

Cvičení 4.: Vícenásobná lineární regrese

Příklad: U 19 vzorků potravinářské pšenice byl zjišťován obsah zinku v zrně (proměnná Y), v kořenech (proměnná X_1), v otrubách (X_2) a ve stonku a listech (X_3). Údaje jsou uvedeny v mg/kg.

| Y | X_1 | X_2 | X_3 |
|-----|-------|-------|-------|
| 175 | 164 | 198 | 162 |
| 169 | 160 | 198 | 159 |
| 175 | 158 | 211 | 164 |
| 181 | 162 | 211 | 162 |
| 539 | 520 | 567 | 523 |
| 526 | 502 | 540 | 491 |
| 344 | 339 | 355 | 334 |
| 475 | 460 | 500 | 446 |
| 820 | 683 | 813 | 695 |
| 841 | 731 | 832 | 714 |
| 828 | 710 | 846 | 697 |
| 775 | 716 | 818 | 709 |
| 622 | 543 | 635 | 563 |
| 661 | 577 | 712 | 580 |
| 579 | 505 | 596 | 531 |
| 936 | 790 | 946 | 814 |
| 903 | 806 | 946 | 834 |
| 927 | 793 | 912 | 824 |
| 889 | 820 | 919 | 807 |

- Normalitu proměnných Y , X_1 , X_2 , X_3 posuďte pomocí Lilieforsova testu s hladinou významnosti 0,05.
- Závislost mezi dvojicemi proměnných (Y, X_1) , (Y, X_2) , (Y, X_3) znázorněte dvourozměrnými tečkovými diagramy.
- Vypočítejte výběrovou korelační matici všech čtyř proměnných a pro $\alpha = 0,05$ otestujte významnost jednotlivých korelačních koeficientů.
- Vypočítejte koeficienty VIF a ukazatele tolerance pro vysvětlující proměnné X_1 , X_2 , X_3 .
- V první fázi zpracování předpokládejte, že je vhodný regresní model $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \varepsilon$. Vypočítejte index determinace a interpretujte ho. Proveďte celkový F-test. Odhadněte parametry regresního modelu. Proveďte dílčí t-testy pro regresní koeficienty. Zjistěte odhad rozptylu. (Hladinu významnosti volte $\alpha = 0,05$.)
- Posuďte pomocí beta koeficientů vliv jednotlivých nezávisle proměnných veličin na regresní model.
- Z regresního modelu odstraňte ty proměnné, jejichž regresní koeficienty se neprokázaly významné pro $\alpha = 0,05$. Sestavte nový regresní model a proveďte v něm tytéž úkoly jako v bodě e).
- Normalitu reziduí v tomto novém regresním modelu posuďte Lilieforsovým testem na hladině významnosti $\alpha = 0,05$.
- V novém regresním modelu najděte 95% interval spolehlivosti pro teoretickou regresní funkci a 95% predikční interval.
- Proveďte regresi metodou STEPWISE, a to jak Forward, tak Backward.

Řešení: Načteme datový soubor zinek.sta.

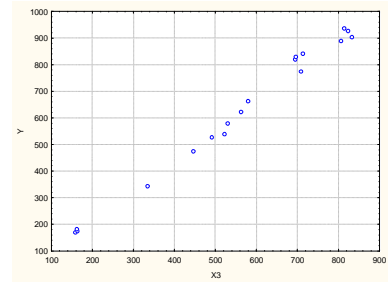
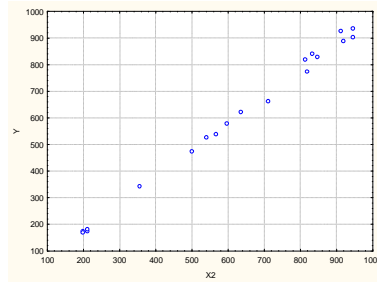
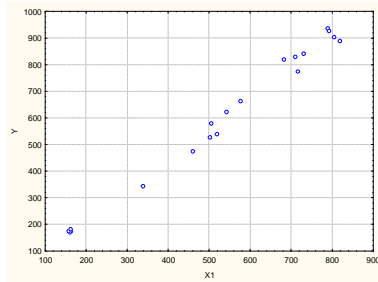
ad a) Výsledky Lilieforsova testu normality

| proměnná | testová statistika | p-hodnota |
|----------------|--------------------|-----------|
| Y | 0,15792 | > 0,2 |
| X ₁ | 0,15613 | > 0,2 |
| X ₂ | 0,18177 | < 0,1 |
| X ₃ | 0,16420 | < 0,2 |

