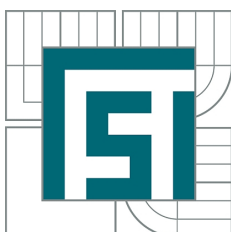


VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ
BRNO UNIVERSITY OF TECHNOLOGY



FAKULTA STROJNÍHO INŽENÝRSTVÍ
ÚSTAV MATEMATIKY
FACULTY OF MECHANICAL ENGINEERING
INSTITUTE OF MATHEMATICS

REGRESNÉ METÓDY ODHADU VYBRANÝCH CHARAKTERISTÍK TAVENÝCH SYROV V ZÁVISLOSTI NA POMERE TAVIACICH SOLÍ

REGRESSION METHODS OF ESTIMATION OF CHOSEN PROPERTIES OF PROCESSED
CHEESE WITH REGARD TO THE RELATIVE AMOUNT OF DIFFERENT TERNARY
MIXTURES OF SODIUM PHOSPHATES

DIPLOMOVÁ PRÁCE
MASTER'S THESIS

AUTOR PRÁCE
AUTHOR

Bc. BRANISLAV PETROVIČ

VEDOUCÍ PRÁCE
SUPERVISOR

doc. RNDr. JAROSLAV MICHÁLEK, CSc.

Vysoké učení technické v Brně, Fakulta strojního inženýrství

Ústav matematiky

Akademický rok: 2012/2013

ZADÁNÍ DIPLOMOVÉ PRÁCE

student(ka): Bc. Branislav Petrovič

který/která studuje v **magisterském navazujícím studijním programu**

obor: **Matematické inženýrství (3901T021)**

Ředitel ústavu Vám v souladu se zákonem č.111/1998 o vysokých školách a se Studijním a zkušebním řádem VUT v Brně určuje následující téma diplomové práce:

Regresní metody odhadu vybraných charakteristik tavených sýrů v závislosti na poměru tavicích solí

v anglickém jazyce:

Regression methods of estimation of chosen properties of processed cheese with regard to the relative amount of different ternary mixtures of sodium phosphates.

Stručná charakteristika problematiky úkolu:

Při výrobě tavených sýrů závisí jejich vlastnosti zejména na použité směsi tavicích solí a na době jejich zralosti. Pro posouzení této statistické vazby byl na Technologické fakultě Univerzity Tomáše Bati ve Zlíně proveden statistický experiment, jehož výsledky umožňují sledované vazby statisticky popsat.

Cíle diplomové práce:

Cílem práce je vybrat, popsat a algoritmizovat vhodné metody regresní analýzy, které jsou pro statistickou analýzu provedeného experimentu vhodné. Zejména se zaměřit na polynomickou regresi a případně metody nelineární. Dále popsat vlastnosti uvedených regresních metod, provést regresní diagnostiku, zaměřit se na výběr regresních parametrů (např. kroková regrese, Cp statistika) a provést příslušné statistické testy. Výsledky analýz demonstrovat graficky.

Seznam odborné literatury:

[1] Brook R.J. and Arnold G.C. Applied regression analysis and experimental design. Marcel Dekker ,INC. 1985

[2] Seber G.A.F. and Wild C.J.: Nonlinear regression. John Wiley and Sons. 1988.

[3] Weiserová E., Doudová L., Galiová L., Žák L., Michálek J., Janiša R., Buňka F.: The effect of combinations of sodium phosphates in binary mixtures on selected texture parameters of processed cheese spreads. International Dairy Journal. Volume 21, Issue 12, December 2011, Pages 979-986 2010

Vedoucí diplomové práce: doc. RNDr. Jaroslav Michálek, CSc.

Termín odevzdání diplomové práce je stanoven časovým plánem akademického roku 2012/2013.

V Brně, dne

L.S.

prof. RNDr. Josef Šlapal, CSc.
Ředitel ústavu

prof. RNDr. Miroslav Doupovec, CSc., dr. h. c.
Děkan fakulty

Abstrakt

Diplomová práca sa zaoberá regresnou analýzou experimentálne nameraných dát tavených syrov. Využitá je polynomiálna regresia a výber regresorov je založený na Krokovej regresii a Mallowsovej štatistike. Následný odhad strednej hodnoty je využitý k nájdeniu najlepšej zmesi taviacich solí, z hľadiska sledovaného parametra taveného syra. Analýza experimentu aj výsledky sú dobre graficky zdokumentované.

