

Necht náh. vektor $(X, Y)^T \sim N_2(\mu, \Sigma)$, kde $\mu = (\mu_1, \mu_2)^T$ je vektor středních hodnot a $\Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}$ je varianční matice. K vyjádření vztahu mezi náh. veličinami X a Y použijeme klasický / lineární Pearsonův korelační koeficient

$$R = \frac{1}{n-1} \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{S_x^2 S_y^2}} = \frac{1}{n-1} \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{S_x S_y}$$

kde $S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ a $S_y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$

Z bakalářského studia víme, že Pearsonův korelační koeficient pochází asymptoticky z normálního rozdělení, konkrétně

$$R \stackrel{A}{\sim} N\left(\underbrace{\rho}_{\text{střední hodnota}}, \underbrace{\frac{(1-\rho^2)^2}{n-1}}_{\text{rozptyl}}\right)$$

Konvergence R k normalitě je však pomalá, proto používáme tzv. Fisherovu Z -proměnnou

$$Z_R = \frac{1}{2} \ln \frac{1+R}{1-R} \sim N\left(\underbrace{\frac{1}{2} \ln \frac{1+\rho}{1-\rho}}_{\text{střední hodnota}}, \underbrace{\frac{1}{n-3}}_{\text{rozptyl}}\right)$$

V příkladech 7.1 a 7.2 se zaměříme na rychlost konvergence statistik R a Z_R k normálnímu rozdělení při $n \rightarrow \infty$ (7.1) a při $\rho \rightarrow 1$ (7.2). Příklady vyřešíme společně vytvořením funkce `rho.stat()`, která pro dané $n, \rho, \mu_1, \mu_2, \sigma_1, \sigma_2$ a M vykreslí dvojici histogramů (pro R + pro Z_R). Tato funkci pak použijeme k vykreslení animací v obou příkladech.

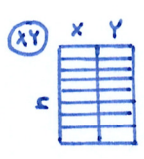
```
rho.stat <- function(n, rho, mu1=..., mu2=..., sigma1=..., sigma2=..., M=...,
                    xlim1=c(0.5, 1), xlim2=c(-0.5, 2), ylim1=c(0, 10), ylim2=c(0, 4)) {
  # rozbah oyx prvního hist.      rozbah oyx druhého hist.      rozbah oyx y prvního hist.      rozbah oyx y druhého hist.
```

`mu <- c(...)` ... vektor středních hodnot $\mu = (\mu_1, \mu_2)$

`Sigma <- matrix(...)` ... varianční matice $\Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}$... ! po vytvoření matice Sigma si ji vypiše a zkontroluje! často se v ní chybuje!

Generování dat: I. způsob - cyklem:

```
R <- ZR <- NULL ... příprava prázdných vektorů R a ZR
for (i in 1:M) {
  XY <- MASS::mvrnorm(n=n, mu=mu, Sigma=Sigma) ...
  X <- XY[, 1] ... vektor X (délka=n)
  Y <- XY[, 2] ... vektor Y (délka=n)
  R[i] <- cor(...) ... R (1 číslo; délka R[i]=1)
  ZR[i] <- ZR ... (1 číslo, délka ZR[i]=1)
}
```



Generování dat: II. způsob - funkcí `replicate()`:

```
R <- replicate(M, cor(mvrnorm(n=n, mu=mu, Sigma=Sigma)) [1, 2]) ... vektor M=1000 korel. koef. R
ZR <- 1/2 * ln(1+R)/(1-R) ... vektor M=1000 Fisherových Z-transformací.
```

Výstupem je korelační matice!!! $\begin{pmatrix} 1 & R \\ R & 1 \end{pmatrix}$

Generování dat (vyzkoušejte si oba způsoby)

Histogram R

```

xfit <- seq(...) ... posloupnost od -1 do 1 s delkou minimalni 500
yfit <- dnorm(...) ... hustota N(g, (1-g^2)^2) nad posloupnaci xfit
par(mfrow=c(2,2), mar=c(4,5,2,2))
hist(..., prob=..., col=..., border=..., xlim=xlim1,
      ylim=ylim1, ...) ... histogram R se srafovanymi slupci, rozsah osy x(xlim) je radan argumentem xlim1,
      y(ylim) -||- ylim1
box(...) □
lines(...) ... křivka hustoty N(g, (1-g^2)^2) se šírkou 2
mtext(...) ... popisak R
mtext(bquote(paste(...)), ...) ... automaticky se měnící popisak n=..., g=...

