

Mnohonásobná lineární regrese – vzorový příklad

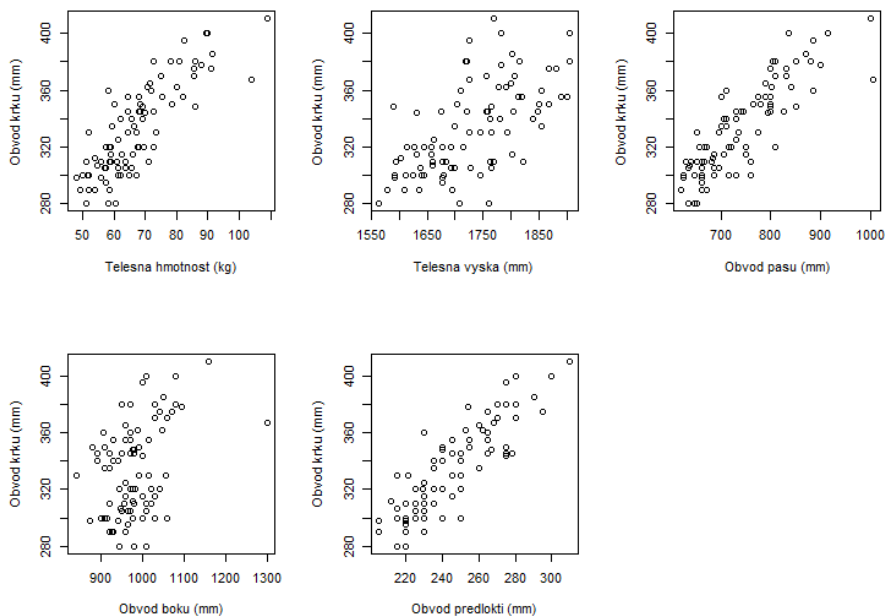
V souboru `neck.txt` máme k dispozici antropometrická data mladých dospělých lidí (převážně studentů vysokých škol z Brna a Ostravy). Chceme modelovat závislost obvodu krku (proměnná `neck.C`) na tělesné hmotnosti (proměnná `body.W`), tělesné výšce (proměnná `body.H`), obvodu pasu (proměnná `waist.C`), obvodu boků (proměnná `hip.C`) a obvodu předloktí (proměnná `antb.C`). Hmotnost byla měřena v kg, délkové míry v mm. Načteme data a podíváme se na ně. Soubor neobsahuje žádná chybějící pozorování.

```
> neck <- read.table("neck.txt", header=T)
> summary(neck)
```

	id	sex	body.w	body.H	waist.C	hip.C	antb.C	neck.C	
Min.	: 1.00	f:38	64,5	: 4	Min. :1563	Min. : 620.0	Min. : 840.0	Min. :205.0	Min. :280.0
1st Qu.	: 44.00	m:49	68	: 4	1st Qu.:1660	1st Qu.: 663.5	1st Qu.: 945.0	1st Qu.:225.0	1st Qu.:306.0
Median	: 91.00		52	: 3	Median:1725	Median: 730.0	Median: 970.0	Median:240.0	Median:330.0
Mean	: 92.23		57,5	: 3	Mean :1729	Mean : 740.1	Mean : 979.9	Mean :244.5	Mean :332.9
3rd Qu.	:139.50		58,5	: 3	3rd Qu.:1792	3rd Qu.: 800.0	3rd Qu.:1010.0	3rd Qu.:263.5	3rd Qu.:355.0
Max.	:188.00		59	: 3	Max. :1906	Max. :1005.0	Max. :1300.0	Max. :310.0	Max. :410.0