Summary

This thesis deals with regression analysis of experimentally measured data of processed cheese. There is a polynomial regression used. The choice of regressors is based on Stepwise Regression and Mallows's Statistics. The estimation of the mean value is used to find the best mixture of the emulsifying salts with regards to the observed characteristic of the processed cheese. Analysis of the experiment and its results are well documented graphically.

Klíčová slova

lineárna regresia, Kroková regresia, Mallowsova C_p štatistika

Keywords

linear regression, Stepwise Regression, Mallows's C_p Statistics

PETROVIČ, B. *Regresné metódy odhadu vybraných charakteristík tavených syrov v závislosti na pomere taviacich solí*. Brno: Vysoké učení technické v Brně, Fakulta strojního inženýrství, 2013. 68 s. Vedoucí diplomové práce doc. RNDr. Jaroslav Michálek, CSc.

Prehlasujem, že som diplomovú prácu „Regresné metódy odhadu vybraných charakteristík tavených syrov v závislosti na pomere taviacich solí“ vypracoval samostatne s použitím odbornej literatúry, uvedenej v zozname, ktorý je súčasťou tejto práce.

Bc. Branislav Petrovič

Ďakujem vedúcemu doc. RNDr. Jaroslavovi Michálkovi, CSc. za cenné rady, venovaný čas a odborný prístup pri tvorbe tejto diplomovej práce.

Bc. Branislav Petrovič

Obsah

| | | |
|----------|--|-----------|
| 1 | Úvod | 3 |
| 1.1 | Charakteristika problematiky | 3 |
| 1.2 | Základné pojmy a označenia | 4 |
| 1.2.1 | Prehľad výsledkov z teórie matic | 7 |
| 2 | Lineárny regresný model úplnej hodnosti | 12 |
| 2.1 | Normálny lineárny regresný model | 19 |
| 2.1.1 | Testy hypotéz a intervaly spoľahlivosti | 20 |
| 3 | Lineárny regresný model neúplnej hodnosti | 24 |
| 3.1 | Odhadnuteľné parametrické funkcie | 25 |
| 3.2 | Normálny lineárny regresný model | 30 |
| 3.2.1 | Testy hypotéz a intervaly spoľahlivosti | 31 |
| 4 | Výber modelu a diagnostika | 33 |
| 4.1 | Výber regresorov | 37 |
| 4.1.1 | Kroková regresia | 37 |
| 4.1.2 | Mallowsova C_p štatistika | 39 |
| 4.2 | Diagnostika | 41 |
| 4.2.1 | Diagnostika pomocou projekčnej matice \mathbf{H} | 41 |
| 4.2.2 | Studentizované reziduá | 42 |
| 5 | Analýza experimentu | 44 |
| 5.1 | Polynóm druhého stupňa | 45 |
| 5.2 | Polynóm tretieho stupňa | 49 |
| 5.3 | Polynóm štvrtého stupňa | 53 |
| 5.3.1 | Modelovanie tvrdosti paraboloidom | 57 |
| 5.4 | Určenie zmesi s najväčšou tvrdosťou | 61 |
| 6 | Záver | 64 |
| | Literatúra | 65 |
| | Príloha | 67 |

OBSAH

1. Úvod

Cieľom tejto práce je pomocou regresnej analýzy vyhodnotiť namerané data v nižšie uvedenom experimente. Práca je rozdelená do šiestich kapitol.

V prvej úvodnej kapitole je odstavce 1.1 venovaný problematike tavených syrov a uskutočnenému experimentu. Text bol vytvorený s použitím článkov [11], [7] a [5]. Druhý odstavce 1.2 je venovaný zavedeniu matematického aparátu. V tomto odstavci bol výklad usporiadaný s použitím literatúr [2], [1] a [9].

Druhá kapitola pojednáva o jednoduchšom prípade, kedy matica modelu nemá plnú stĺpcovú hodnotu. Výklad pre model úplnej hodnoty je usporiadaný s použitím zdrojov [3] a [13].

Tretia kapitola hovorí o lineárnom regresnom modeli, kedy nie je splnená úplná stĺpcová hodnota matice modelu. V tejto kapitole sa občas porovnáva model s úplnou stĺpcovou hodnotou s modelom neúplnou stĺpcovou hodnotou. Výklad je vytvorený s použitím zdrojov [9] a [2].

Štvrtá kapitola pojednáva o dvoch metódach, ktoré sa využívajú pri výbere regresorov v lineárnom modeli. Ďalej hovorí o diagnostike lineárneho regresného modelu. Výklad o výbere modelu vychádza z literatúr [13], [4] a [6]. Výklad o diagnostike je usporiadaný s použitím [10] a [13].

