

# statistika-testy

2023-03-07

1.) Načtěte soubor vokaly.xls a vytvořte kontingenční tabulku počtu zastoupení vokál vs. pohlaví. Kolikrát je zastoupen každý vokál u mužů a kolikrát u žen?

```
tab <- read_xls("vokaly.xls")
tab_cont <- table(tab$Vowel, tab$Sex)
tab_cont
```

```
##
##      f      m
## a 384 216
## e 384 216
## i 384 216
## o 384 216
## u 384 216
```

2.) Vyberte dvě skupiny (každou uložte do samostatné proměnné): Vokál [a] muži, vokál [a] ženy. Sdílejte počty řádků těchto tabulek s počty zjištěnými v bodě 1?

```
men_a <- filter(tab, Vowel == "a", Sex == "m")
women_a <- filter(tab, Vowel == "a", Sex == "f")
nrow(men_a)
```

```
## [1] 216
```

```
nrow(women_a)
```

```
## [1] 384
```

3.) Proveďte test, zda hodnoty F2\_Hz obou skupin pocházejí z normálního rozdělení.

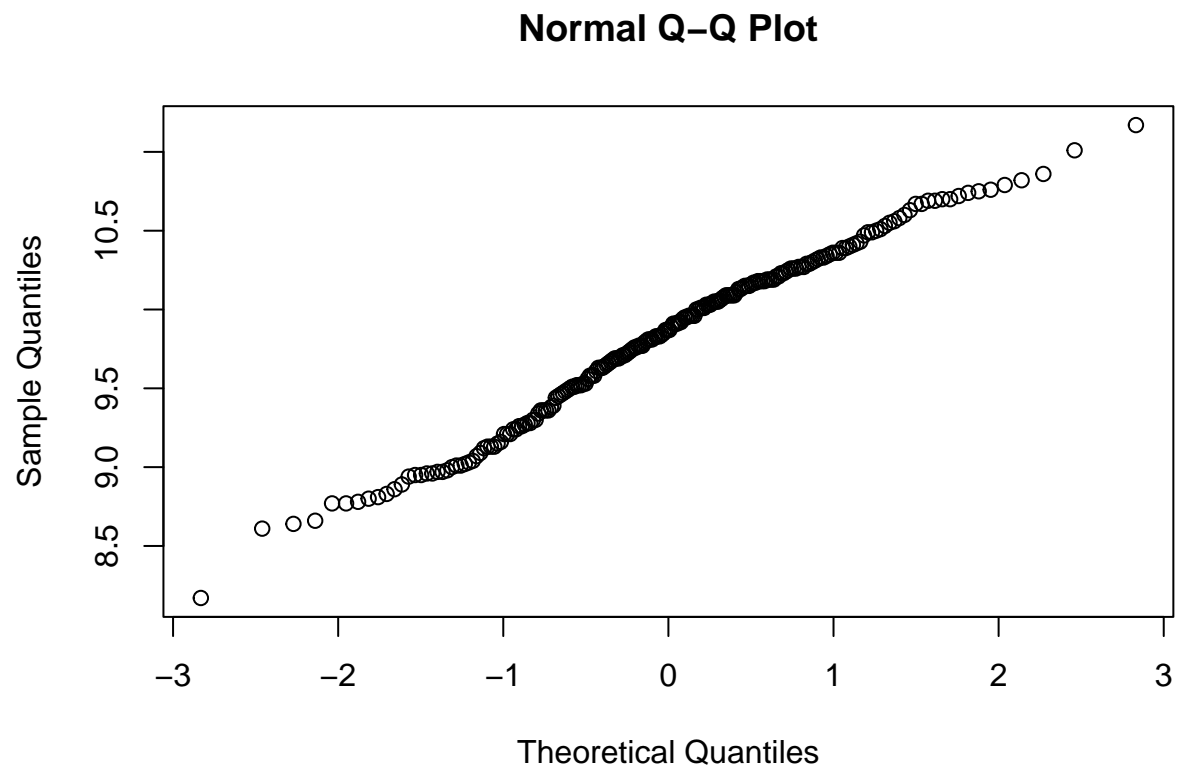
```
shapiro.test(men_a$F2_B)
```

```
##
##  Shapiro-Wilk normality test
##
## data:  men_a$F2_B
## W = 0.98684, p-value = 0.04327
```