Na hladině významnosti 0,05 nelze ani v jednom případě zamítnout hypotézu o normalitě.

ad b)

Dvourozměrné tečkové diagramy dvojic (Y,X₁), (Y,X₂), (Y,X₃) svědčí o existenci dosti silné přímé lineární závislosti.



ad c) Výběrová korelační matice proměnných Y, X₁, X₂, X₃ spolu s odpovídajícími p-hodnotami:

| Proměnná | Y | X1 | X2 | X3 |
|----------|--------|--------|--------|--------|
| Y | 1,0000 | ,9947 | ,9981 | ,9959 |
| | p= --- | p=,000 | p=0,00 | p=0,00 |
| X1 | ,9947 | 1,0000 | ,9954 | ,9980 |
| | p=,000 | p= --- | p=,000 | p=0,00 |
| X2 | ,9981 | ,9954 | 1,0000 | ,9962 |
| | p=0,00 | p=,000 | p= --- | p=0,00 |
| X3 | ,9959 | ,9980 | ,9962 | 1,0000 |
| | p=0,00 | p=0,00 | p=0,00 | p= --- |

Na hladině významnosti 0,05 zamítáme hypotézu o nevýznamnosti jednotlivých korelačních koeficientů.

ad d) Výpočet koeficientů VIF a ukazatelů tolerance:

Statistiky - Pokročilé lineární/nelineární modely – Obecné regresní modely – OK – Proměnné – Závislá Y, Spojité nezávisle proměnné X1, X2, X3 – OK – Matice – Parciální korelace.

| Statistiky kolineace za daných podmínek (zinek.sta) | | | | | | | | |
|---|----------|---------------------|----------------|-------------|---------------|---------------|-----------|----------|
| Sigma-omezená parametrizace | | | | | | | | |
| Efekt | Toler. | Rozptyl Infl fak | R ² | Y Beta v | Y Parciál. | Y Semipar. | Y t | Y p |
| X1 | 0,003802 | 262,9861 | 0,996198 | -0,037425 | -0,038960 | -0,002308 | -0,151006 | 0,881983 |
| X2 | 0,007214 | 138,6290 | 0,992786 | 0,793836 | 0,751501 | 0,067422 | 4,411716 | 0,000505 |
| X3 | 0,003120 | 320,5035 | 0,996880 | 0,242409 | 0,223005 | 0,013540 | 0,886006 | 0,389598 |

O existenci multikolinearity svědčí extrémně vysoké koeficienty VIF a velmi malé ukazatele tolerance.

ad e) Výsledky pro regresní model $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \varepsilon$ získáme takto:
 Statistiky – Vícenásobná regrese – Proměnné – Závislá proměnná Y, seznam nezáv. proměnných X1, X2, X3 – OK – OK.

| Výsledky regrese se závislou proměnnou : Y (zinek.sta) R= ,99824679 R2= ,99649665 Upravené R2= ,99579598 F(3,15)=1422,2 p<,00000 Směrod. chyba odhadu : 18,094 | | | | | | |
|--|-----------|---------------|----------|------------|----------|----------|
| N=19 | Beta | Sm.chyba beta | B | Sm.chyba B | t(15) | Úroveň p |
| Abs.člen | | | -28,7607 | 10,60478 | -2,71205 | 0,016066 |
| X1 | -0,037425 | 0,247835 | -0,0439 | 0,29089 | -0,15101 | 0,881983 |
| X2 | 0,793836 | 0,179938 | 0,8079 | 0,18312 | 4,41172 | 0,000505 |
| X3 | 0,242409 | 0,273598 | 0,2802 | 0,31623 | 0,88601 | 0,389598 |

Adjustovaný index determinace je 0,9958, tedy zvolený regresní model s proměnnými X₁, X₂, X₃ vysvětluje variabilitu proměnné Y z 99,58 %. Testová statistika pro celkový F-test nabývá hodnoty 1422,2, odpovídající p-hodnota je velmi blízká 0, tedy model jako celek je významný na hladině 0,05.

