

Osnova přednášky Vícenásobná lineární regrese

1. Popis modelu

2. Specifika modelu vícenásobné lineární regrese

2.1. Kroky před provedením regresní analýzy

2.2. Sedm hlavních předpokladů modelu

2.3. Ověřování předpokladů modelu

2.4. Posouzení vlivu nezávisle proměnných veličin v modelu

2.5. Pravidla pro stanovení počtu pozorování v závislosti na počtu prediktorů

3. Dvě hlavní metody při provádění vícenásobné lineární regrese

3.1. Metoda ENTER

3.2. Metoda STEPWISE

3.3. Postup při budování modelu vícenásobné lineární regrese

4. Příklad

1. Popis modelu vícenásobné lineární regrese

Budeme zkoumat lineární závislost veličiny Y na p nezávisle proměnných veličinách (regresorech) X_1, \dots, X_p .

Omezíme se pouze na model tvaru

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \varepsilon_i, \quad i = 1, \dots, n.$$

Interpretace parametrů:

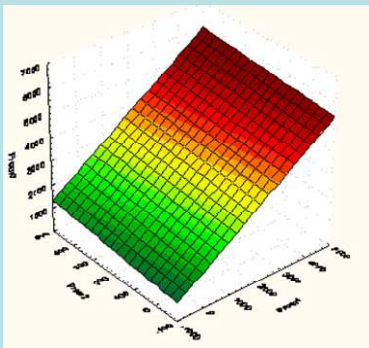
β_0 ... teoretická hodnota závisle proměnné veličiny při nulových hodnotách všech nezávisle proměnných veličin,

β_j ... přírůstek teoretické hodnoty závisle proměnné veličiny odpovídající jednotkové změně j -té nezávisle proměnné veličiny při konstantní úrovni ostatních nezávisle proměnných, $j = 1, \dots, p$.

Parametry β_1, \dots, β_p se nazývají **parciální regresní koeficienty**.

Geometricky tento model představuje regresní nadrovinu.

Ilustrace pro dva regresory:



Model $Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \varepsilon_i$, $i = 1, \dots, n$ lze formálně ztotožnit s lineárním regresním modelem z přednášky „Jednoduchá lineární regrese“:

$$Y_i = \beta_0 + \beta_1 f_1(x_i) + \dots + \beta_p f_p(x_i) + \varepsilon_i, \quad i = 1, \dots, n,$$

kde položíme $f_1(x_i) = x_{i1}$, ..., $f_p(x_i) = x_{ip}$, $i = 1, \dots, n$.

Dostáváme tedy maticový tvar $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, kde regresní matice

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix}, \text{ přičemž } h(\mathbf{X}) = p+1 < n \text{ a } \boldsymbol{\varepsilon} \sim N_n(0, \sigma^2 \mathbf{I}).$$

Všechny výsledky uvedené v přednášce „Jednoduchá lineární regrese“ zůstávají v platnosti.

Příklady vícenásobné regrese

Lékaře zajímá, jak krevní tlak Y závisí na věku pacienta X_1 , na jeho BMI X_2 a na množství vypitého alkoholu X_3 .

Majitele realitní kanceláře zajímá, jak cena bytu Y závisí na velikosti bytu X_1 , na počtu pokojů X_2 , vzdálenosti bytu od centra města X_3 a existenci vlastního parkovacího místa X_4 (1 – ano, 0 – ne).

Pěstitele brambor zajímá, jak výnos Y jisté odrůdy brambor závisí na množství dodaného hnojiva X_1 , na množství srážek X_2 ve vegetačním období a na teplotě půdy X_3 .

Ekonomu zajímá, jak výdaje domácnosti za potraviny a nápoje Y závisí na čistém příjmu domácnosti X_1 a na počtu členů domácnosti X_2 .

Příklad:

Při zkoumání závislosti hodinové výkonnosti dělníka (veličina Y – v kusech) na jeho věku (veličina X_1 – v letech) a době zapracovanosti (veličina X_2 – v letech) byly u 10 náhodně vybraných dělníků zjištěny tyto údaje:

Y	67	65	75	66	77	84	69	60	70	66
X_1	43	40	49	46	41	41	48	34	32	42
X_2	6	8	14	14	8	12	16	1	5	7

Najděte regresní matici a vektor regresních parametrů.