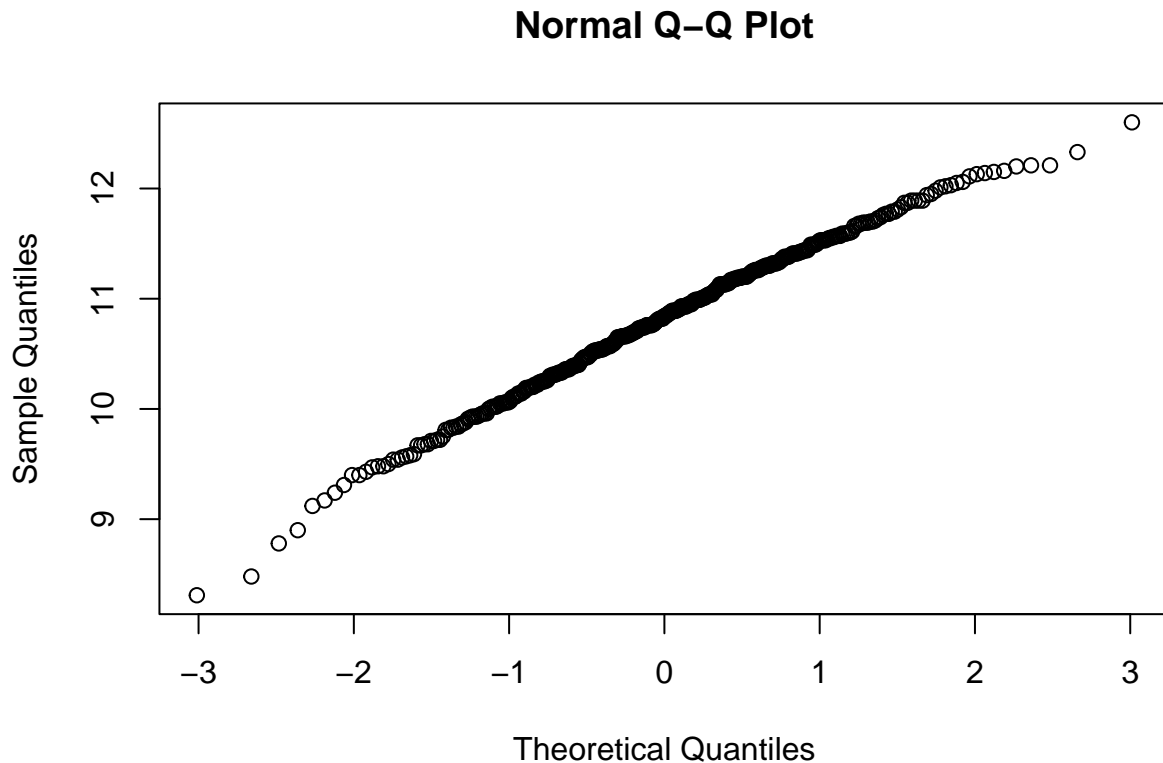
```
shapiro.test(women_a$F2_B)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data:  women_a$F2_B  
## W = 0.99244, p-value = 0.04903
```

```
qqnorm(men_a$F2_B)
```



```
qqnorm(women_a$F2_B)
```



4.) Pomocí vhodně zvoleného testu rozhodněte mezi hypotézami při  $\alpha = 0.05$ . Vypočtěte i p-hodnotu.  
 $H_0$ : střední hodnoty F2 žen a mužů u tohoto vokálu jsou stejné.  $H_1$ : střední hodnoty jsou různé.

```
x <- mean(women_a$F2_B)
y <- mean(men_a$F2_B)

t.test(women_a$F2_B, men_a$F2_B, paired = FALSE, conf.level = 0.95)

##
##  Welch Two Sample t-test
##
## data:  women_a$F2_B and men_a$F2_B
## t = 18.747, df = 533.94, p-value < 2.2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  0.8775731 1.0830171
## sample estimates:
## mean of x mean of y
## 10.802656  9.822361
```

