

Változók típusai

Diszkrét (kategorikus)	<p>Véges vagy legfeljebb megszámlálhatóan végtelen pl.: nem (férfi/nő), vércsoport, iskolai végzettség, ...</p> <p><u>Eloszlása</u> Mik a változó lehetséges értékei, és azokat milyen gyakorisággal veszi fel</p>
Folytonos	<p>Adott intervallumon végtelen sok érték Pl.: testmagasság, csapadékmennyiség, ...</p> <p><u>Eloszlásának jellemzői</u></p> <ol style="list-style-type: none"> Középpont <ul style="list-style-type: none"> mintaátlag (\bar{x}): várható érték becslése medián (Q2): aminél az adatok fele \leq módusz: leggyakoribb mintaelem <p> min---Q1---Q2---Q3---max </p> Szóródás <ul style="list-style-type: none"> terjedelem = max - min IQR = Q3 - Q1 szórás: átlagosan milyen távol vannak az adatok a mintaátlagtól Alak <ul style="list-style-type: none"> szimmetrikus (átlag = medián) jobbra ferde (átlag > medián) $skewness() > 0$ balra ferde (átlag < medián) $skewness() < 0$

T-próba

n = elemszám

\bar{x} = mintaátlag

c = vizsgált érték

SD = minta szórása

SE = standard hiba = SD / \sqrt{n}

alfa = szignifikancia szint

```
elemszam = 80
mintaatlag = 118
mintaszoras = 12
teszt = 110
alpha = 0.05
```

		H0-t elutasítjuk \Leftrightarrow
Próbastatisztika alapján	$t = (\bar{x} - c) / SE$ $t_{\alpha/2}$: n-1 szabadsági fokú t-eloszlás 1-($\alpha/2$) kvartilise	$ t > t_{\alpha/2}$
<pre>egymintas_t = function(mu, n, xvonas, sd) { (xvonas-mu) / (sd/sqrt(n)) } t = egymintas_t(teszt, elemszam, mintaatlag, mintaszoras) t_alpha = qt(1-alpha/2, df = elemszam-1)</pre>		

Konfidencia-intervallum alapján	$KI = \bar{x} \pm t_{\alpha/2} * SE$	$c \notin KI$
<pre># 95% megbízhatóságú konfidencia intervallum a várható értékre konf.int_also = function(n, xvonas, sd, alpha) { xvonas-qt(1-alpha/2, df = n-1)*sd/sqrt(n) } konf.int_felső = function(n, xvonas, sd, alpha) { xvonas+qt(1-alpha/2, df = n-1)*sd/sqrt(n) } konf.int_also(elemszam, mintaatlanlag, mintaszoras, alpha) konf.int_felső(elemszam, mintaatlanlag, mintaszoras, alpha)</pre>		
p-érték alapján		$p < \alpha$
<pre>p = pt(abs(t), df = elemszam-1, lower.tail = F)*2</pre>		

Elsőfajú, másodfajú hiba

	H0-t "elfogadjuk"	H0-t elutasítjuk
H0 igaz	ok	P(Elsőfajú hiba) = α
H0 hamis	P(Másodfajú hiba) = β	ok
próba ereje: $1 - \beta = P(\text{elutasítjuk } H_0\text{-t} \mid H_A)$		

R tips & tricks

Adat beolvasása X táblából az Y csomagból	<pre>library(Y) input = data.frame(X) attach(input) # így elérjük a mezőit könnyen</pre>
Adott állomány beolvasása	<pre>input = read.table("salary.txt", header = TRUE) bank = read.csv2("bankloan.csv") kings = scan("http://.../kings.dat", skip = 3) admission = read.csv(url)</pre>
Pontbecslések	<pre>mean(data) # átlag sd(data) # szórás var(data) # variancia median(data) # medián</pre>
Alaki jellemzők	<pre>library(moments) skewness(data) # ferdeség kurtosis(data) # lapultság</pre>
Kvartilisek, min, max, terjedelem, IQR	<pre>quantile(data) quantile(data, probs = c(0.25, 0.5, 0.75)) max(data) min(data) max(data) - min(data) IQR(data) summary(data)</pre>
Üres elemek kiszűrése	<pre>result = data[!is.na(data)]</pre>
Üres elemek számossága	<pre>sum(is.na(data))</pre>
Új táblázat létrehozása meglévő adatok alapján	<pre>result = table(data1, data2) # hasznosabb result = data.frame(data1, data2)</pre>
Pontbecslések valami szerint csoportosítva (Bármilyen függvény alkalmazása a cellákra)	<pre>tapply(data, group, function, na.rm = T) # például az életkorok átlaga nemek szerint tapply(Age, Sex, mean)</pre>
Kiugró értékek	<pre>boxplot.stats(data)\$out</pre>

Sorbrarendezés	<pre>res = data[order(data\$by),] # vagy vektor esetén res = data[order(data)]</pre>
Szűrés valamilyen feltétel alapján	<pre>data[data\$by == val,]</pre>
Standard normális eloszlás generálása	<pre>set.seed(1234) minta = rnorm(1000)</pre>
Standardizálás, QQ-ábra	<pre>std_salary = scale(salary) qqnorm(std_salary) qqline(std_salary)</pre>

