

cv5.R

Marek

2023-03-16

```
library('ltxexpdfr')
library('ltxexp2exp')
```

Jednovyberove testy hypotez pre parametre normalneho rozdelenia.

Cize musi platit predpoklad, ze data su normalne rozdelenie (Shapiro Wilk test) Hypoteza je tvrdenie, ktorého pravdivosť je treba overiť statistickými metódami. Výsledok setrenia je zatazený dvomi chybami (pravdepodobnosťou týchto chýb). Prvá chyba, chyba prvého druhu (pravdepodobnosť chyby) α je hladina významnosti. Zamietam nulovú hypotézu H_0 (tvrdenie) a ona platí. Chyba druhého druhu je β , nezamietam H_0 a H_0 neplatí. Jednu chybu ľoneme a druhu minimalizujeme tak, že správne vykonáme test. K H_0 sa pridáva alternatívna hypotéza H_1 . Ak hypotézu H_0 zamietnem, tak prijmem alternatívu H_1

Test pre strednu hodnotu μ ak σ je známe

$$H_0 \quad \mu = \mu_0 \qquad H_1 \mu \neq \mu_0$$

Zvolíme vhodnú štatistiku, vypočítame ju, rozdelíme R na oblasť, kde zamietam H_0 (to sa volá kritická oblasť testu) a kde nezamietam H_0 (dokopy musia dať obor hodnôt) Priklad 1 Vyrobcia uvádza, že v kopírke je nutné meniť toner v priemere po 2500 skopirovaných stranách so smerodajnou odchýlkou $= 30$. Na hladine významnosti $\alpha = 0.05$ testujeme, či tvrdenie výrobcu je v súlade so skutočnosťou. Najprv vypočítom štatistiky, potom prikazom v R. Každý test sa začína stanovením nulovej a alternatívnej hypotézy.

$$H_0 \quad \mu = 2500 \qquad H_1 \quad \mu \neq 2500$$

```
x <- c(2445, 2450, 2453, 2462, 2463, 2463, 2466, 2471, 2474, 2475, 2475,
      2484, 2485, 2486, 2487, 2490, 2491, 2493, 2499, 2501, 2501, 2503,
      2504, 2505, 2505, 2506, 2506, 2507, 2509, 2511, 2511, 2513, 2514,
      2515, 2519, 2523, 2523, 2524, 2525, 2527, 2529, 2530, 2530, 2533,
      2535, 2536, 2537, 2539, 2560, 2571)

alfa <- 0.05
mean <- mean(x)
sigma <- 30
n <- length(x)
mi0 <- 2500
t <- (mean - mi0) * sqrt(n) / sigma
t

## [1] 0.7683894

k1 <- qnorm(alfa / 2, 0, 1)
k2 <- qnorm(1 - alfa / 2, 0, 1)
c(t, k1, k2)

## [1] 0.7683894 -1.9599640 1.9599640
```

Štatistika padla do intervalu (-1.96, 1.96), nezamietame H_0 , tvrdenie výrobcu je pravdivé alebo v súlade so skutočnosťou, pomocou knižnice DescTools a prikazu v R.

```
library('DescTools')
ZTest(x, mu = 2500, sd_pop = 30) #chyba je defaultne nastavena na 0.05

##
## One Sample z-test
##
## data: x
## z = 0.76839, Std. Dev. Population = 30, p-value = 0.4423
## alternative hypothesis: true mean is not equal to 2500
## 95 percent confidence interval:
## 2494.945 2511.575
## sample estimates:
## mean of x
## 2503.26
```

rozhodnute bud pozriem, či testovaná hodnota je z príslušného intervalu spoľahlivosti, alebo rozhodneme podľa P-value, pravdepodobnostnej hodnoty P-Hodnoty. Ak je p-hodnota menšia ako alfa (0.05), tak zamietam nulovú hypotézu H_0 .

```
ZTest(x, mu = 2500, sd_pop = 30)$p.value

## [1] 0.4422559
```

P-Hodnota je 0.44 > 0.05, teda nezamietam H_0

```
pvalue <- pnorm(-abs(t)) * 2
pvalue

## [1] 0.4422559
```