V piatej kapitole sa nachádza aplikačná časť práce. Je tu popísané vyhodnocovanie experimentu a doložené grafickou dokumentáciou. Sú porovnávané obe metódy na výber regresorov. Kapitola uzaviera odhad zmesi taviacich solí, ktorá najlepšie spĺňa zadané kritérium.

1.1. Charakteristika problematiky

Tavené syry

Tavený syr (v angl. processed cheese) je mliečny výrobok, ktorého hlavnou zložkou je tvrdý alebo mäkký syr (alebo ich kombinácia) a príp. tvaroh, sušené mlieko alebo maslo. Výroba taveného syra spočíva v zmiešaní syrov (aj rôznej doby zrenia alebo skladovania) s tavnými soľami (v angl. emulsifying salts) a prípadne inými zložkami a za neustáleho miešania a ohrievania privedené na vhodnú konzistenciu. Pôvodnou myšlienkou výroby syra s kombináciou taviacej soli (pôvodne citrátu sodného), bolo predĺženie trvanlivosti kvôli možnosti lepšieho predaja tohto výrobku. Neskôr začali byť používané tiež taviace soli tvorené fosfátmi alebo ich kombináciami.

Tavený syr bol objavený v roku 1911, vo Švajčiarsku, Walterom Gerberom a Fritzom Stettlerom z Gerber and Co., ktorý rozpúšťali Švajčiarský syr¹ použitím citrátu sodného ako taviacej soli na výrobu jemného a homogénneho výrobku ([7], str. 194).

Hlavná úloha taviacich solí spočíva v procese v ktorom vápnik, ktorý viaže kazeínové bielkoviny (alebo iné hydrolyzované frakcie kazeínu²) vytvárajúc trojdimenzionálnu štruktúru syra, je odlúčený (z kazeínového matrixu) a nahradený sodíkovými iónmi. Prostredníctvom výmeny vápnikových iónov za sodíkové ióny, nerozpustný vápnikový para-kazeín sa zmení na viac rozpustiteľný sodíkový para-kazeín, ktorý môže pôsobiť ako

¹Švajčiarsky syr (v angl. Swiss cheese) je známejší pod názvom Ementál.

²Kazeín je hlavný proteín v cicavčom mlieku.

1.2. ZÁKLADNÉ POJMY A OZNAČENIA

emulgátor ([11], str. 979).

Z predchádzajúceho citovaného odseku teda vyplýva, že taviace soli môžu mať vplyv na vlastnosti výsledného produktu - taveného syra. Toto potvrdzujú aj štúdie, ktoré sú uvedené vo vyššie uvedených článkoch. Výsledné vlastnosti taveného syru sú ovplyvnené ako výrobným procesom (napr. vlastnosťami miešania, dobou ohrievania, teplotou ohrievania a následného ochladenia), tak i skladovaním (napr. teplotou a dobou skladovania). Preto je možné modelovať pomocou regresnej analýzy vybrané parametre tavených syrov (obzvlášť tvrdosti) získaných experimentálne v závislosti na pridaných taviacich soliach.

Experiment

Experiment bol uskutočnený na Fakulte technologickej Univerzity Tomáše Bati v Zlíne. Zmes taviacej soli bola zložená z troch zložiek fosforečnanových taviacich solí:

- Na_2HPO_4 - hydrogenfosforečnan sodný (v angl. disodium hydrogen phosphate)
- $\text{Na}_4\text{P}_2\text{O}_7$ - difosforečnan sodný (v angl. tetrasodium diphosphate)
- $(\text{NaPO}_3)_m$ - polyfosforečnan sodný (v angl. sodium salt of polyphosphate)

Každá zložka zmesi bola dávkovaná po 10 % zo zmesi. Vytvorením všetkých takýchto kombinácií vynikne 66 rôznych zmesí $\left(\binom{12}{10} = 66\right)$. Pre každé meranie boli získané dve hodnoty, preto celkový počet meraní je 132 pre jednu sledovanú vlastnosť, dobu zrenia a dobu skladovania.

Boli sledované tri vlastnosti tavených syrov: tvrdosť, kohezivita a relatívna lepivosť. Pre každú vlastnosť bolo nameraných 132 hodnôt po uplynutí dväťnástich rôznych dób. Najpr sa tavené syry nechali zrieť 2, 4 alebo 8 týždňov a po tejto dobe boli skladované 2, 9 alebo 30 dní. Potom sa uskutočnili jednotlivé merania troch vlastností (každá vlastnosť vždy na osobitnom prvku).