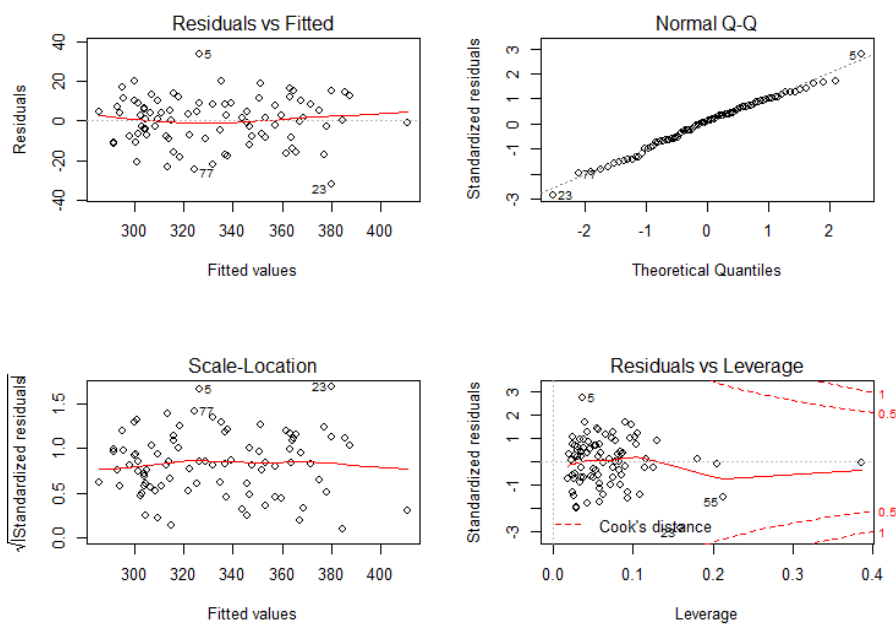
Vykreslíme si bodové diagramy pro dvojice (obvod krku, hmotnost); (obvod krku, výška); (obvod krku, obvod pasu); (obvod krku, obvod boků) a (obvod krku, obvod předloktí).

```
> par(mfrow=c(2,3))
> plot(neck$body.w, neck$neck.C, xlab='Telesna hmotnost (kg)', ylab='Obvod krku (mm)')
> plot(neck$body.H, neck$neck.C, xlab='Telesna vyska (mm)', ylab='Obvod krku (mm)')
> plot(neck$waist.C, neck$neck.C, xlab='Obvod pasu (mm)', ylab='Obvod krku (mm)')
> plot(neck$hip.C, neck$neck.C, xlab='Obvod boku (mm)', ylab='Obvod krku (mm)')
> plot(neck$antb.C, neck$neck.C, xlab='Obvod predlokti (mm)', ylab='Obvod krku (mm)')
```



Bodové diagramy naznačují, že je mezi dvojicemi lineární závislost. Sestavíme regresní model a pomocí analýzy reziduí ověříme předpoklady modelu.

```
> model1 <- lm(neck.C ~ body.W + body.H + waist.C + hip.C + antb.C, data=neck)
> par(mfrow=c(2,2))
> plot(model1)
```



Interpretace grafu je stejná jako u jednoduchého regresního modelu. Ověříme předpoklady i pomocí vhodných testů. Pomocí t-testu otestujeme hypotézu, že rezidua mají nulovou střední hodnotu. Normalitu reziduí ověříme pomocí Shapiro-Wilkova testu a nezávislost reziduí ověříme pomocí Durbinova-Watsonova testu (z knihovny car).

```
> t.test(model1$residuals)
One Sample t-test
data: model1$residuals
t = 9.3512e-17, df = 86, p-value = 1
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 -2.555173  2.555173
sample estimates:
 mean of x
1.201944e-16
```

```
> shapiro.test(model1$residuals)
Shapiro-wilk normality test
data: model1$residuals
W = 0.9902, p-value = 0.7645
```

```
> library(car)
> durbinWatsonTest(model1)
lag Autocorrelation D-W Statistic p-value
 1      0.1541587      1.678153      0.12
Alternative hypothesis: rho != 0
```

Hypotézu o nulové střední hodnotě reziduí, protože t-test nabývá hodnoty s p-hodnotou, z grafického posouzení také nevidíme problém.

Shapiro-Wilkův test nabývá hodnoty s p-hodnotou, v kvantil-kvantilovém grafu jsou rezidua, předpoklad normality tedy považujeme za

Předpoklad rovnosti rozptylů se na základě grafického posouzení zdá

Durbin-Watsonův test nabývá hodnoty s p-hodnotou, tedy nezávislost reziduí.
 Předpoklady modelu jsou tedy

Podívejme se, jestli v našem modelu není problém s multikolinearitou. Vypočítáme si korelační koeficienty mezi nezávislými proměnnými a také hodnoty koeficientu V IF pro proměnné sestaveného modelu.

```
> cor(neck[,c('body.w', 'body.H', 'waist.C', 'hip.C', 'antb.C')])
      body.w  body.H  waist.C  hip.C  antb.C
body.w  1.000000  0.6086383  0.9047087  0.7604090  0.8810742
body.H  0.6086383  1.0000000  0.4591687  0.2303759  0.5851208
waist.C 0.9047087  0.4591687  1.0000000  0.6539080  0.8520787
hip.C   0.7604090  0.2303759  0.6539080  1.0000000  0.5251877
antb.C  0.8810742  0.5851208  0.8520787  0.5251877  1.0000000
```

```
> vif(model1)
      body.w  body.H  waist.C  hip.C  antb.C
18.895276  2.307445  6.812388  3.904779  6.116750
```