Toto bolo o výrobcovi. Mňa ako používateľa skor zaujíma overenie

tvrdenia, že nakopírujem aspon 2500 strán

$$H_0 \quad \mu = 2500 (\geq 2500) \qquad H_1 \quad \mu < 2500$$

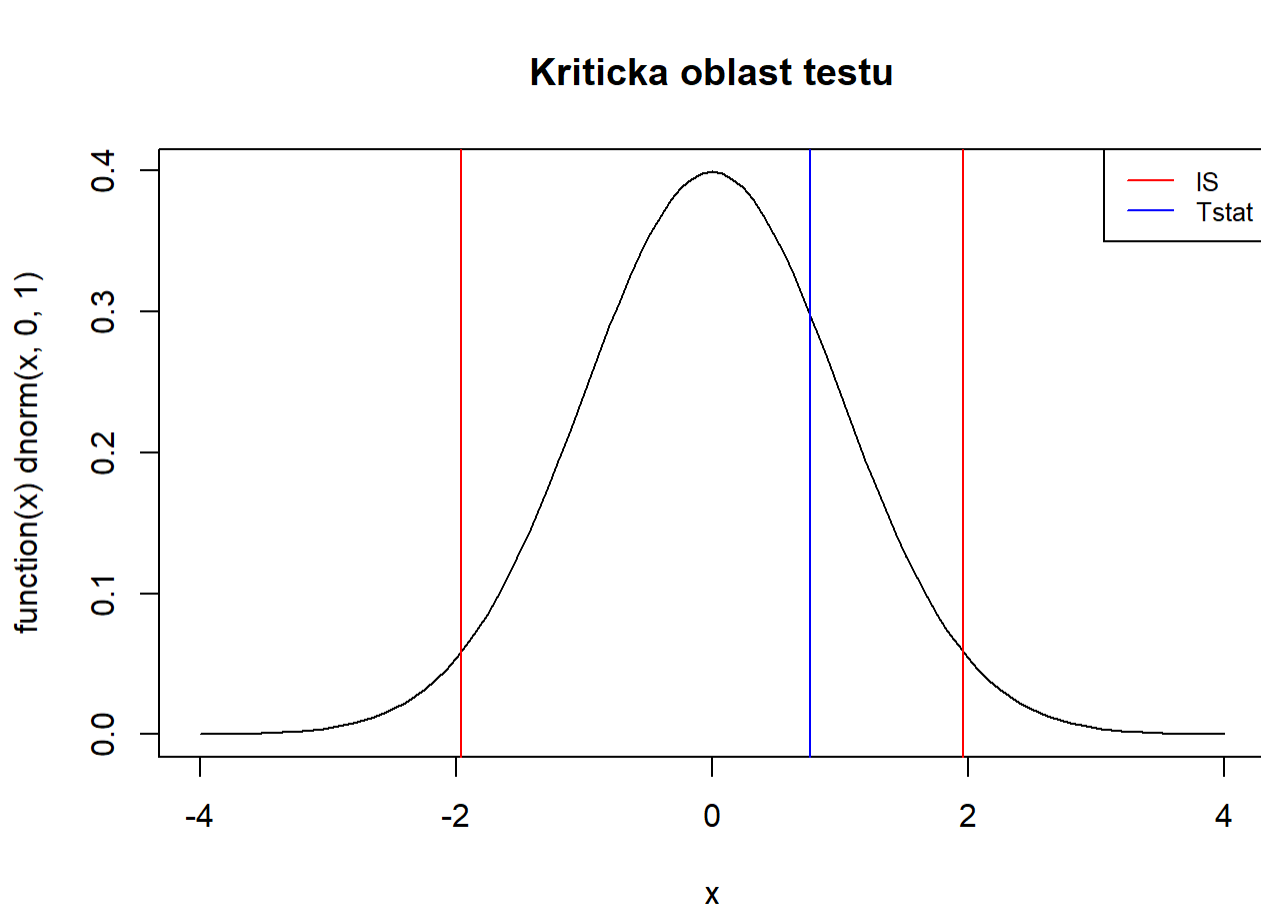
alternatíva je less

```
ZTest(x, mu=2500, sd_pop=30, alternative = 'l')

##
## One Sample z-test
##
## data: x
## z = 0.76839, Std. Dev. Population = 30, p-value = 0.7789
## alternative hypothesis: true mean is less than 2500
## 95 percent confidence interval:
## -Inf 2510.239
## sample estimates:
## mean of x
## 2503.26
```