Odhad rozptylu získáme z tabulky analýzy rozptylu, kterou dostaneme pomocí cesty
 Výsledky – vícenásobná regrese – Detailní výsledky – ANOVA.

| Efekt | Součet čtverců | sv | Průměr čtverců | F | Úroveň p |
|---------|----------------|----|----------------|----------|----------|
| Regres. | 1396846 | 3 | 465615,2 | 1422,205 | 0,000000 |
| Rezid. | 4911 | 15 | 327,4 | | |
| Celk. | 1401757 | | | | |

Vidíme, že $s^2 = 327,4$

Odhadnutá regresní funkce má tvar: $\hat{Y} = -28,7607 - 0,0439x_1 + 0,8079x_2 + 0,2802x_3$.

Dílčí t-testy pro jednotlivé regresní koeficienty:

testová statistika pro test hypotézy $H_0: \beta_0 = 0$ je -2,71205, p-hodnota je 0,016066, tedy H_0 zamítáme na hladině významnosti 0,05;

testová statistika pro test hypotézy $H_0: \beta_1 = 0$ je -0,15101, p-hodnota je 0,881983, tedy H_0 nezamítáme na hladině významnosti 0,05;

testová statistika pro test hypotézy $H_0: \beta_2 = 0$ je 4,41172, p-hodnota je 0,000505, tedy H_0 zamítáme na hladině významnosti 0,05;

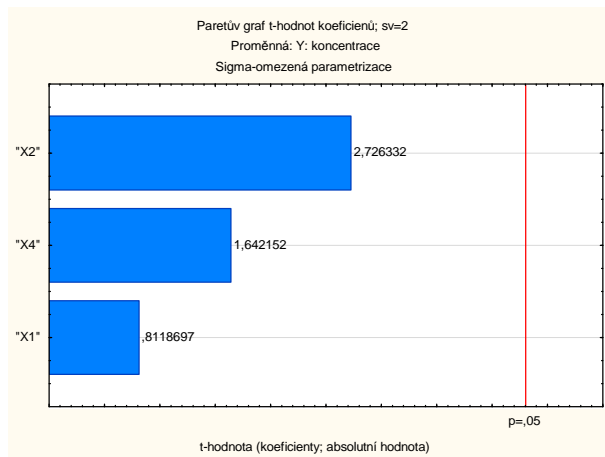
testová statistika pro test hypotézy $H_0: \beta_3 = 0$ je 0,88601, p-hodnota je 0,389598, tedy H_0 nezamítáme na hladině významnosti 0,05.

ad f) Interpretace beta koeficientů:

beta1 = -0,037425, beta2 = 0,793836, beta3 = 0,242409. V absolutní hodnotě je největší beta2, tedy obsah zinku v otrubách má největší vliv na obsah zinku v znu.

Znázornění beta koeficientů pomocí Paretova diagramu:

Statistiky - Pokročilé lineární/nelineární modely – Obecné regresní modely – OK – Proměnné – Závislá Y, Spojité nezávisle proměnné X1, X2, X3 – OK – Paretův graf.



ad g) Protože dílčí t-testy prokázaly, že na hladině 0,05 nejsou proměnné X_1 a X_3 významné, sestavíme nový regresní model $Y = \beta_0 + \beta_2 X_2 + \varepsilon$.

| Výsledky regrese se závislou proměnnou : Y (zinek.sta) | | | | | | |
|--|----------|---------------|----------|------------|----------|----------|
| R= ,99807615 R2= ,99615600 Upravené R2= ,99592988 | | | | | | |
| F(1,17)=4405,5 p<0,0000 Směrod. chyba odhadu : 17,803 | | | | | | |
| N=19 | Beta | Sm.chyba beta | B | Sm.chyba B | t(17) | Úroveň p |
| Abs. člen | | | -30,2507 | 10,31117 | -2,93378 | 0,009274 |
| X2 | 0,998076 | 0,015037 | 1,0157 | 0,01530 | 66,37372 | 0,000000 |

Adjustovaný index determinace je 0,9959, tedy zvolený regresní model s proměnnou X_2 vysvětluje variabilitu proměnné Y z 99,59 %. Testová statistika pro celkový F-test nabývá hodnoty 4405,5, odpovídající p-hodnota je velmi blízká 0, tedy model jako celek je významný na hladině 0,05.

Vidíme, že $\hat{Y} = -30,2507 + 1,0157x_2$.