Vidíme, že jak korelační koeficienty, tak koeficienty V IF nabývají vysokých hodnot, lze tedy soudit na existenci multikolinearity. Vypíšeme si podrobné informace o modelu:

```
> summary(model1)
Call:
lm(formula = neck.C ~ body.w + body.H + waist.C + hip.C + antb.C,
    data = neck)
```

Residuals:

Min	1Q	Median	3Q	Max
-32.266	-8.030	1.169	8.493	33.577

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	165.63910	64.79600	2.556	0.0124 *
body.w	1.02594	0.46867	2.189	0.0315 *
body.H	0.04039	0.02314	1.745	0.0847 .
waist.C	0.18260	0.04025	4.537	1.96e-05 ***
hip.C	-0.18166	0.04070	-4.463	2.58e-05 ***
antb.C	0.29120	0.14144	2.059	0.0427 *

 Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 12.35 on 81 degrees of freedom
 Multiple R-squared: 0.8573, Adjusted R-squared: 0.8485
 F-statistic: 97.33 on 5 and 81 DF, p-value: < 2.2e-16

MNČ odhady regresních koeficientů a jejich interpretace:

- $\beta_0 = \dots\dots\dots$
- $\beta_1 = \dots\dots\dots$
- $\beta_2 = \dots\dots\dots$
- $\beta_3 = \dots\dots\dots$
- $\beta_4 = \dots\dots\dots$
- $\beta_5 = \dots\dots\dots$

Odhadnutá regresní funkce má tvar

Index determinace ID2 =

Adjustovaný index determinace ID2adj =

Celkový F-test na hladině významnosti 0:05:

F =

p-hodnota =
 závěr

Díličí t-testy

parametr	hodnota testovací statistiky	p-hodnota	závěr
β_0			
β_1			
β_2			
β_3			
β_4			
β_5			

Intervaly spolehlivosti pro regresní koeficienty:

```
> confint(model1)
                2.5 %      97.5 %
(Intercept) 36.715381614 294.56281096
body.w       0.093443025  1.95844313
body.H      -0.005657005  0.08643304
waist.C      0.102513283  0.26268666
hip.C       -0.262652856 -0.10067593
antb.C       0.009770142  0.57262713
```

Z výsledku díličích testu vidíme, že proměnná body.H není na hladině 0,05 významná, sestavíme model, který ji neobsahuje.

```
> model2 <- lm(neck.C ~ body.W + waist.C + hip.C + antb.C, data=neck)
> summary(model2)
```

Call:

```
lm(formula = neck.C ~ body.w + waist.C + hip.C + antb.C, data = neck)
```

Residuals:

```
      Min       1Q   Median       3Q      Max
-36.799  -7.585  -0.460    8.523   33.903
```

Coefficients:

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 255.80441    39.59226   6.461 6.96e-09 ***
body.w       1.49670     0.38801   3.857 0.000227 ***
waist.C      0.15765     0.03809   4.139 8.42e-05 ***
hip.C       -0.21493     0.03641  -5.903 7.74e-08 ***
antb.C       0.28727     0.14318   2.006 0.048114 *
---

```

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```
Residual standard error: 12.51 on 82 degrees of freedom
Multiple R-squared: 0.8519, Adjusted R-squared: 0.8447
F-statistic: 118 on 4 and 82 DF, p-value: < 2.2e-16
```

Z podrobných informací o druhém modelu vidíme, že adjustovaný index determinace je o něco nižší, zřejmě tedy i proměnná body.H přispívá k vysvětlení variability obvodu krku.

Nyní použijeme proceduru STEPWISE, abychom získali co nejlepší model. K tomu v R slouží funkce step(). Vstupními argumenty jsou maximální model a směr, kterým chceme model budovat – zpětný nebo dopředný.

Začneme zpětnou krokovou metodou. Metoda v prvním kroku vypouští vždy jen jednu proměnnou a zjišťuje, jestli se vypuštěním konkrétní jedné proměnné model zlepšuje. Poté vypouští tu, která vede k největšímu zlepšení modelu. V dalším kroku pracuje s modelem,

který ji neobsahuje, a u něj opět vypouští jednotlivé proměnné. Opět vybere tu, jejíž vypuštění vede k největšímu zlepšení modelu, a takto pokračuje dále. Pokud metoda během kroku zjistí, že vypuštění jakékoli proměnné nevede k zlepšení, proces končí.

```
model.back <- step(lm(neck.C ~ body.W + body.H + waist.C + hip.C + antb.C,
data=neck),
+ direction='backward')
Start: AIC=443.21
neck.C ~ body.W + body.H + waist.C + hip.C + antb.C
```

	Df	Sum of Sq	RSS	AIC
<none>			12361	443.21
- body.H	1	464.81	12826	444.42
- antb.C	1	646.82	13008	445.64
- body.W	1	731.29	13092	446.21
- hip.C	1	3039.71	15401	460.33
- waist.C	1	3140.65	15502	460.90

