

Experimental Design and Statistics - AGA46E

M. Maciak (Czech University of Life Sciences, Prague)

Lab Session 7 - Summer Term 2015

1 Dvouvyberove problemy - teoreticke priklady

- Predpokladajme, ze vyska muzu v populaci ma normalne rozdeleni $N(\mu_1, 40)$ a vyska zen v populaci ma normalne rozdeleni $N(\mu_2, 40)$ (predpokladame tedy homoskedasticitny pripad). Na zaklade nahodnych vyberu $X_1, \dots, X_{25} \sim N(\mu_1, 40)$ a $Y_1, \dots, Y_{30} \sim N(\mu_2, 40)$ jsme spocetli vyberove prumery $\bar{X}_{25} = 182 \text{ cm}$ a $\bar{Y}_{30} = 176.8 \text{ cm}$. Sestrojte konfidencny interval pro rozdil $\mu_1 - \mu_2$ na hladine spolehlivosti 95 %. Na stejne hladine take testujte nulovou hypotezu $H_0 : \mu_1 = \mu_2$ oproti alternative $H_A : \mu_1 \neq \mu_2$.
- Uvazujme stejný pripad s vyskou muzů a žen v populaci, avšak parametr rozptylu $\sigma^2 > 0$ je neznámý (nadále však predpokladame homoskedasticitu). Vyberove rozptyly jsou $s_{25}^2 = 42$ a $s_{30}^2 = 38$. Sestrojte 95 % konfidencny interval za techto predpokladu a take testujte nulovou hypotezu $H_0 : \mu_1 = \mu_2$ oproti alternative $H_A : \mu_1 \neq \mu_2$.
- Jak se reseni zmeni, jestliže budeme predpokladat heteroskedasticitny pripad, tudiz nestejne rozptyly, avšak vyberove rozptyly s_{25}^2 a s_{30}^2 budu stejne, jako v predchozim pripadu?
- Otestujte nulovou hypotezu o homoskedasticite v predchozim pripade ($H_0 : \sigma_1^2 = \sigma_2^2$). Uvazujte hladinu spolehlivosti 90 %.

2 Dvouvyberove problemy v Rku

Pouzite datovy soubor `passengerData3.RData` a nactete data do softwaru R pomoci prikazu `read.csv()`; Využite nektere zakladne popisne statistiky v Rku (summary statistics) aby ste získali zakladny přehled o celkovych datech. Pak udelejte nasledujici:]

- Použite `help` v Rku z zjistete, jak funguje funkce (prikaz) `t.test()`; Zjistete, jaké dodatecne parametre jsou potrebne a jakym způsobem se specifikují jednotlivé případy, které pro statisticky dvouvyberovy *t*-test rozlišujeme.
- Najdete 95% interval spolehlivosti pro stredni očekavanou výšku mužských a ženských pasazerů; Zároveň otestujte nulovou hypotezu, že mezi očekavanou výškou mužů a žen není žádný rozdíl (použijte $\alpha = 0.05$ a funkci `t.test()` která je v Rku). Jak by vypadal příslušný 95% interval spolehlivosti pro rozdíl v očekávané výšce mužů a žen?
- Uvazujte průměrnou dobu letu a průměrnou čekací dobu spocetenu za období jednoho měsíce (12 měsíců dohromady). Lze říct, že očekávaná doba čekání je kratší, než očekávaná doba letu? Použite hladinu významnosti $\alpha = 0.05$.
- Uvazujte čekací doby samostatně pro mužské a ženské pasazery a testujte nulovou hypotezu, že očekávaná doba čekání je stejná pro muže i pro ženy. Zvolte rozumnou hladinu významnosti $\alpha \in (0, \frac{1}{2})$; Jak vypada příslušný konfidencny interval spolehlivosti pro rozdíl v očekávané době čekání mužů a žen
- Stejně příklady se pokuste spocítat manuálně, pomoci vzorečku, které byli na přednášce. Pro příslušné kvantilové hodnoty použijte tabulky příslušných rozdelení, nebo statisticky software R.

Strucny prehled prednasky

- pro nahodne vybery $X_1, \dots, X_n \sim N(\mu_1, \sigma_1^2)$ a $Y_1, \dots, Y_m \sim N(\mu_2, \sigma_2^2)$ plati:

$$\begin{aligned}\bar{X}_n &= \frac{1}{n} \sum_{i=1}^n X_i & \bar{Y}_m &= \frac{1}{m} \sum_{i=1}^m Y_i \\ s_{X'}^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 & s_{Y'}^2 &= \frac{1}{m-1} \sum_{i=1}^m (Y_i - \bar{Y}_m)^2\end{aligned}$$

- kdyz jsou nahodne vybery zavisle (paired samples), pak nutne plati, ze $n = m$ a zaroven

$$Z_i = X_i - Y_i \sim N(\mu_Z, \sigma_Z^2), \text{ for } \mu_Z = \mu_1 - \mu_2 \text{ and } \sigma_Z^2 = \sigma_1^2 + \sigma_2^2.$$

Prislusny pocet stupnu volnosti je $df = n - 1$.

- kdyz jsou oba nahodne vybery vzajemne nezavisle a s nerovnymi rozptyly $\sigma_1^2 \neq \sigma_2^2$, pak plati

$$\bar{X}_n - \bar{Y}_m \sim N\left(\mu_1 - \mu_2, \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}\right), \text{ so stupnemi volnosti } df = \frac{\left(\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}\right)^2}{\frac{\sigma_X^4}{n^2(n-1)} + \frac{\sigma_Y^4}{m^2(m-1)}}$$

- kdyz jsou nahodne vybery nezavisle se stejnými rozptyly $\sigma_1^2 = \sigma_2^2$ (pouze predpoklad, zjiskane odhady rozptylu muzu byt obecne ruzne), pak plati

$$\bar{X}_n - \bar{Y}_m \sim N\left(\mu_1 - \mu_2, \sigma_{XY}^2 \left(\frac{1}{n} + \frac{1}{m}\right)\right), \text{ for } \sigma_{XY}^2 = \frac{(n-1)\sigma_X^2 + (m-1)\sigma_Y^2}{m+n-2},$$

co se take nazyva **pooled estimate** s prislusnými stupnemi volnosti $df = n + m - 2$.