Dílčí t-testy pro jednotlivé regresní koeficienty:

testová statistika pro test hypotézy $H_0: \beta_0 = 0$ je -2,93378, p-hodnota je 0,009274, tedy H_0 zamítáme na hladině významnosti 0,05;

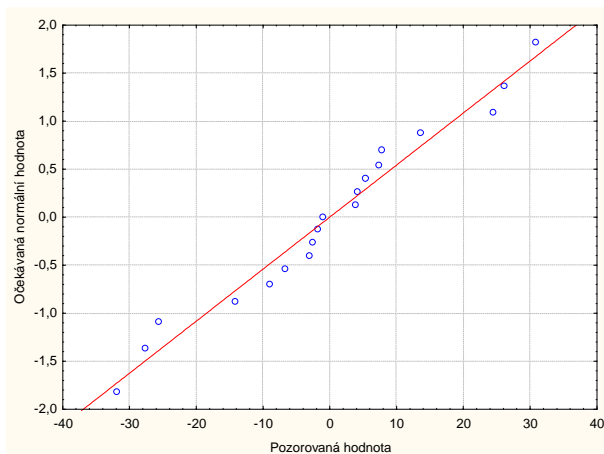
testová statistika pro test hypotézy $H_0: \beta_2 = 0$ je 66,37372, p-hodnota je 0,000000, tedy H_0 zamítáme na hladině významnosti 0,05.

ad h) Ověření normality reziduí

Abychom mohli analyzovat rezidua, musíme je uložit. Ve výstupní tabulce zvolíme Rezidua/předpoklady/předpovědi – Reziduální analýza – Uložit – Uložit rezidua& předpovědi - OK.

Testová statistika pro Lilieforsův test nabývá hodnoty 0,1163, odpovídající p-hodnota je větší než 0,20, tedy hypotézu o normalitě reziduí nezamítáme na hladině významnosti 0,05.

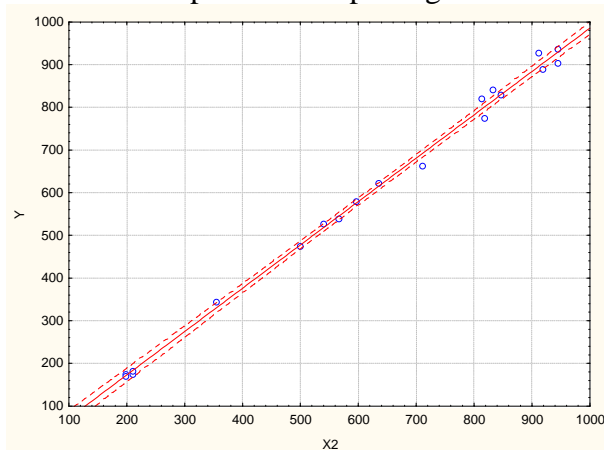
Pro úplnost ještě posoudíme vzhled N-P plotu:



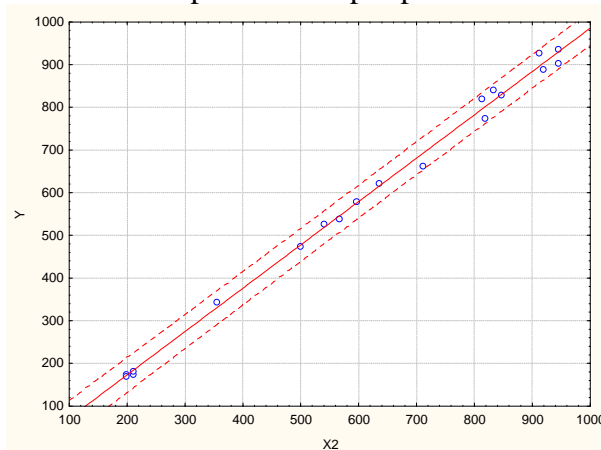
N-P plot svědčí o tom, že rozložení reziduí se příliš neliší od normálního rozložení.

ad i) Intervaly spolehlivosti pro regresní funkci a pro predikci získáme pomocí dvourozměrných tečkových diagramů, kde v Detailech vybereme lineární proložení a zvolíme regresní pásy.