Řešení:

$$\mathbf{X} = \begin{pmatrix} 1 & 43 & 6 \\ 1 & 40 & 8 \\ 1 & 49 & 14 \\ 1 & 46 & 14 \\ 1 & 41 & 8 \\ 1 & 41 & 12 \\ 1 & 48 & 16 \\ 1 & 34 & 1 \\ 1 & 32 & 5 \\ 1 & 42 & 7 \end{pmatrix} \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}$$

2. Specifika modelu vícenásobné lineární regrese

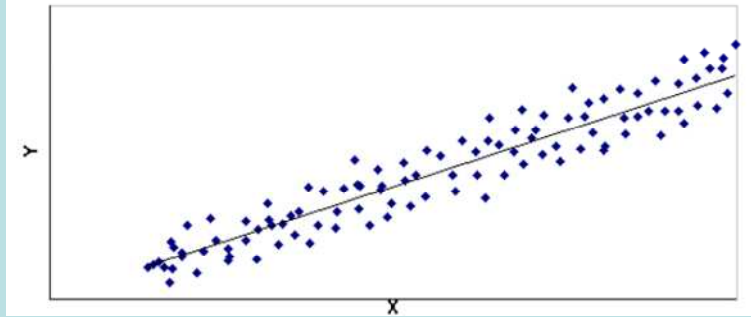
2.1. Kroky před prováděním vícenásobné lineární regrese

- a) Musíme prozkoumat, zda naše data splňují předpoklady pro regresní analýzu.
- b) Pokud je nesplňují, posoudíme, jak vážné je porušení těchto předpokladů.
- c) Je-li porušení předpokladů vážné, musíme s daty provést některé operace, abychom porušení předpokladů odstranili (nebo aspoň zmírnili).

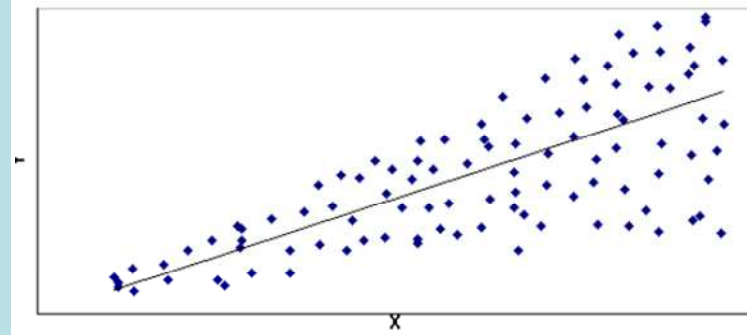
2.2. Sedm hlavních předpokladů regresní analýzy

1. Závisle proměnná Y musí být proměnná aspoň intervalového typu. (Pokud není, musíme použít logistickou regresi.)
2. Nezávisle proměnné X_1, \dots, X_p jsou rovněž aspoň intervalového typu. Mohou to být i proměnné alternativní.
3. Nezávisle proměnné by neměly být mezi sebou příliš vysoce korelovány. Pokud v datech existuje multikolinearita, výsledky regrese jsou nespolehlivé. Vysoká multikolinearita zvyšuje pravděpodobnost, že důležitá nezávisle proměnná bude shledána statisticky nevýznamná a bude vyřazena z modelu.
4. V datech nesmějí být odlehlé či extrémní hodnoty, neboť na ty je regresní analýza citlivá. Odlehlé hodnoty mohou vážně narušit kvalitu odhadů regresních parametrů.
5. Proměnné musejí být v lineárním vztahu. Vícenásobná lineární regrese je založena Pearsonově korelačním koeficientu, takže neexistence linearity způsobuje, že i důležité vztahy mezi proměnnými, pokud nejsou lineární, zůstanou neodhaleny.
6. Proměnné mají normální rozložení. Význam tohoto předpokladu ustupuje do pozadí, máme-li dostatečně velký datový soubor, kde se již uplatňuje působení centrální limitní věty.
7. Proměnné vykazují homoskedasticitu, tedy homogenitu rozptylu. (Opakem homoskedasticity je heteroskedasticita.)

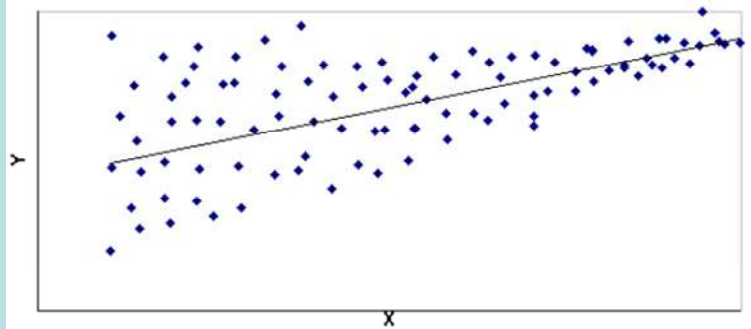
Ukázka homoskedastických dat:



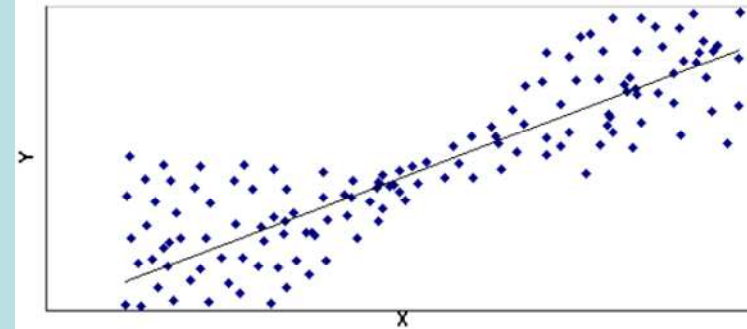
Ukázka dat s rostoucí heteroskedasticitou:



Ukázka dat s klesající heteroskedasticitou:



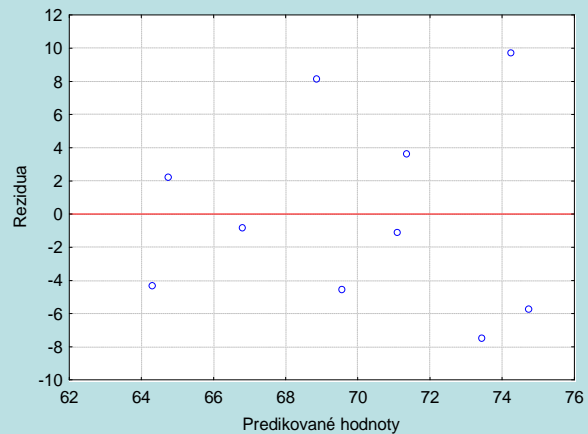
Ukázka dat s proměnlivou heteroskedasticitou:



2.3. Ověřování předpokladů modelu

Ověřování normality:

- jednorozměrná: použijeme např. N-P plot a S-W test či Lilieforsův test.
- vícerozměrná: sestrojíme graf závislosti reziduí na predikovaných hodnotách. Tečky by měly být rovnoměrně rozptýleny po obou stranách vodorovné osy.



Odhalení multikolinearity:

- Vysoké absolutní hodnoty výběrových korelačních koeficientů nezávisle proměnných (orientačně $> 0,75$).
- Velké rozdíly mezi párovými a parciálními korelačními koeficienty.
- Celkový F-test je významný, ale dílčí t-testy nikoliv.

Informace o multikolinearitě poskytuje koeficientu VIF (Variance inflation factor). Má-li koeficient VIF hodnotu 1, pak příslušná nezávisle proměnná není korelovaná s ostatními nezávisle proměnnými, jestliže $1 < VIF < 5$, pak existuje mírná korelace, pro $VIF > 5$ vysoká korelace a pro $VIF > 10$ extrémní multikolinearita. Užitečný je též ukazatel tolerance. Hodnoty pod 0,2 svědčí o multikolinearitě.

Odstranění multikolinearity:

- Je-li multikolinearita způsobena silnou lineární závislostí dvou proměnných, vypustíme jednu z nich z analýzy. Tím se nedopustíme žádné závažné chyby, neboť když máme dvě vysoce vzájemně korelované proměnné, velmi často to znamená, že obě indikují podobný jev. Tím, že jednu z těchto proměnných z regresního modelu vyřadíme, nijak jej neoslabíme.
- Je-li multikolinearita zapříčiněna vzájemnou korelovaností několika proměnných, nabízí se řešení zkombinovat je do jedné nové proměnné. Tu vytvoříme např. s pomocí analýzy hlavních komponent.

Příklad: Při zkoumání závislosti hodinové výkonnosti dělníka (veličina Y – v kusech) na jeho věku (veličina X_1 – v letech) a době zapracovanosti (veličina X_2 – v letech) byly u 10 náhodně vybraných dělníků zjištěny tyto údaje:

Y	67	65	75	66	77	84	69	60	70	66
X_1	43	40	49	46	41	41	48	34	32	42
X_2	6	8	14	14	8	12	16	1	5	7

Posuďte pomocí koeficientu VIF a ukazatel tolerance, zda proměnné věk a doba zapracovanosti mohou způsobit multikolinearitu v modelu $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$.

Řešení:

Efekt	Toler.	Rozptyl Infl fak	R ²	Y Beta v	Y Parciál.	Y Semipar.	Y t	Y p
"X1"	0,282545	3,539258	0,717455	-0,550937	-0,328630	-0,292850	-0,920604	0,387883
"X2"	0,282545	3,539258	0,717455	0,920415	0,502564	0,489246	1,537994	0,167937

Koeficient VIF je 3,54, ukazatel tolerance je větší než 0,2, tedy mezi věkem a dobou zapracovanosti existuje jen mírná korelace.