P-Hodnota je 0.77>0.05, nezamietam H_0 nakopírujem toľko alebo viac, ale určite nie menej Grafický výstup

```
plot(function(x)dnorm(x, 0, 1),
      xlim=c(-4, 4),
      main = 'Kritická oblasť testu')
abline(v = k1, col = 'red')
abline(v = k2, col = 'red')
abline(v = t, col = 'blue')
legend('topright', c('l1', 'Tstat'), col = c('red', 'blue'), lt = 1, cex = 0.8)
```



Test pre strednu hodnotu, ak σ nie je známe

Riešte rovnakú úlohu ako hore, ale za predpokladu, že σ nepoznáme. Najskor vypočtom, potom prikazom.

```
sd <- sd(x)
t <- (mean - mi0) * sqrt(n) / sd
k1 <- qt(alfa / 2, n - 1)
k2 <- qt(1 - alfa / 2, n - 1)
c(t, k1, k2)

## [1] 0.8214202 -2.0095752 2.0095752
```

Hodnota vypočítanej štatistiky padne do intervalu Spoľahlivosti, nezamietam H_0 .

```
#abline(v = k1, col = 'green')
#abline(v = k2, col = 'green')
#abline(v = t, col = 'purple')
```

prikazom v R, rozhodneme podľa P hodnoty

```
t.test(x, mu = 2500)

##
## One Sample t-test
##
## data: x
## t = 0.82142, df = 49, p-value = 0.4154
## alternative hypothesis: true mean is not equal to 2500
## 95 percent confidence interval:
## 2495.285 2511.235
## sample estimates:
## mean of x
## 2503.26

t.test(x, mu = 2500)$p.value

## [1] 0.4153858
```

P-Hodnota = 0.41 > 0.05, nezamietame H_0 . Tvrdenie o počte skopirovaných strán je pravdivé Zmena hladiny významnosti príkladom Malý rodinný podnik predáva 100% bio jablkovú stavu v balení 0.5l. Po oprave pniacej linky sme namerali tieto hodnoty objemu balenia(ml) na hladine významnosti $\alpha = 0.1$ testujeme hypotézu, že plniacia linka je dobre nastavená

$$H_0 \quad \mu = 500 \qquad H_1 \mu \neq 500$$

```
js <- c(495.2,496.8,502.1,498.5,501.503,500.7,
      501.5,501.8,499.1,500.9,502.2,501.7, 500.4,
      500.2,501.1,499.9,500.2,501.1,500.8,499.3)
t.test(js, mu = 500, conf.level = 0.9)

##
## One Sample t-test
##
## data: js
## t = 0.89828, df = 20, p-value = 0.3797
## alternative hypothesis: true mean is not equal to 500
## 90 percent confidence interval:
## 499.6714 501.0429
## sample estimates:
## mean of x
## 500.3571
```

P-Hodnota je 0.37, čo je viac ako 0.1, takže H_0 nezamietame, plniacia linka je dobre nastavená Ďalší príklad Firma ABC, ktorá vyrába batérie do notebookov tvrdí, že jednu vyrobí v priemere do 13 minút (max 13 minút). Na hladine významnosti $\alpha = 0.05$ overte toto tvrdenie, ak máte k dispozícii dataset dĺžok výroby. greater

$$H_0 \quad \mu = 13 (\leq 13) \qquad H_1 \quad \mu > 13$$

```
bat <- c(12,19,16,11,7,12,14,18,15,19,17,19,13,9,
      11,20,12,19,8,13)
t.test(bat, mu = 13, alternative = 'g')

##
## One Sample t-test
##
## data: bat
## t = 1.3346, df = 19, p-value = 0.09889
## alternative hypothesis: true mean is greater than 13
## 95 percent confidence interval:
## 12.6453 Inf
## sample estimates:
## mean of x
## 14.2

t.test(bat, mu = 13, alternative = 'g')$p.value

## [1] 0.09888518
```