```

Histogram ZR

```

xfit <- seq(...) ... posloupnost od min(ZR)-1 do max(ZR)+1 s delkou minimalni 500
yfit <- dnorm(...) ... hustota N(1/2 ln (1+g)/(1-g), 1/(n-3)) nad posl. xfit
hist(..., xlim=xlim2, ylim=ylim2, ...) ... histogram ZR
box(...) □
lines(...) ... křivka hustoty N(1/2 ln (1+g)/(1-g), 1/(n-3))
mtext(...) ... popisak ZR
mtext(...) ... automaticky se měnící popisak n=..., g=...

```

S5.1: Animace:

```

n <- seq(...) ... posloupnost n hodnot podle radani príkladu 5,10,15,...,70.
... úvodní nastavení animace
saveLatex(for(...){
  rho.stat(n=n[i], rho=0.8)
},...)

```

Vytvořím animaci si pomalu projít a porovnejte kvalitu konvergence R a Z_R k normálnímu rozdělení. Do komentáře uveďte rozhodnutí kvality konvergence a rovně, pro jak velké n je podle vašeho názoru již vhodné použít Pearsonův korelační koeficient R a Fisherovu Z-transformaci Z_R.

S5.2: Animace pro n=5:

```

rho <- seq(...) ... radaná posloupnost g = 0.1, 0.2, ..., 0.9
... úvodní nastavení animace
saveLatex(for(...){
  rho.stat(n=5, rho=rho[i], xlim1=c(0,1), xlim2=c(-2,3),
           ylim1=c(0,5), ylim2=c(0,1))
},...)

```

S5.2: Animace pro n=50:

Analogicky jako pro n=5. Hodnoty ostatních argumentů rovle: xlim1=c(0,1), xlim2=c(-2,3), ylim1=c(0,12), ylim2=c(0,3).

Vytvořím animaci si pomalu projít a porovnejte kvalitu konvergence R a Z_R k normálnímu rozdělení. Do komentáře popište, jak se mění kvalita konvergence R a Z_R při g → 0.9 (mění-li se nějak) a je-li nějaký rozdíl mezi situací když n=5 a n=50.

V rámci tohoto příkladu si vybereme a aplikaci testovacích statistik ZW a ULR na reálná data. Nejprve ověříme předpoklad dvourozměrné normality. Následně dostaneme H_0 reordání, a to řešíme třemi způsoby (krit. oborem, IS a p-hodnotou) pro každou test. statistiku ZW a ULR. Nerapomeneme nikdy uvést návrh H_0 a nakonec interpretaci návrhu testování. Nakonec vykrájíme hranici a oblast rozhodnostního 95% empirického IS.

Načtení a příprava dat

```
data <- read.delim(...)
data.M <- data[date$sex == ..., c('skull.pH', 'face.H')]
data.M <- na.omit(...)
skull.M <- data.M$...
face.M <- data.M$...
```

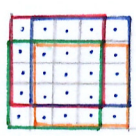
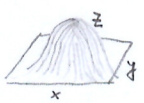
Test(y) dvourozměrné normality

```
# H0: ...
# H1: ...
MVN::mvn(data.M, mvnTest = 'mardia')$multi ... MARDIUV TEST ... testuje náležitě očekávané úsklonosti a špicovitosti
MVN::mvn(..., mvnTest = 'hz')$... ... HENZE-ZIRKLERUV TEST
MVN::mvn(..., mvnTest = 'royston')$... ... ROYSTONUV TEST
# Závěr: Probaře..., H0..., na hladině významnosti...
# Interpretace:
```

VYVÝSKOVÉ
S
-
Y
M
E
R
N
Y
T
M
E
T
Y

3D (persp) graf

```
kd <- MASS::kde2d(skull.M, face.M, n=50, lims = c(120, 155, 90, 145)) ... jádrový odhad dvourozměrné hustoty (mod osou X n normální 120-155 a mod osou y n normální 90-145).
x <- kd$x ... některé hodnoty x vybrané pro kde2d
y <- kd$y ... některé hodnoty y vybrané pro kde2d
z <- kd$z ... matice dvourozměrné hustoty f(x,y) vybraná pro kde2d
N <- dim(z)[1] (50)
stredy <- z[-1,-1] + z[-1,-n] + z[...,...] + z[...,...] ... analogická funkce na určování středů jako v SI 1.
vyska <- cut(stredy, 12) ... rozdělení středů do 12 kategorií
par(mar=...) ... stanovení hraničí kategorií jedním z 12 barev šedé heat.colors().
persp(..., theta = -20, phi = 30, col = heat.colors(12)[vyska], ...) ... 3D jádrový odhad hustoty
```



z[-N,-N]
z[-N,-1]
z[-1,-N]
z[-1,-1]