Odhalení nelinearity vztahů:

Pomocí tečkového diagramu prozkoumáme závislost reziduí na hodnotách závisle proměnné veličiny Y. Pokud tečky vytvoří nelineární obrazec, pak buď jedna z nezávisle proměnných nebo kombinace nezávisle proměnných mají nelineární vztah se závisle proměnnou veličinou Y. Tento graf nám také pomůže odhalit případnou heteroskedasticitu v datech.

Odstranění nelinearity vztahů:

Doporučuje se ty proměnné, u nichž jsme detekovali nelinearitu, transformovat pomocí logaritmické nebo odmocninové transformace. Pokud tento postup nepomůže, musíme použít nelineární regresi.

Odhalení odlehlých hodnot:

Použijeme krabicové grafy nebo pravidlo 3 sigma. Odlehlé hodnoty mají velký vliv na kvalitu odhadu regresních parametrů.

Způsoby řešení problému odlehlých hodnot:

Ověříme, zda při zadávání hodnot dané proměnné nedošlo k překlepu;

proměnnou transformujeme;

upravíme hodnotu odlehlého případu;

odstraníme případy s odlehlou hodnotou;

proměnnou vymažeme.

2.4. Posouzení vlivu jednotlivých nezávisle proměnných v modelu

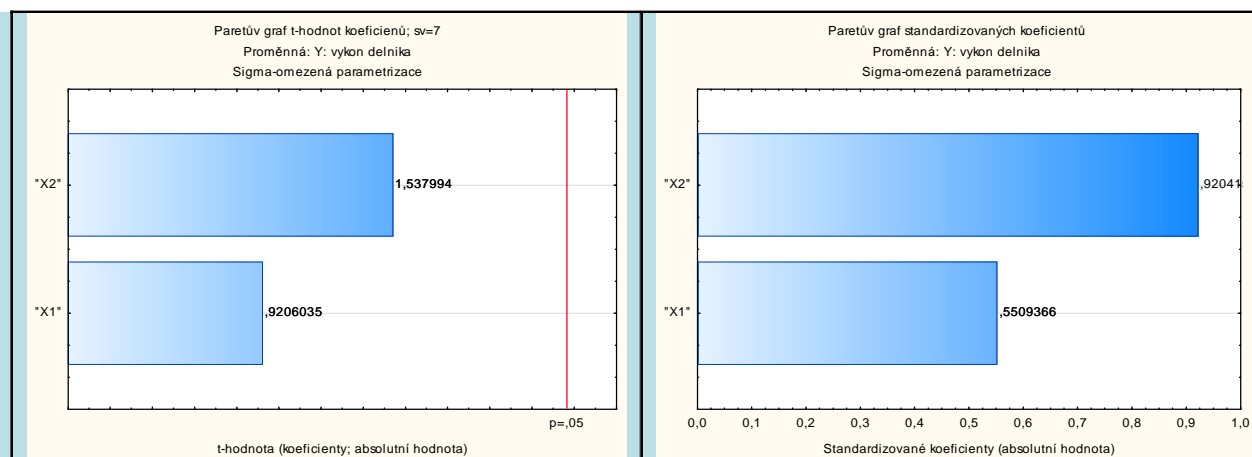
Chceme-li porovnávat vliv, jaký mají proměnné x_1, \dots, x_p v modelu $Y = X\beta + \varepsilon$, můžeme spočítat tzv. standardizované regresní parametry, kterým se také říká B-koefficienty (nebo také beta koefficienty). Zavedeme proto standardizované veličiny

$$Z_i = \frac{Y_i - m_Y}{s_Y}, v_{ij} = \frac{x_{ij} - m_{x_j}}{s_{x_j}}, j = 1, \dots, p, i = 1, \dots, n$$

a vytvoříme regresní model s těmito standardizovanými proměnnými. Odhady regresních parametrů v tomto novém modelu jsou B-koefficienty, které pak vyjadřují intenzitu vlivu jednotlivých nezávisle proměnných veličin na veličinu Y .

V systému STATISTICA jsou B-koefficienty značeny b^* .

Graficky lze absolutní hodnoty standardizovaných regresních parametrů (nebo absolutní hodnoty testových statistik dílčích t-testů) znázornit pomocí Paretových grafů.