Vizsgálatok és megvalósításuk R-ben

Várható érték konkrét szám $\mu = c$	Egymintás t-próba	<code>t.test(var, mu = val)</code>
<pre># 5 százalékos szignifikancia szinten teszteljük azt a nullhipotézist, hogy # az átlagos pulzusszám 75-tel egyenlő. t.test(Pulse, mu = 75)</pre>		
Összefüggő minták (= 2) $\mu_1 = \mu_2$	Páros t-próba	<code>t.test(var1, var2, paired = T)</code> <code>t.test(var1-var2, mu = 0)</code>
<pre># Hasonlítsuk össze a két kéz fesztávolságát. 5 százalékos szignifikancia szinten # teszteljük azt a nullhipotézist, hogy # a két kéz átlagos fesztávolsága megegyezik. t.test(Wr.Hnd, NW.Hnd, paired = T) t.test(Wr.Hnd-NW.Hnd, mu = 0)</pre>		
Független minták (= 2) $\mu_1 = \mu_2$	Kétmintás t-próba (F-próba + t-próba)	<code>var.test(var ~ group, conf.level = 0.99)</code> <code>t.test(var ~ group, var.equal = T, conf.level = 0.99)</code> <code>t.test(var1, var2, var.equal = T, conf.level = 0.99)</code>
<pre># 1 százalékos szignifikancia szinten teszteljük azt a nullhipotézist, hogy # az átlagos pulzusszám azonos férfiak és nők esetén. # Adjunk meg egy 99 százalékos megbízhatóságú konfidencia intervallumot # a két várható érték különbségére. var.test(Pulse ~ Sex, conf.level = 0.99) t.test(Pulse ~ Sex, var.equal = T, conf.level = 0.99)</pre>		
Független minták (> 2) $\mu_1 = \mu_2 = \dots = \mu_n$	Egyszempontos ANOVA (Levene-próba + ANOVA)	<code>leveneTest(var ~ group, center = mean)</code> <code>oneway.test(var ~ group, var.equal = TRUE)</code> <code>factor(group)</code>
<pre># Vizsgáljuk meg a pulzusszámokat a testmozgás gyakorisága (Exer változó) szerint. # Teszteljük azt a nullhipotézist, hogy a pulzusszámok várható értéke # minden csoportban ugyanakkora. library(car) leveneTest(Pulse ~ Exer, center = mean) summary(aov(Pulse ~ Exer)) oneway.test(Pulse ~ Exer, var.equal = TRUE)</pre>		
Egy változó normális eloszlást követ-e?	Egymintás Kolmogorov-Szmirnov-próba	<code>ks.test(var, "pnorm", mean(var), sd(var))</code>
<pre># A Kolmogorov-Szmirnov-próba alkalmazásával teszteljük le # 5 százalékos szignifikanciaszinten azt a nullhipotézist, hogy # a 'salary' változó normális eloszlást követ. ks.test(salary, "pnorm", mean(salary), sd(salary)) # p < 0.05 => salary nem normális eloszlású</pre>		

Egymástól független minták eloszlásának összehasonlítása	Kétmintás Kolmogorov-Szmirnov-próba	ks.test(var[cond1], var[cond2])
<pre># Kétmintás Kolmogorov-Szmirnov-próba segítségével teszteljük le # azt a nullhipotézist, hogy a férfiak és a nők körében azonos # a fizetés eloszlása. ks.test(salary[gender == "Male"], salary[gender == "Female"])</pre>		
A megfigyelt vs. várt gyakoriságok	khi ² -próba binomiális-próba	chisq.test(gyak, p = vals) binom.test(egyik, ossz, vals)
<pre># Egy dobókockával százszor egymás után dobva a következő gyakoriságokat kaptuk: # 1: 15, 2: 16, 3: 14, 4: 15, 5: 20, 6: 20 # Teszteljük azt a nullhipotézist, hogy a dobókocka szabályos. vals = rep(1/6, 6) gyak = c(15, 16, 14, 15, 20, 20) chisq.test(gyak, p = vals) # p = 0.8323 > 0.05 => nem tudjuk elutasítani a H0-t, vagyis # nem tudjuk igazolni, hogy a dobókocka nem szabályos # ===== # Egy felmérés során a megkérdezett 127 ember közül 51 az aktuális polgármestert # támogatja, 76 pedig a másik jelöltet. # Mondhatjuk-e, hogy azonos a két jelölt támogatottsága? vals = c(0.5, 0.5) gyak = c(51, 76) chisq.test(gyak, p = vals) # p = 0.02653 < 0.05 => elutasítjuk H0-t, vagyis # a két jelölt támogatottsága szignifikánsan eltér # binomiális próba (egyik, össz, valszin) binom.test(51, 127, p = 1/2) # ===== # Az 'ertekek20' nevű vektor 20 kockadobás eredményét tartalmazza. # Teszteljük azt a nullhipotézist, hogy a dobókocka szabályos. # Milyen figyelmeztető üzenetet kapunk, és ennek mi az oka? # Ha az eredmények vektora adott, akkor abból előbb egy táblázatot kell készíteni ertekek20 = c(1,1,1,2,2,2,3,3,3,4,4,4,5,5,5,5,6,6,6,6) gyak20 = table(ertekek20) vals = rep(1/6, 6) chisq.test(gyak20, p = vals) # nem elég nagy az elemszám # van olyan gyakoriság, ami < 5 # túl alacsony a mintavétel, nem tudunk pontos eredményt adni # 2 elem esetén lehet binomiális próbát alkalmazni!</pre>		