P hodnota = 0.098 < 0.05, nezamietam H_0 , tvrdenie firmy je v súlade so skutočnosťou. Ak by sa zmenila hladina významnosti na 0.1, tak 0.098 < 0.1 na tejto hladine významnosti by sme zamietli hypotézu Ďalší príklad Firma XYZ nakupuje batérie do elektronických hračiek. Vyrábajú garantuje, že vydržia minimálne 19 hodín nepretržitej prevádzky (19 a viac). Kontrolor náhodne vyberie 12 batériek a získa údaje o životnosti. Testujeme na hladine významnosti $\alpha = 0.05$.

$$H_0 \quad \mu = 19 (\geq 19) \qquad H_1 \quad \mu < 19$$

alternatíva bude less

```
bat1 <- c(20,2,19,6,18,6,19,4,17,18,5,18,18,4,19,
      18,17,9,18,1)
t.test(bat1, mu = 19, alternative = 'l')$p.value

## [1] 0.05397316
```

P hodnota 0.053 > 0.05, nezamietam H_0 , vydržia 19 a viac, nie menej

Test pre disperziu overte tvrdenie výrobcu (prvý príklad) aj čo sa týka disperzie, najprv vypočtom a potom prikazom v R

```
alfa <- 0.05
t <- (n - 1) * var(x) / 30^2
k1 <- qchisq(alfa / 2, n - 1)
k2 <- qchisq(1 - alfa / 2, n - 1)
c(t, k1, k2)

## [1] 42.87736 31.55492 70.22241
```

Pomocou knižnice EnvStats

```
library('EnvStats')

##
## Attaching package: 'EnvStats'

##
## The following objects are masked from 'package:stats':
##
## predict, predict.lm

##
## The following object is masked from 'package:base':
##
## print.default

VarTest(x, sigma.squared = 30^2)

##
## One Sample Chi-Square test on variance
##
## data: x
## X-squared = 42.877, df = 49, p-value = 0.5632
## alternative hypothesis: true variance is not equal to 900
## 95 percent confidence interval:
## 549.5342 1222.9353
## sample estimates:
## variance of x
## 787.5433
```

Štatistika padla do oblasti, kde nezamietam H_0 , tvrdenie výrobcu o disperzii je pravdivé # dvojitý testy o parametroch normalného rozdelenia parove (závisle merania) neaparove (nezávisle merania) Parove testy sú určené pre dvojice dát, ktoré sú spojené jedným objektom, práh pocutenosti, prava ľave ucho, ojazdenosť pneu predná/zadná data by sme mali dostať ako dvojice, najčastejšie sa používajú v prípadoch pred liečbou, po liečbe, pred a potom, pred skúškou, pred skúškou, po skúške, zaujíma nás zmena (X, Y), Z = X - Y, a sa nie zmenilo, tak vlastne testujeme jednovýberový testom, že stredná hodnota Z je nula. Sú dôležité v sekundách, počas ktorých vytvárali kontrolné úlohy pred a po špeciálnych cvičeniach z pamätového počtania. Zlepšili cvičenia schopnosť znakov rýchlejšie riešiť úlohy? (zlepšili sa, ak to zvládli rýchlejšie), pred - po => 0, alternatíva je menej ako 0

```
pred <- c(87,61,98,90,93,74,83,72,81,75,83)
po <- c(30,45,79,90,88,65,52,79,84,61,52)
z <- pred - po
t.test(z, mu = 0, alternative = 'l')
```

```
##
## One Sample t-test
##
## data: z
## t = 3.1128, df = 10, p-value = 0.0945
## alternative hypothesis: true mean is less than 0
## 95 percent confidence interval:
## -Inf 21.86391
## sample estimates:
## mean of x
## 13.81818
```

P-Hodnota = 0.99 nezamietam H_0 , nezmenili alebo zlepšili, nie zhoršili

```
t.test(pred, po, mu = 0, paired = T, alternative = 'l')

##
## Paired t-test
##
## data: pred and po
## t = 3.1128, df = 10, p-value = 0.0945
## alternative hypothesis: true mean difference is less than 0
## 95 percent confidence interval:
## -Inf 21.86391
## sample estimates:
## mean difference
## 13.81818
```