2.5. Pravidla pro stanovení počtu pozorování v závislosti na počtu prediktorů

- a) Pokud nás přednostně zajímá variabilita vysvětlená modelem (tj. index determinace), mělo by platit $n \geq 50 + 8p$. Např. pro 4 prediktory bychom měli mít aspoň 82 pozorování.
- b) Zajímají-li nás odhady regresních parametrů (tj. jde nám především o predikci), mělo by platit $n \geq 104 + p$. Např. pro 4 prediktory bychom měli mít aspoň 108 pozorování.
- c) Nejnižší možný poměr proměnná/počet případů je **1:5**. V tom případě ale platí silný požadavek na normalitu – rozložení reziduí by mělo být normální. Při čtyřech prediktorech bychom měli mít aspoň 20 pozorování.

Příklad: Při zkoumání závislosti hodinové výkonnosti dělníka (veličina Y – v kusech) na jeho věku (veličina X_1 – v letech) a době zapracovanosti (veličina X_2 – v letech) byly u 10 náhodně vybraných dělníků zjištěny tyto údaje:

Y	67	65	75	66	77	84	69	60	70	66
X_1	43	40	49	46	41	41	48	34	32	42
X_2	6	8	14	14	8	12	16	1	5	7

Posuďte vliv věku a doby zapracovanosti na výkon dělníka pomocí odhadů standardizovaných regresních parametrů a interpretujte nestandardizované odhady regresních parametrů.

Řešení:

N=10	Výsledky regrese se závislou proměnnou : Y (vykony delniku.sta) R= ,54005243 R2= ,29165662 Upravené R2= ,08927280 F(2,7)=1,4411 p<,29913 Směrod. chyba odhadu : 6,6491					
	b*	Sm.chyba z b*	b	Sm.chyba z b	t(7)	p-hodn.
Abs.člen			86,74217	25,32397	3,425299	0,011056
X1	-0,550937	0,598452	-0,70031	0,76071	-0,920604	0,387883
X2	0,920415	0,598452	1,35062	0,87817	1,537994	0,167937

Odhady standardizovaných regresních parametrů jsou uvedeny ve sloupci b^* . Pro věk má tento parametr hodnotu -0,5509 a pro dobu zapracovanosti 0,9204. V absolutní hodnotě je vyšší parametr pro dobu zapracovanosti, tedy tato proměnná má vyšší vliv na výkon než věk.

Odhad konstanty je 86,74, což znamená, že dělníka s věkem 0 a dobou zapracovanosti 0 by hodinový výkon byl 86,74.

Odhad parametru pro věk je -0,7, tedy při konstantní době zapracovanosti by se při zvýšení věku o rok výkon snížil o 0,7.

Odhad parametru pro dobu zapracovanosti je 1,35, což znamená, že při konstantním věku by se při zvýšení doby zapracovanosti o rok zvýšil výkon o 1,35.

3. Dvě hlavní metody při provádění mnohonásobné lineární regrese

3.1. Metoda ENTER

Tato metoda je standardní metoda, do modelu vstupují všechny nezávisle proměnné najednou.

Metodu ENTER použijeme v případě,

- kdy chceme popsat, jak velký podíl rozptylu závisle proměnné veličiny Y je vysvětlen nezávisle proměnnými veličinami X_1, \dots, X_p (zajímá nás index determinace),
- kdy chceme zjistit, jak velký vliv má každá z nezávisle proměnných na proměnnou závislou při kontrole vlivu působení ostatních proměnných (interpretujeme nestandardizované odhady regresních parametrů),
- kdy nás zajímá, jaká je relativní důležitost každé z nezávisle proměnných (posuzujeme standardizované odhady regresních parametrů).

3.2. Metoda STEPWISE

Metoda STEPWISE (postupná regrese) je metoda nalezení „nejlepšího“ modelu (co nejmenší počet nezávisle proměnných veličin, co nejkvalitnější predikce).

Uživatel nekontroluje pořadí proměnných, jak postupně vstupují do modelu, to provádí samotný program, který pracuje podle jistého algoritmu.

Používá se ve dvou variantách – dopředná (forward) a zpětná (backward).

Při metodě forward se prediktory postupně přidávají, při metodě backward se nejdříve zařadí všechny prediktory a pak se postupně odebírají.

Pořadí vkládání nezávisle proměnných je důležité, neboť může vést k různým odhadům jejich důležitosti v modelu. Proto je při mnohonásobné regresi vždy nutné si dobře rozmyslet, jakou metodu vkládání proměnných zvolíme.

Princip postupné regrese spočívá v tom, že regresní model je budován krok po kroku tak, že v každém kroku zkoumáme všechny prediktory a zjišťujeme, který z nich nejlépe vystihuje variabilitu závisle proměnné veličiny.

Zařazování prediktoru do modelu či jeho vylučování se děje pomocí sekvenčních F-testů.

Sekvenční F-test je založen na statistice F, která je podílem přírůstku regresního součtu čtverců při zařazení daného prediktoru do modelu a reziduálního součtu čtverců.