Tečkový graf s 95% elipsou spolehlivosti

```
par(...) ... 4x5 grid
car::dataEllipse(skull.M, face.M, level=0.95, xlim=c(122, 150), ylim=c(90, 141), pch=21, bg='khaki1', col=c('orange2', 'khaki4'), ...) ... Graf elipsy spolehlivosti
mtext(...) ... popisek 'největší výška mozku (vmm)' ... barva šedé bodů
mtext(...) ... popisek 'Mardiův test ... (nemusí být automatizovaný)'. ... barva druhé barvy elipsy
```

Poznámka: Interpretace normality na náklade tečkového grafu s elipsou spolehlivosti. Předpokládáme, že pokud data patří z dvourozměrného normálního rozdělení N_2 , potom 95% elipsa spolehlivosti pokrývá alespoň 95% dat. (tj. nejvýše 5% dat leží mimo elipsu spolehlivosti).

Zjistíte tedy rozsah náhodného výběru n, spočítáte 5% z n, dále spočítáte, kolik bodů leží vně elipsy spolehlivosti a do komentářů uvedete relevantní návrh o dvourozměrné normalitě dat na náklade tohoto grafu.

Test o korelačním koeficientu g:

Ho:

H1:

Příprava hodnot

```
rho0 <- ... g0
n <- length(...) ... počet pozorování (rozsah mat. výběru)
alpha <- ... l
r <- cor(...) ... r ... realizace Pearsonova korelačního koeficientu R pro skull. Ma face. M
zR <- 1/2 * ln((1+R)/(1-R)) ... zR
ksi0 <- 1/2 * ln((1+g0)/(1-g0)) ... xi0
```

Funkce počítající hodnotu test. stat. ULR

```
ULR.stat <- function(r, rho0, n){
  ULR <- n * (ln((1 - rho0^2)^2 / ((1 - r^2)(1 - g0^2))))
  return(...)}

```

Testovací statistiky

```
Testování kritickým oborem
zW <- sqrt(n-3) * (zR - xi0) ... zW
ULR <- ULR.stat(...) ... hodnota test. statistiky ULR
```

Hranice kritických oborů

```
q1 <- ... u_{alpha/2}
q2 <- ... u_{1-alpha/2}
q <- ... chi^2_{alpha}(1-l)
```

Pro každý test uveďte nulový návrh a Ho

Testování IS:

Intervaly spolehlivosti

```
dh.zR <- zR - u_{alpha/2} / sqrt(n-3) ... dolní hranice Waldova IS pro xi0
hh.zR <- zR + u_{alpha/2} / sqrt(n-3) ... horní hranice Waldova IS pro xi0
dh.R <- tanh(dh.zR) ... zpětně transformovaná dolní hranice Waldova IS (pro g)
hh.R <- tanh(hh.zR) ... zpětně transformovaná horní hranice Waldova IS (pro g)
rho.i <- seq(...) ... posloupnost g; od 0.1 do 0.6 se vzdáleností mezi body 0.0001 (5001 hodnot)
ULR.i <- ULR.stat(r, ...) ... vektor 5001 test. statistik ULR,i pro posloupnost g; (rho.i)
dh.ULR <- min(...) ... dolní hranice věroh. 95% empirického IS
hh.ULR <- max(...) ... horní hranice -11-
} princip je analogický jako u 6.6 a 5.8

```

Pro každý test uveďte nulový návrh a Ho:

Testování p-hodnotou

p-hodnoty

```
p.zW <- 2 * min(pnorm(...), 1-pnorm(...))
p.ULR <- 1 - pchisq(...)
```

Pro každý test uveďte nulový návrh a Ho.

Interpretace výsledků testování: Uveďte antropologický návrh (interpretaci výsledků testování). Co jsme se dověděli o vztahu největší výšky mužovny a morfologické výšky kráň u mužů starověké egyptské populace?

Na návrh uveďte také kompletní interpretaci Pearsonova korelačního koeficientu r=0.3306 tak, jak jsme si ji uváděli na začátku semestru.