Na výrobu taveného syra pre tento experiment boli použité suroviny: Eidamská tehla, maslo, voda a hore uvedená zmes taviacich solí (3 % z celkovej hmotnosti). Pre polyfosforečnan sodný v zmesi bolo m v rozsahu 15 – 20. Eidamský syr s maslom boli nakrájané na menšie kusky a zohrievané a rozmieľané počas doby 1 min približne na otáčkach 4000 ot/min. Zmes taviacich solí a voda boli pridané pri otáčkach 2000 ot/min. Ohrev trval 10-12 min a počas jednej minuty bola teplota udržiavaná na 90°C . Následne bola tavenina zabalená do malých foriem a ochladzovaná na teplotu 6°C počas 2 hodín. Ďalej sa výrobok nechal zrieť v miestnosti na to určenej pri teplote $9-10^\circ\text{C}$ (po dobu 2, 4 alebo 8 týždňov). Následné skladovanie prebiehalo pri teplote 6°C po dobu 2, 9 alebo 30 dní.

Úlohou tejto práce je štatistický rozbor na tvrdosti taveného syru. Preto bude uvedená definícia tvrdosti a technické prevedenie. Valcová tyč s priemerom 20mm bola pri rýchlosti 2 mm/s vtlačovaná do taveného syra do hĺbky 10 mm. Mierou tvrdosti je tu maximálna sila, ktorá bola potrebná pri prieniku do syra.

Podrobnejšie informácie o uvedenom experimente sú uvedené v článku [5].

1.2. Základné pojmy a označenia

V tejto práci budeme pri výklade matematickej teórie predpokladať základné znalosti z kurzov pravdepodobnosti a štatistiky. V tomto odstavci budú uvedené niektoré matema-

tické pojmy, ktoré budú použité v následnom teoretickom výklade. Ďalej odstavce bude doplnený matematickými vetami z teórie matic a vektorov. Okrem miernych zmien, ktoré budú detailne uvedené, sa autor drží označenia zavedeného Andělom v literatúre [3].

Priestor všetkých elementárnych javov bude označený symbolom Ω . Nech je na priestore Ω nejaká σ -algebra \mathcal{A} podmnožin Ω . Ak P je nejaká pravdepodobnostná miera na σ -algebre \mathcal{A} , potom trojica (Ω, \mathcal{A}, P) bude označovať pravdepodobnostný priestor. Z teórie pravdepodobnosti bude uvedená známa Bonferroniho nerovnosť, ktorá znie nasledovne: Pre náhodné javy A_1, \dots, A_n platí vzťah

$$P(\cup_{i=1}^n A_i) \leq \sum_{i=1}^n P(A_i). \quad (1.1)$$

Zobrazenie $X : \Omega \rightarrow \mathbb{R}$ bude nazývané *náhodnou veličinou* v (Ω, \mathcal{A}, P) ak platí

$$\{\omega \in \Omega \mid X(\omega) \leq a\} \in \mathcal{A} \quad \text{pre } \forall a \in \mathbb{R}.$$

Náhodné veličiny a vektory budú označené veľkými písmenami z konca latinky. Hustota rozdelenia náhodnej veličiny bude označená f a distribučná funkcia bude označená F . Náhodný vektor oproti náhodnej veličine bude označený tučne. Strednú hodnotu náhodnej veličiny X bude označená EX a rozptyl DX . Stredná hodnota náhodného vektora \mathbf{X} bude označená $E\mathbf{X}$ a variančná matica $\text{Var}\mathbf{X}$. Variančná matica je symetrická a pozitívne semidefinitná. Kovariancia dvoch náhodných vektorov alebo veličín \mathbf{X} a \mathbf{Y} bude označená $\text{Cov}(\mathbf{X}, \mathbf{Y})$. Následne budú uvedené vzťahy, ktoré platia pre variančnú a kovariančnú maticu

$$\text{Var}\mathbf{X} = E\mathbf{X}\mathbf{X}' - (E\mathbf{X})(E\mathbf{X})', \quad \text{Cov}(\mathbf{X}, \mathbf{Y}) = E\mathbf{X}\mathbf{Y}' - (E\mathbf{X})(E\mathbf{Y})'.$$