Jestliže je tato statistika větší než hodnota zvaná „F to enter“ (česky „F na zahrnutí“, ve STATISTICE implicitně 1 pro dopřednou metodu, 11 pro zpětnou), je prediktor zařazen.

Je-li statistika F menší než hodnota zvaná „F to remove“ (česky „F na vyjmutí“, ve STATISTICE implicitně 0 pro dopřednou metodu, 10 pro zpětnou), je již dříve zařazený prediktor z modelu vyloučen.

Po vybrání proměnných do modelu jsou odhadnuty parametry lineární regresní funkce a kvalita regrese je posouzena indexem determinace.

Do modelu se postupně přidávají další proměnné, pokud se zvyšuje podíl vysvětlené variability hodnot veličiny Y.

3.3. Postup při budování modelu vícenásobné lineární regrese

Metoda ENTER

1. Ověříme předpoklady modelu: normalitu, homoskedasticitu, prozkoumáme existenci případné multikolinearity, prověříme linearitu vztahů, detekujeme případná vybočující pozorování.
2. V modelu $Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i, i = 1, \dots, n$ získáme bodové a intervalové odhady regresních parametrů $\beta_0, \beta_1, \dots, \beta_p$, index determinace, odhad rozptylu. Provedeme dílčí t-testy a celkový F-test. Vliv jednotlivých proměnných posoudíme pomocí B-koeficientů.
3. Z modelu vyloučíme ty nezávisle proměnné, pro něž byly dílčí t-testy nevýznamné a odhadneme parametry výsledného modelu.
4. Provedeme reziduální analýzu.

Metoda STEPWISE

1. Ověření předpokladů modelu.
2. Zvolíme dopřednou nebo zpětnou metodu Stepwise, nastavíme hladinu významnosti, hodnoty F na zahrnutí a F na vyjmutí (nebo ponecháme implicitně nastavené hodnoty 0,05, 1, 0).
3. Pro výsledný model provedeme reziduální analýzu.

4. Příklad:

Šest studentů gymnázia absolvovalo čtyři testy, které měří následující veličiny: X_1 - přírodovědné vědomosti, X_2 - literární vědomosti, X_3 - schopnost koncentrace, X_4 - logické myšlení. Testy se hodnotí na škále od 1 do 10 (1 = špatný výsledek, 10 = výborný výsledek).

student	X_1	X_2	X_3	X_4
1	7	9	10	8
2	9	8	8	10
3	4	3	1	2
4	2	3	2	2
5	3	1	2	4
6	1	1	1	4

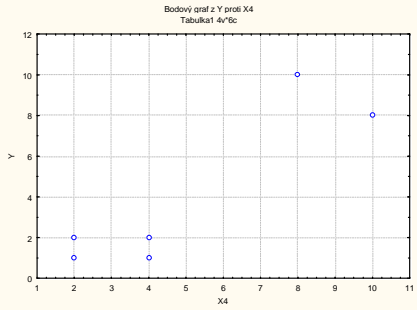
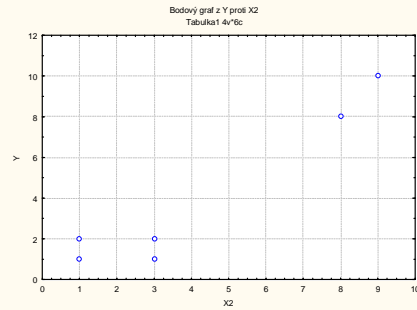
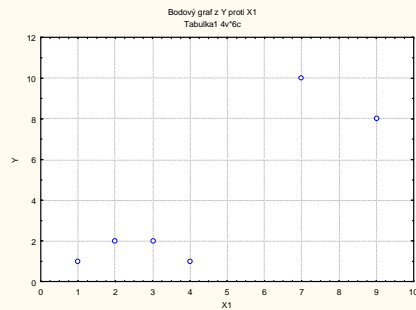
Zajímá nás, kolik bodů můžeme očekávat v testu koncentračních schopností studenta, jestliže známe výsledky testů pro literární schopnosti, přírodovědné schopnosti a logické myšlení.

Řešení:

V tomto problému je proměnná X_3 závislá (označíme ji Y) a ostatní proměnné jsou nezávislé.

Sestavíme regresní model $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_4 x_{i4} + \varepsilon_i$, $i = 1, \dots, 6$.

Nejprve sestrojíme dvourozměrné tečkové diagramy vyjadřující závislost Y na X_1 , X_2 a X_4 .