95% interval spolehlivosti pro regresní funkci



95% interval spolehlivosti pro predikci



ad j) Nejprve aplikujeme metodu Forward:

Statistiky – Vícerozměrná regrese – Proměnné – Závisle proměnná Y, Nezávisle proměnné X1, X2, X3 – OK – Detailní nastavení – zaškrtneme Další možnosti – OK – Metoda – zvolíme Kroková dopředná – na záložce Metoda zvolíme Zobrazit výsledky Po každém kroku – OK (V kroku 0 nejsou v regresní rovnici žádné proměnné.) Klikneme na Další – Výpočet: Výsledky regrese.

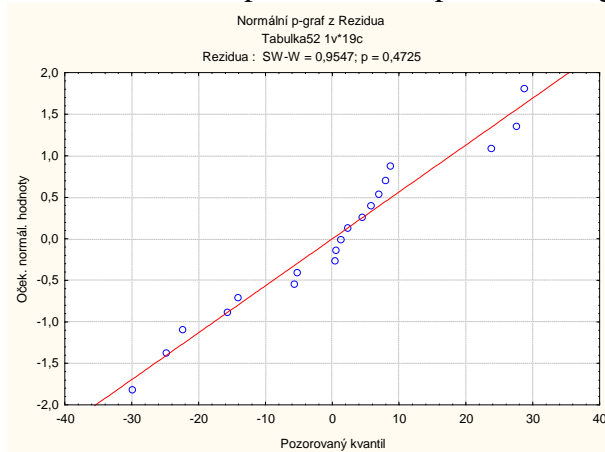
| | | | | | | |
|--|----------|---------------|----------|--------------|----------|----------|
| Výsledky regrese se závislou proměnnou : Y (zinek.sta) | | | | | | |
| R= ,99807615 R2= ,99615600 Upravené R2= ,99592988 | | | | | | |
| F(1,17)=4405,5 p<0,0000 Směrod. chyba odhadu : 17,803 | | | | | | |
| N=19 | b* | Sm.chyba z b* | b | Sm.chyba z b | t(17) | p-hodn. |
| Abs.člen | | | -30,2507 | 10,31117 | -2,93378 | 0,009274 |
| X2 | 0,998076 | 0,015037 | 1,0157 | 0,01530 | 66,37372 | 0,000000 |

V prvním kroku byla vybrána proměnná X2. Opět klikneme na Další a dostaneme výsledky kroku 2, který je již konečný:

| Výsledky regrese se závislou proměnnou : Y (zinek.sta) R= ,99824412 R2= ,99649132 Upravené R2= ,99605274 F(2,16)=2272,1 p<,0000 Směrod. chyba odhadu : 17,533 | | | | | | |
|---|----------|------------------|----------|-----------------|----------|----------|
| N=19 | b* | Sm.chyba z b* | b | Sm.chyba z b | t(16) | p-hodn. |
| Abs.člen | | | -28,9426 | 10,20929 | -2,83493 | 0,011948 |
| X2 | 0,788109 | 0,170440 | 0,8020 | 0,17345 | 4,62396 | 0,000282 |
| X3 | 0,210764 | 0,170440 | 0,2436 | 0,19700 | 1,23659 | 0,234086 |

Empirická regresní funkce má tvar $\hat{Y} = -28,9426 + 0,802x_2 + 0,2436x_3$.
Model jako celek je významný na hladině 0,05, avšak nezávisle proměnná X_3 významná není. Přispívá však k vysvětlení variability hodnot závisle proměnné veličiny Y. Adjustovaný index determinace je 0,9961. V modelu s nezávisle proměnnou X_2 byl 0,9959 a v modelu se všemi třemi nezávisle proměnnými byl 0,9958.

Normalitu reziduí prozkoumáme pomocí N-P grafu a S-W testu:



Rezidua neporušují předpoklad normality.

Nyní provedeme metodu Backward:

Statistiky – Vícerozměrná regrese – Proměnné – Závisle proměnná Y, Nezávisle proměnné X_1, X_2, X_3 – OK – Detailní nastavení – zaškrtneme Další možnosti – OK – Metoda – zvolíme Kroková zpětná – na záložce Metoda zvolíme Zobrazit výsledky Po každém kroku – OK – Výpočet: Výsledky regrese.