Ďalej budú uvedené základné pravdepodobnostné rozdelenia, ktoré budú použité. Normálne rozdelenie so strednou hodnotou μ a rozptylom σ^2 bude označené $N(\mu, \sigma^2)$. Formálne sem patrí aj singulárne normálne rozdelenie $N(\mu, 0)$. n -rozmerné normálne rozdelenie so strednou hodnotou $\boldsymbol{\mu}$ a variančnou maticou $\sigma^2\mathbf{V}$ bude označené $N_n(\boldsymbol{\mu}, \sigma^2\mathbf{V})$. Ďalej bude uvedené označovanie pre odvodené rozdelenia z normálneho rozdelenia ako Studentovo rozdelenie s n stupňami voľnosti $t(n)$, Chí-kvadrát rozdelenie s n stupňami voľnosti $\chi^2(n)$ a Fisherovo-Snedecorovo rozdelenie s n a k stupňami voľnosti $F(n, k)$. Následne budú zavedené označenia jednotlivých kvantilových funkcií. Keďže všetky tieto rozdelenia majú rastúcu (pre $\chi^2(n)$ a $F(n, k)$ uvažované na intervale $(0, \infty)$) a spojitú distribučnú funkciu F , bude postačovať definícia kvantilovej funkcie ako inverznej funkcie F^{-1} k distribučnej funkcii. Napríklad hodnota kvantilovej funkcie $F^{-1}(\alpha)$ (pre $0 < \alpha < 1$) a prípadne príslušné stupne voľnosti budú označené nasledovne: u_α kvantil normovaného normálneho rozdelenia, $t_\alpha(n)$ kvantil studentovho rozdelenia, $\chi_\alpha^2(n)$ kvantil chí-kvadrát rozdelenia, $F_\alpha(n, k)$ kvantil Fisherovho-Snedecorovho rozdelenia.

Skutočnosť, že náhodná veličina X má normálne rozdelenie $N(\mu, \sigma^2)$ bude jednoducho označená zápisom $X \sim N(\mu, \sigma^2)$. Nech $\mathbf{X} = (X_1, \dots, X_n)'$ je náhodný vektor, \mathbf{V} pozitívne semidefinitná matica a $\boldsymbol{\mu}$ je daný vektor. Ak pre ľubovoľný vektor $\mathbf{c} \in \mathbb{R}^n$ platí

$$\mathbf{c}'\mathbf{X} \sim N(\mathbf{c}'\boldsymbol{\mu}, \sigma^2\mathbf{c}'\mathbf{V}\mathbf{c}), \quad (1.2)$$

potom sa povie, že \mathbf{X} má n -rozmerné normálne rozdelenie $N_n(\boldsymbol{\mu}, \sigma^2\mathbf{V})$ (zápisom $\mathbf{X} \sim N_n(\boldsymbol{\mu}, \sigma^2\mathbf{V})$). V prípade, že matica \mathbf{V} je singulárna, sa pridáva prívlastok singulárne

1.2. ZÁKLADNÉ POJMY A OZNAČENIA

rozdelenie. Pre vektor $\mathbf{a} \in \mathbb{R}^n$, maticu \mathbf{B} typu $n \times m$ a náhodný vektor $\mathbf{X} \sim \mathbf{N}_m(\boldsymbol{\mu}, \sigma^2 \mathbf{V})$ platí

$$\mathbf{a} + \mathbf{B}\mathbf{X} \sim \mathbf{N}_n(\mathbf{a} + \mathbf{B}\boldsymbol{\mu}, \sigma^2 \mathbf{B}\mathbf{V}\mathbf{B}'), \quad (1.3)$$

ako ukázal napríklad Anděl v literatúre ([2], str. 75).

Nech je uvažované $\mathbf{X} \sim \mathbf{N}_n(\boldsymbol{\mu}, \sigma^2 \mathbf{V})$, kde $h(\mathbf{X}) = r \geq 1$. Ďalej nech $\mathbf{B}_{n \times r}$ je matica spĺňajúca $\mathbf{B}\mathbf{B}' = \mathbf{V}$. Potom pre náhodný vektor $\mathbf{Z} \sim \mathbf{N}_r(\mathbf{0}, \mathbf{I})$ platí, že rozdelenie vektora $\boldsymbol{\mu} + \sigma \mathbf{B}\mathbf{Z}$ je rovnaké ako náhodného vektora \mathbf{X} (ako ukázal napríklad Anděl v literatúre ([2], str. 76)).