Dále spočteme výběrové korelační koeficienty r_{Y,X_1} , r_{Y,X_2} , r_{Y,X_4} a výběrové korelační koeficienty r_{X_1,X_2} , r_{X_1,X_4} , r_{X_2,X_4} :

Proměnná	X1	X2	X4
Y	0,872258	0,958133	0,893684

Proměnná	X1	X2	X4
X1	1,000000	0,906937	0,870307
X2	0,906937	1,000000	0,820031
X4	0,870307	0,820031	1,000000

Vidíme, že korelace dvojic (Y, X_1) , (Y, X_2) , (Y, X_4) , (X_1, X_2) , (X_1, X_4) a (X_2, X_4) jsou vysoké.

Metodou nejmenších čtverců získáme odhady regresních parametrů.

Výsledky regrese se závislou proměnnou : y (ctyri_testy.sta) R= ,98240301 R2= ,96511567 Upravené R2= ,91278918 F(3,2)=18,444 p<,05187 Směrod. chyba odhadu : 1,1664						
N=6	b*	Sm.chyba z b*	b	Sm.chyba z b	t(2)	p-hodn.
Abs.člen			-1,08961	0,941927	-1,15679	0,366858
X1	-0,299065	0,368366	-0,38391	0,472872	-0,81187	0,502130
X2	0,864242	0,316998	0,97862	0,358949	2,72633	0,112320
X4	0,445257	0,271142	0,53513	0,325873	1,64215	0,242263

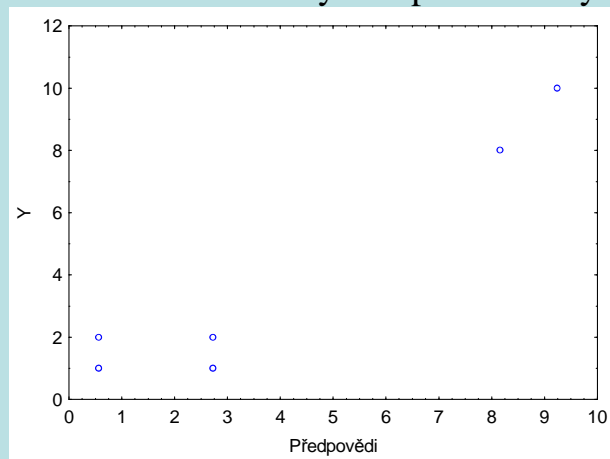
Empirická regresní funkce má tedy tvar $\hat{Y} = -1,09 - 0,38x_1 + 0,98x_2 + 0,54x_4$. Variabilita proměnné Y je z 96,5 % vysvětlená zvoleným regresním modelem. Pro $\alpha = 0,05$ je celkový F-test nevýznamný, všechny dílčí t-testy rovněž. Podíváme-li se na beta koeficienty, vidíme, že největší vliv má proměnná X₂. Sestavíme tedy nový model $Y_i = \beta_0 + \beta_2x_{i2} + \varepsilon_i, i = 1, \dots, 6$. Metodou nejmenších čtverců opět získáme odhady regresních parametrů.

Výsledky regrese se závislou proměnnou : X3 (ctyri testy.sta) R= ,95813306 R2= ,91801897 Upravené R2= ,89752371 F(1,4)=44,792 p<,00259 Směrod. chyba odhadu : 1,2644						
N=6	b*	Sm.chyba z b*	b	Sm.chyba z b	t(4)	p-hodn.
Abs.člen			-0,520548	0,850099	-0,612338	0,573413
X2	0,958133	0,143162	1,084932	0,162108	6,692666	0,002593

Nyní má empirická regresní funkce tvar $\hat{Y} = -0,52 + 1,08x_2$, model jako celek je významný a nezávisle proměnná X₂ rovněž. Pro kontrolu kvality regrese porovnáme zjištěné a predikované hodnoty veličiny Y.

	1 student1	2 student2	3 student3	4 student4	5 student5	6 student6
Skutečnost	10.0	8.0	1.0	2.0	2.0	1.0
Predikce	9.2	8.2	2.7	2.7	0.6	0.6

Vztah mezi naměřenými a predikovanými hodnotami znázorníme pomocí dvourozměrného tečkového diagramu.



Nyní aplikujeme dopřednou metodu postupné regrese:
V nultém kroku je do modelu zařazena pouze konstanta.