| Výsledky regrese se závislou proměnnou : Y (zinek.sta) R= ,99824679 R2= ,99649665 Upravené R2= ,99579598 F(3,15)=1422,2 p<,00000 Směrod. chyba odhadu : 18,094 | | | | | | |
|--|-----------|------------------|----------|-----------------|----------|----------|
| N=19 | b* | Sm.chyba z b* | b | Sm.chyba z b | t(15) | p-hodn. |
| Abs.člen | | | -28,7607 | 10,60478 | -2,71205 | 0,016066 |
| X1 | -0,037425 | 0,247835 | -0,0439 | 0,29089 | -0,15101 | 0,881983 |
| X2 | 0,793836 | 0,179938 | 0,8079 | 0,18312 | 4,41172 | 0,000505 |
| X3 | 0,242409 | 0,273598 | 0,2802 | 0,31623 | 0,88601 | 0,389598 |

V prvním kroku byly zařazeny všechny proměnné.

Klikneme na Další – Výpočet: Výsledky regrese.

| Výsledky regrese se závislou proměnnou : Y (zinek.sta) R= ,99824412 R2= ,99649132 Upravené R2= ,99605274 F(2,16)=2272,1 p<0,0000 Směrod. chyba odhadu : 17,533 | | | | | | |
|--|----------|------------------|----------|-----------------|----------|----------|
| N=19 | b* | Sm.chyba z b* | b | Sm.chyba z b | t(16) | p-hodn. |
| Abs.člen | | | -28,9426 | 10,20929 | -2,83493 | 0,011948 |
| X2 | 0,788109 | 0,170440 | 0,8020 | 0,17345 | 4,62396 | 0,000282 |
| X3 | 0,210764 | 0,170440 | 0,2436 | 0,19700 | 1,23659 | 0,234086 |

V tomto kroku byla vyloučena proměnná X1.

Opět klikneme na Další – Výpočet: Výsledky regrese a dostaneme konečnou tabulku:

| Výsledky regrese se závislou proměnnou : Y (zinek.sta) R= ,99807615 R2= ,99615600 Upravené R2= ,99592988 F(1,17)=4405,5 p<0,0000 Směrod. chyba odhadu : 17,803 | | | | | | |
|--|----------|------------------|----------|-----------------|----------|----------|
| N=19 | b* | Sm.chyba z b* | b | Sm.chyba z b | t(17) | p-hodn. |
| Abs.člen | | | -30,2507 | 10,31117 | -2,93378 | 0,009274 |
| X2 | 0,998076 | 0,015037 | 1,0157 | 0,01530 | 66,37372 | 0,000000 |

Vidíme, že metoda STEPWISE, Backward poskytla stejné výsledky jako metoda ENTER.

Příklad k samostatnému řešení:

V datovém souboru ozon.sta jsou uloženy tyto údaje:

Y ... obsah ozónu v ovzduší [promile],

X₁ ... intenzita slunečního záření v oblasti 400 - 700 nm, [W/m²]

X₂ ... průměrná rychlost větru [km/h]

X₃ ... maximální denní teplota [°C]

Vypočtete koeficienty korelace proměnné Y se všemi nezávisle proměnnými a dále koeficienty korelace všech dvojic nezávisle proměnných.

Výsledek:

| Korelace (ozon.sta) Označ. korelace jsou významné na hlad. p < ,05000 N=23 (Celé případy vynechány u ChD) | | | |
|---|----------|-----------|----------|
| Proměnná | X1 | X2 | X3 |
| Y | 0,668009 | -0,421457 | 0,674464 |

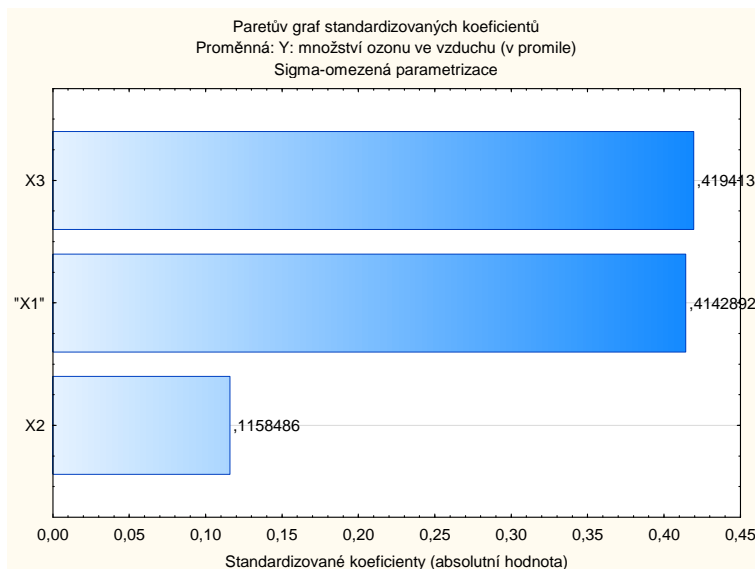
| Korelace (ozon.sta) Označ. korelace jsou významné na hlad. p < ,05000 N=23 (Celé případy vynechány u ChD) | | | |
|---|-----------|-----------|-----------|
| Proměnná | X1 | X2 | X3 |
| X1 | 1,000000 | -0,349474 | 0,508409 |
| X2 | -0,349474 | 1,000000 | -0,383451 |
| X3 | 0,508409 | -0,383451 | 1,000000 |

