \quad \text{in} \quad \text{Var}(X) = \sigma_X^2$$

Vzorčno povprečje $\langle X \rangle_n$ populacije X ima normalno porazdelitev z:

$$E[\langle X \rangle_n] = m_X \quad \text{in} \quad \text{Var}(\langle X \rangle_n) = \frac{\sigma_X^2}{n}$$

Na osnovi **centralnega limitnega teorema** ima statistika oziroma naključna spremenljivka:

$$Z = \frac{(\langle X \rangle_n - m_X)}{\sigma_X / \sqrt{n}}$$

standardno normalno porazdelitev.

Predpostavimo, da variance σ_X populacije X ne poznamo. Kaj se zgodi z porazdelitvijo spremenljivke Z če v njej σ_X nadomestimo z vzorčno varianco S :

$$T = \frac{(\langle X \rangle_n - m_X)}{S / \sqrt{n}}$$

V splošnem lahko pokažemo, da:

Če je Z **normalna** spremenljivka in V **hi-kvadrat** spremenljivka z n prostostnimi stopnjami in če sta Z in V statistično neodvisni, potem ima spremenljivka:

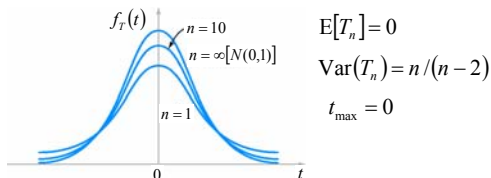
$$T = \frac{Z}{\sqrt{V/n}}$$

gostoto verjetnosti:

$$f_T(t) = \frac{\Gamma[(n+1)/2]}{\sqrt{\pi n} \Gamma(n/2)} \cdot \frac{1}{[(t^2/n)+1]^{(n+1)/2}} \quad -\infty < t < \infty,$$

ki se imenuje **Studentova** ali t_n porazdelitev z n prostostnim stopnjami.

Primeri grafov Studentove gostote verjetnosti:



Vrednosti ostalih značilnih parametrov:

$$E[T_n^{2r}] = m_{2r} = \frac{1 \cdot 3 \cdots (2r-1)}{(n-2)(n-4) \cdots (n-2r)} n^r$$

$$g_1 = 0$$

$$g_2 = 3 \frac{n-2}{n-4} - 3, \quad n > 4$$

V našem primeru imamo spremenljivko:

$$T = \frac{(\langle X \rangle_n - m_X)}{S / \sqrt{n}} = \frac{(\langle X \rangle_n - m_X)}{\sqrt{S^2/n}}$$

Kjer je $\langle X \rangle_n - m_X$ **normalna** in S^2 **hi-kvadrat** spremenljivka z $n-1$ prostostnimi stopnjami, ki sta **statistično neodvisni**.

Zato ima gornja statistika T **Studentovo** oziroma t porazdelitev z $n-1$ prostostnimi stopnjami:

$$f_T(t) = \frac{\Gamma[(n)/2]}{\sqrt{\pi n-1} \Gamma((n-1)/2)} \cdot \frac{1}{[(t^2/(n-1))+1]^{n/2}}$$

Tudi za spremenljivko:

$$T_i = \frac{X_i - m_X}{S}$$

lahko pokažemo, da ima Studentovo porazdelitev z $n-1$ stopnjami.

Studentova porazdelitev je podana v tabeli, ki podaja:

$$P(T_n \geq t_{n,\alpha}) = \int_{t_{n,\alpha}}^{\infty} f_T(t) dz = \alpha$$