Výsledky 1. kroku:

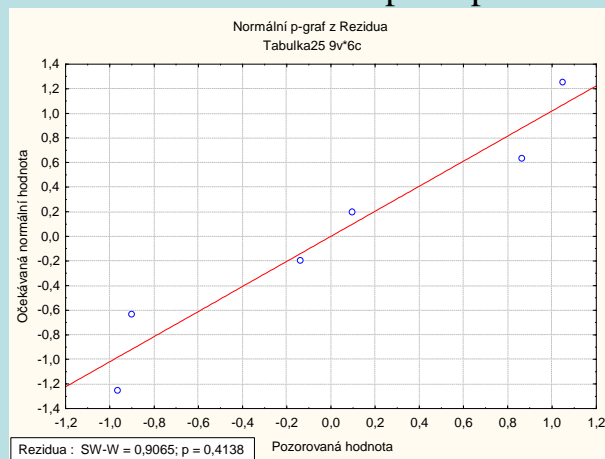
Výsledky regrese se závislou proměnnou : X3 (čtyři testy.sta)						
R= ,95813306 R2= ,91801897 Upravené R2= ,89752371						
F(1,4)=44,792 p<,00259 Směrod. chyba odhadu : 1,2644						
N=6	b*	Sm.chyba z b*	b	Sm.chyba z b	t(4)	p-hodn.
Abs.člen			-0,520548	0,850099	-0,612338	0,573413
X2	0,958133	0,143162	1,084932	0,162108	6,692666	0,002593

V prvním kroku byla vybrána proměnná X₂. Krok 2 je již konečný:

Výsledky regrese se závislou proměnnou : X3 (čtyři testy.sta)						
R= ,97653416 R2= ,95361897 Upravené R2= ,92269829						
F(2,3)=30,841 p<,00999 Směrod. chyba odhadu : 1,0981						
N=6	b*	Sm.chyba z b*	b	Sm.chyba z b	t(3)	p-hodn.
Abs.člen			-1,22615	0,872554	-1,40524	0,254603
X2	0,687789	0,217256	0,77881	0,246007	3,16580	0,050644
X4	0,329675	0,217256	0,39622	0,261109	1,51745	0,226436

Empirická regresní funkce má tvar $\hat{Y} = -1,23 + 0,78x_2 + 0,4x_4$, model jako celek je významný na hladině 0,05, avšak nezávisle proměnné X₂ a X₄ nikoliv. Přispívají však k vysvětlení variability hodnot závisle proměnné veličiny Y. Adjustovaný index determinace je 0,9227. V modelu s nezávisle proměnnou X₂ byl 0,8975 a v modelu se všemi třemi nezávisle proměnnými byl 0,9128.

V tomto výsledném modelu uložíme rezidua a predikované hodnoty:
Pomocí S-W testu a N-P plotu prozkoumáme normalitu reziduí:



Vidíme, že rozložení reziduí je blízké normálnímu rozložení.

Zkusíme ještě zpětnou metodu postupné regrese:

V nultém kroku jsou do modelu zařazeny všechny nezávisle proměnné:

Výsledky regrese se závislou proměnnou : X3 (čtyři testy.sta) R= ,98240301 R2= ,96511567 Upravené R2= ,91278918 F(3,2)=18,444 p<,05187 Směrod. chyba odhadu : 1,1664						
N=6	b*	Sm.chyba z b*	b	Sm.chyba z b	t(2)	p-hodn.
Abs.člen			-1,08961	0,941927	-1,15679	0,366858
X1	-0,299065	0,368366	-0,38391	0,472872	-0,81187	0,502130
X2	0,864242	0,316998	0,97862	0,358949	2,72633	0,112320
X4	0,445257	0,271142	0,53513	0,325873	1,64215	0,242263

V 1. kroku je z modelu vyřazena proměnná X₁:

Výsledky regrese se závislou proměnnou : X3 (čtyři testy.sta) R= ,97653416 R2= ,95361897 Upravené R2= ,92269829 F(2,3)=30,841 p<,00999 Směrod. chyba odhadu : 1,0981						
N=6	b*	Sm.chyba z b*	b	Sm.chyba z b	t(3)	p-hodn.
Abs.člen			-1,22615	0,872554	-1,40524	0,254603
X2	0,687789	0,217256	0,77881	0,246007	3,16580	0,050644
X4	0,329675	0,217256	0,39622	0,261109	1,51745	0,226436

Ve 2. kroku, který je současně poslední, je vyřazena proměnná X₄:

Výsledky regrese se závislou proměnnou : X3 (čtyři testy.sta) R= ,95813306 R2= ,91801897 Upravené R2= ,89752371 F(1,4)=44,792 p<,00259 Směrod. chyba odhadu : 1,2644						
N=6	b*	Sm.chyba z b*	b	Sm.chyba z b	t(4)	p-hodn.
Abs.člen			-0,520548	0,850099	-0,612338	0,573413
X2	0,958133	0,143162	1,084932	0,162108	6,692666	0,002593

Metoda zpětné postupné regrese tedy jako optimální našla model regresní přímky s nezávisle proměnnou X₂.

Upozornění: Pokud bychom změnili hodnoty „F na zahrnutí“ a „F na vyjmutí“, mohli bychom dostat jiné výsledky.