Pomocí koeficientů VIF proveďte, zda v modelu, který vysvětluje proměnnou Y pomocí proměnných X1 až X3, mezi nezávisle proměnnými veličinami existuje multikolinearita.

| Statistiky kolineace za daných podmínek (ozon.sta) | | | | | | | | |
|--|-----------|---------------------|----------------|-------------|---------------|---------------|------------|-----------|
| Sigma-omezená parametrizace | | | | | | | | |
| Efekt | Toler. | Rozptyl Infl fak | R ² | Y Beta v | Y Parciál. | Y Semipar. | Y t | Y p |
| "X1" | 0,7135268 | 1,4014891 | 0,2864732 | 0,4142892 | 0,4881146 | 0,3499522 | 2,4377773 | 0,0247744 |
| X2 | 0,8207641 | 1,2183769 | 0,1792359 | -0,1158486 | -0,1654184 | -0,1049543 | -0,7311144 | 0,4736301 |
| X3 | 0,6932861 | 1,4424060 | 0,3067139 | 0,4194138 | 0,4873359 | 0,3492199 | 2,4326755 | 0,0250400 |

V uvedeném regresním modelu posuďte pomocí beta koeficientů vliv jednotlivých nezávisle proměnných na Y. Použijte také Paretův diagram.

| Výsledky regrese se závislou proměnnou : Y (ozon.sta) | |
|---|-----------|
| R= ,78003445 R ² = ,60845374 Upravené R ² = ,54663064 | |
| F(3,19)=9,8419 p<,00039 Směrod. chyba odhadu : 6,8481 | |
| N=23 | b* |
| Abs.člen | |
| X1 | 0,414289 |
| X2 | -0,115849 |
| X3 | 0,419414 |



Nyní pro výstavbu modelu použijte dopřednou i zpětnou krokovou metodu a jejich výsledky porovnejte.

Dopředná metoda:

| Výsledky regrese se závislou proměnnou : Y (ozon.sta) | | | | | | |
|---|----------|------------------|----------|-----------------|----------|----------|
| R= ,77294136 R ² = ,59743834 Upravené R ² = ,55718217 | | | | | | |
| F(2,20)=14,841 p<,00011 Směrod. chyba odhadu : 6,7679 | | | | | | |
| N=23 | b* | Sm.chyba z b* | b | Sm.chyba z b | t(20) | p-hodn. |
| Abs.člen | | | -18,3306 | 9,984323 | -1,83594 | 0,081280 |
| X3 | 0,451562 | 0,164755 | 1,6168 | 0,589899 | 2,74080 | 0,012599 |
| X1 | 0,438431 | 0,164755 | 0,0512 | 0,019228 | 2,66110 | 0,014999 |

Zpětná metoda:

| Výsledky regrese se závislou proměnnou : Y (ozon.sta) R= ,67446432 R2= ,45490212 Upravené R2= ,42894508 F(1,21)=17,525 p<,00042 Směrod. chyba odhadu : 7,6856 | | | | | | |
|---|----------|------------------|----------|-----------------|----------|----------|
| N=23 | b* | Sm.chyba z b* | b | Sm.chyba z b | t(21) | p-hodn. |
| Abs.člen | | | -25,6816 | 10,89563 | -2,35706 | 0,028204 |
| X3 | 0,674464 | 0,161112 | 2,4149 | 0,57685 | 4,18631 | 0,000416 |

Mezi těmito dvěma modely rozhodněte na základě reziduální analýzy a adjustovaného indexu determinace.