#### Úloha B

V percepčním testu se ti samí lidé po proškolení ve svých výsledcích v 5 případech zhoršili, ve 35 zlepšili. Je to statisticky významná změna při  $\alpha = 0.05$ ?  $H_0$ : není to změna  $H_1$ : lidé se zlepšili

```
p <- min(1, 2*dbinom(5, 40, 0.5))
p
```

```
## [1] 1.382612e-06
```

### Úloha C

V roce 2013 nám na ústav přišlo 5 dotazů týkajících se intonace, 15 zaměřených na hlásky a 16 ohledně IPA transkripce. V roce 2014 to bylo 7 dotazů na intonaci, 7 na hlásky a 10 ohledně IPA. Je evidentní, že v roce 2014 přišlo celkově méně dotazů, nás ale spíše zajímá poměr mezi četnostmi v jednotlivých třech kategoriích. Jde o to, zda se poměry v jednotlivých dvou rocích liší, neboli zda je nějaká souvislost mezi rokem a typem dotazu. H0: poměry mezi kategoriemi jsou v obou rocích stejné – rozdíly jsou jen vlivem náhody, neboli nezávisí na konkrétním roce. H1: poměry se statisticky významně změnily, neboli hodnoty v jednotlivých kategoriích jsou závislé na tom, zda jsou v jednom či druhém roce.

```
r2013 <- c(5, 15, 16)
r2014 <- c(7, 7, 10)

tabulka <- rbind(r2013, r2014)
cnames <- c("intonace", "hlasky", "IPA")
colnames(tabulka) <- cnames
tabulka
```

```
##      intonace hlasky IPA
## r2013         5     15  16
## r2014         7      7  10
```

```
chisq.test(tabulka)
```

```
## Warning in chisq.test(tabulka): Chi-squared approximation may be incorrect
```

```
##
## Pearson's Chi-squared test
##
## data:  tabulka
## X-squared = 2.3198, df = 2, p-value = 0.3135
```

```
fisher.test(tabulka)
```

```
##
## Fisher's Exact Test for Count Data
##
## data:  tabulka
## p-value = 0.318
## alternative hypothesis: two.sided
```

### Úloha D: BONUS na zamyšlenou (pro významné plus k zápočtu)

Představte si situaci, kdy je prováděn percepční test (položka zní stejně / zní odlišně), a to pro tři nezávislé skupiny slov (každá skupina obsahuje 240 slov). Byly zjištěny tyto četnosti úspěchů (počet neúspěchů je v rámci každé skupiny vždy zbytek do 240). skupina 1: 88 skupina 2: 95 skupina 3: 118 Jak vyhodnotit,

zda mají jednotlivé skupiny rozdílné poměry úspěšnosti nebo ne? Nápad A) Nasadit na to test chí-kvadrát s hypotézou  $H_0$ : jednotlivé skupiny mají stejnou četnost úspěšných rozpoznání,  $H_1$ : četnosti jsou různé. Samozřejmě, bavíme se o odhadu vlastností základního souboru, jde nám tedy o to, zda pozorované rozdíly v tomto výběru jsou jen dílem náhody nebo zda jsou statisticky významné. Nápad B) Pro každou kategorii zvlášť vypočítat interval spolehlivosti pro binomické rozdělení (viz kapitola 12. Intervalové odhady). V každé kategorii samostatně takto počítáme poměr úspěšnosti (např. skupina 1:  $88/240 = 36.6\%$ ), ale díky intervalu spolehlivosti k tomuto bodovému odhadu přidáme statistické vyhodnocení nejistoty, tedy `binom.test(88, 240)$conf.int` => poměr základního souboru se na 95 % nachází mezi 30.5 % a 43.1 %. Takto vypočteme intervaly pro každou kategorii zvlášť a pokud se některé intervaly nebudou překrývat, prohlásíme dané skupiny za významně rozdílné.

```
# skupina 1: 88
# skupina 2: 95
# skupina 3: 118
# velikost souboru 240

binom.test(88, 240)$conf.int
```

```
## [1] 0.3056046 0.4310672
## attr(,"conf.level")
## [1] 0.95
```

```
binom.test(95, 240)$conf.int
```

```
## [1] 0.3335122 0.4607611
## attr(,"conf.level")
## [1] 0.95
```

```
binom.test(118, 240)$conf.int
```

```
## [1] 0.4267816 0.5567601
## attr(,"conf.level")
## [1] 0.95
```