Tabulka výsledků

```

tab <- data.frame(...
  IS.dh=c(dh.R, dh.ULR),
  ...)

```

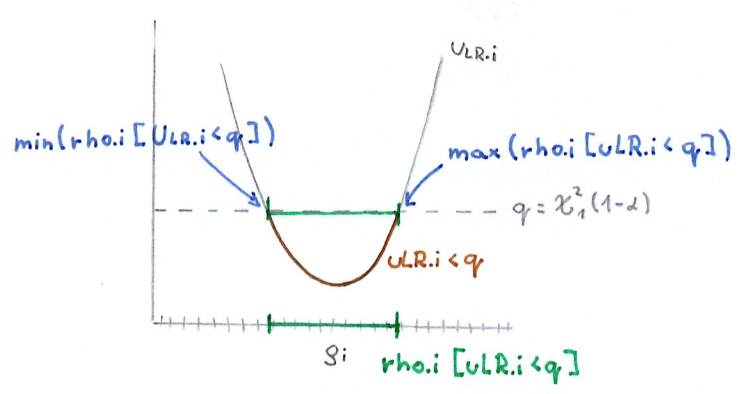
round(tab, 4) ... zaokrouhlení na 4 des. místa

Graf s vřehodnostním IS

```

par(...)  $\begin{matrix} 4 \\ \square \\ 4 \end{matrix}$ 
plot(rho.i, uLR.i, ylim=c(0,5),...) ... křivka ULR √
mtext(...) ... popisek  $\varrho$ 
abline(...) ... vodorovná čára přesuřovaná referenční čára v hodnotě  $\chi^2_{1-d}$  ---
lines(rho.i[...], uLR.i[...],...) ... barevně vyznačená oblast ULR v IS √
mtext(...) ... popisek IS = (... , ...), analogický nápis jako v 5.8 a 6.6.

```



Odvrození:

(Odvrození je analogické jako ve cvičení 04. Doporučuji napsat ke zde uvedených vzorců síly, vykoušet si odvození samostatně a poté uvedení odvození porovnat jako kontrolu. =))

1. $H_{01}: g = g_0$ $H_{11}: g \neq g_0$ (obousměrná alb.)

Exaktní síla: $\beta_{11}^*(\xi) = \Phi(u_{\alpha/2} + \sqrt{n-3} (\xi_0 - \xi)) + \Phi(u_{\alpha/2} - \sqrt{n-3} (\xi_0 - \xi))$

Aproximativní síla: $\tilde{\beta}_{11}^*(\xi) = \Phi(u_{\alpha/2} + \sqrt{n-3} |\xi_0 - \xi|)$

Minimální rozsah n:

$\tilde{\beta}_{11}^*(\xi) = \Phi(u_{\alpha/2} + \sqrt{n-3} |\xi_0 - \xi|) / u(\cdot)$

$u_{\tilde{\beta}_{11}^*}(\xi) = u_{\alpha/2} + \sqrt{n-3} |\xi_0 - \xi|$

$$\left(\frac{u_{\tilde{\beta}_{11}^*}(\xi) - u_{\alpha/2}}{|\xi_0 - \xi|}\right)^2 + 3 = n \Rightarrow n \geq \left(\frac{u_{\tilde{\beta}_{11}^*}(\xi) - u_{\alpha/2}}{|\xi_0 - \xi|}\right)^2 + 3 \quad / u_{\alpha/2} = -u_{1-\alpha/2}$$

$$n \geq \left(\frac{u_{\tilde{\beta}_{11}^*}(\xi) + u_{1-\alpha/2}}{|\xi_0 - \xi|}\right)^2 + 3$$

2. $H_{02}: g \leq g_0$ $H_{12}: g > g_0$ (pravosměrná alb.)

Exaktní síla: $\beta_{12}^*(\xi) = 1 - \Phi(u_{1-\alpha} + \sqrt{n-3} (\xi_0 - \xi)) = \Phi(u_{\alpha} - \sqrt{n-3} (\xi_0 - \xi))$

Minimální rozsah n:

$\beta_{12}^*(\xi) = 1 - \Phi(u_{1-\alpha} + \sqrt{n-3} (\xi_0 - \xi))$

$1 - \beta_{12}^*(\xi) = \Phi(u_{1-\alpha} + \sqrt{n-3} (\xi_0 - \xi)) / u(\cdot)$