`model.back`

```
Call:
lm(formula = neck.C ~ body.W + body.H + waist.C + hip.C + antb.C,
data = neck)
```

Coefficients:

(Intercept)	body.W	body.H	waist.C	hip.C	antb.C
165.63910	1.02594	0.04039	0.18260	-0.18166	0.29120

Informace o výsledném modelu vypíšeme pomocí funkce `summary`:

```
summary(model.back)
```

Call:

```
lm(formula = neck.C ~ body.W + body.H + waist.C + hip.C + antb.C,
data = neck)
```

Residuals:

Min	1Q	Median	3Q	Max
-32.266	-8.030	1.169	8.493	33.577

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	165.63910	64.79600	2.556	0.0124 *
body.W	1.02594	0.46867	2.189	0.0315 *
body.H	0.04039	0.02314	1.745	0.0847 .
waist.C	0.18260	0.04025	4.537	1.96e-05 ***
hip.C	-0.18166	0.04070	-4.463	2.58e-05 ***
antb.C	0.29120	0.14144	2.059	0.0427 *

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 12.35 on 81 degrees of freedom

Multiple R-squared: 0.8573, Adjusted R-squared: 0.8485

F-statistic: 97.33 on 5 and 81 DF, p-value: < 2.2e-16

Po kolika krocích zpětná kroková metoda skončila?

Liší se její výsledek od výsledku metody ENTER?

Budeme pokračovat dopřednou krokovou metodou. Vstupními argumenty jsou: minimální model (tedy model konstanty), dále v argumentu scope maximální model, který připadá v úvahu (zadá se pravá strana modelu, tedy v našem případě scope= body.W + body.H + waist.C + hip.C + antb.C), a upřesnění směru direction='forward'. Metoda v prvním kroku zkouší přidat vždy jen jednu proměnnou a zjišťuje, jestli se přidáním konkrétní jedné proměnné model zlepší. Poté do modelu zahrne tu, která vede k největšímu zlepšení modelu. V dalším kroku pracuje s modelem, který ji obsahuje, a u něj opět zkouší přidávat ostatní proměnné. Opět vybere tu, jejíž vypuštění vede k největšímu zlepšení modelu, a takto pokračuje dále. Pokud metoda během kroku zjistí, že přidání jakékoli další proměnné nevede k zlepšení, proces končí.

```
model.for <- step(lm(neck.C ~ 1, data=neck),
+                 scope= ~ body.W + body.H + waist.C + hip.C + antb.C,
+                 direction='forward')
Start:  AIC=602.6
neck.C ~ 1
```

	Df	Sum of Sq	RSS	AIC
+ antb.C	1	64034	22594	487.68
+ waist.C	1	62506	24123	493.37
+ body.W	1	58753	27875	505.95
+ body.H	1	33549	53080	561.99
+ hip.C	1	13611	73017	589.73
<none>			86628	602.60

Step: AIC=487.68

```
neck.C ~ antb.C
```

	Df	Sum of Sq	RSS	AIC
+ waist.C	1	4317.8	18277	471.23
+ body.H	1	1873.3	20721	482.15
+ body.W	1	1688.5	20906	482.92
<none>			22594	487.68
+ hip.C	1	363.9	22231	488.27

Step: AIC=471.23

```
neck.C ~ antb.C + waist.C
```

	Df	Sum of Sq	RSS	AIC
+ hip.C	1	3123.45	15153	456.92
+ body.H	1	2459.65	15817	460.65
<none>			18277	471.23
+ body.W	1	0.08	18277	473.23

Step: AIC=456.92

```
neck.C ~ antb.C + waist.C + hip.C
```

	Df	Sum of Sq	RSS	AIC
+ body.W	1	2327.3	12826	444.42
+ body.H	1	2060.8	13092	446.21
<none>			15153	456.92

Step: AIC=444.42

```
neck.C ~ antb.C + waist.C + hip.C + body.W
```

	Df	Sum of Sq	RSS	AIC
+ body.H	1	464.81	12361	443.21
<none>			12826	444.42

Step: AIC=443.21

```
neck.C ~ antb.C + waist.C + hip.C + body.W + body.H
```

model.for

call:

```
lm(formula = neck.C ~ antb.C + waist.C + hip.C + body.w + body.H,  
    data = neck)
```

Coefficients:

(Intercept)	antb.C	waist.C	hip.C	body.w	body.H
165.63910	0.29120	0.18260	-0.18166	1.02594	0.04039

Po kolika krocích dopředná kroková metoda skončila?

Které proměnné byly postupně do modelu přidávány?

.....

Liší se její výsledek dopředné metody od výsledku zpětné metody?