$u_{1-\beta_{12}^*}(\xi) = u_{1-\alpha} + \sqrt{n-3} (\xi_0 - \xi)$

$$\left(\frac{u_{1-\beta_{12}^*}(\xi) - u_{1-\alpha}}{\xi_0 - \xi}\right)^2 + 3 = n \Rightarrow n \geq \left(\frac{u_{1-\beta_{12}^*}(\xi) - u_{1-\alpha}}{\xi_0 - \xi}\right)^2 + 3 \quad / u_{\alpha} = -u_{1-\alpha}$$

$$n \geq \left(\frac{u_{1-\beta_{12}^*}(\xi) + u_{\alpha}}{\xi_0 - \xi}\right)^2 + 3$$

3. $H_{03}: g \geq g_0$ $H_{13}: g < g_0$ (levosměrná alb.)

Exaktní síla: $\beta_{13}^*(\xi) = \Phi(u_{\alpha} + \sqrt{n-3} (\xi_0 - \xi))$

Minimální rozsah n:

$\beta_{13}^*(\xi) = \Phi(u_{\alpha} + \sqrt{n-3} (\xi_0 - \xi)) / u(\cdot)$

$u_{\beta_{13}^*}(\xi) = u_{\alpha} + \sqrt{n-3} (\xi_0 - \xi)$

$$\left(\frac{u_{\beta_{13}^*}(\xi) - u_{\alpha}}{\xi_0 - \xi}\right)^2 + 3 = n \Rightarrow n \geq \left(\frac{u_{\beta_{13}^*}(\xi) - u_{\alpha}}{\xi_0 - \xi}\right)^2 + 3 \quad / u_{\alpha} = -u_{1-\alpha}$$

$$n \geq \left(\frac{u_{\beta_{13}^*}(\xi) + u_{1-\alpha}}{\xi_0 - \xi}\right)^2 + 3$$

1. Nejprve si naprogramujeme funkci `min.n()`, která pro zadání hodnoty $\rho, \rho_0, \alpha, \beta^*$ a pro zvolený typ alternativy vrátí minimální rozsah máh. výzkru n .

```

min.n <- function(rho, rho0, alpha=..., sila=0.8, alternative=...) {
  ksi <- ...  $\xi = \frac{1}{2} \ln \frac{1+\rho}{1-\rho}$ 
  ksi0 <- ...  $\xi_0 = \frac{1}{2} \ln \frac{1+\rho_0}{1-\rho_0}$ 
  if(alternative == 'two.sided') {n <-  $\left( \frac{U_{\beta^*}(\xi) - U_{\alpha/2}}{|\xi_0 - \xi|} \right)^2 + 3$  }
  if(      == 'greater' ) {n <-  $\left( \frac{U_{1-\beta^*}(\xi) - U_{1-\alpha}}{\xi_0 - \xi} \right)^2 + 3$  }
  if(      == 'less' ) {n <-  $\left( \frac{U_{\beta^*}(\xi) - U_{\alpha}}{\xi_0 - \xi} \right)^2 + 3$  }
  <-  $\left( (qnorm(sila) - qnorm(alpha)) / (ksi0 - ksi) \right)^2 + 3$ 
  return(round(n))
}

```

2. Vykreslime graf pro (a) dvoustrannou alternativu:

```

rho <- c(seq(-0.95, 0.1, by=0.05), seq(0.1, 0.95, by=0.05))
n11 <- min.n(rho=rho, rho0=0, alternative='two.sided')
par(mfrow=c(1,2))
plot(rho, n11, xlim=c(-1,1), ...)
mtext(expression(rho - rho[0]), ...)
mtext(..., ...)

```

3. Vykreslime graf ((b) pro pravostannou alt. resp. (c) pro levostannou alt.)

```

rho <- seq(-0.95, 0.1, by=0.05)
n12 <- min.n(...)
plot(rho, n12, xlim=c(-1,1), ...)
segments(-2, 0, 0, 0, ...)
mtext(expression(...))
mtext(..., ...)

```

4. Spočítáme tabulku:

```

n11 <- min.n(rho=0.60, rho0=0.85, alpha=..., sila=..., alternative=...)
n12 <- min.n(...)
n13 <- min.n(...)
alpha <- c(0.10, 0.05, 0.01)
silal <- c(...)
beta <- 1 - silal
rho0 <- c(...)
rho <- c(...)
n <- c(n11, n12, n13)
tab <- data.frame(alpha, silal, ..., n, row.names=c('pravostanna', ...))
round(..., 4)

```