Následne budú zavedené výberové charakteristiky pre náhodný výber X_1, \dots, X_n rozsahu n . Výberový priemer označený \bar{X} a výberový rozptyl označený S_X^2 budú definované vzťahmi

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad S_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2 \right).$$

Nech je uvažovaný náhodný výber $(X_1, Y_1)', \dots, (X_n, Y_n)'$ rozsahu n z dvojrozmerného rozdelenia. Budú uvažované charakteristiky výberový priemer \bar{X} a výberový rozptyl S_X^2 náhodného výberu X_1, \dots, X_n . Obdobne sú uvažované charakteristiky \bar{Y} a S_Y^2 pre náhodný výber Y_1, \dots, Y_n . Potom je charakteristika výberová kovariancia označená S_{XY} a definovaná vzťahom

$$S_{XY} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = \frac{1}{n-1} \left(\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y} \right).$$

A výberový korelačný koeficient označený r_{XY} je definovaný vzťahom $r_{XY} = S_{XY} / \sqrt{S_X^2 S_Y^2}$, ak $S_X^2 S_Y^2 \neq 0$.

Nech je uvažovaná postupnosť náhodných veličín X_1, X_2, \dots a náhodná veličina X . Nech sú tieto veličiny uvažované na rovnakom pravdepodobnostnom priestore $(\Omega, \mathcal{A}, \mathbf{P})$. Ďalej nech X_n má distribučnú funkciu $F_n(x)$ a X má distribučnú funkciu $F(x)$. Povie sa, že X_n konverguje k X v distribúcii, ak $F_n(x)$ konverguje k $F(x)$ v každom bode x , v ktorom je $F(x)$ spojitá. Rovnako je možné povedať, X_n má asymptoticky rozdelenie náhodnej veličiny X .

Niekedy súčet nezávislých a rovnako rozdelených náhodných veličín môže mať rozdelenie asymptoticky normálne. O tom pojednáva nasledujúca centrálna limitná veta.

Veta 1.1. (Lévyho-Lindebergova) *Nech X_1, X_2, \dots sú nezávislé náhodné veličiny s rovnakým rozdelením, ktoré má strednú hodnotu μ a konečný rozptyl σ^2 . Potom*

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n X_i - \mu \right) \text{ konverguje v distribúcii k rozdeleniu } \mathbf{N}(0, \sigma^2). \quad (1.4)$$

Dôkaz. Vid' Anděl ([2], str. 184).

Delta metóda

Nech $\hat{\theta}_n$ je štatistika, ktorá je závislá na rozsahu výberu n a má asymptoticky normálne rozdelenie $\mathbf{N}(\theta, \sigma^2/n)$. Ako bolo uvedené vyššie, platí

$$\sqrt{n} (\hat{\theta}_n - \theta) \text{ konverguje v distribúcii k } \mathbf{N}(0, \sigma^2).$$

Nech je uvažovaná funkcia g , ktorá je v bode θ najmenej dvakrát diferencovateľná. Otázka je, aké rozdelenie bude mať $g(\hat{\theta}_n)$. Ako ukazuje Agresti ([1], str. 578), platí

$$\sqrt{n} \left(g(\hat{\theta}_n) - g(\theta) \right) \text{ konverguje v distribúcii k } \mathbf{N}(0, \sigma^2 [g'(\theta)]^2).$$

Tento výsledok je nazývaný *delta metódou* pre určenie asymptotického rozdelenia. Väčšinou σ^2 aj $g'(\theta)$ je závislé na θ , asymptotický rozptyl je neznámy. Využije sa teda výberový odhad $\hat{\theta}$ parametra θ . A za predpokladu spojitosti g' i σ v bode θ , je $\sigma(\hat{\theta})g'(\hat{\theta})$ konzistentným odhadom $\sigma(\theta)g'(\theta)$. A toto má významný dôsledok pri testovaní, že platí

$$\sqrt{n} \left(g(\hat{\theta}_n) - g(\theta) \right) / \sigma(\hat{\theta}) \left| g'(\hat{\theta}) \right| \text{ konverguje v distribúcii k } \mathbf{N}(0, 1).$$

Pre veľký výber je možné vytvoriť nasledujúci interval spoľahlivosti pre $g(\theta)$ na hladine významnosti α

$$\left(g(\hat{\theta}) + u_{\alpha/2} \sigma(\hat{\theta}) \left| g'(\hat{\theta}) \right| / \sqrt{n}, g(\hat{\theta}) + u_{1-\alpha/2} \sigma(\hat{\theta}) \left| g'(\hat{\theta}) \right| / \sqrt{n} \right).$$

Ďalej Agresti zovšeobecňuje delta metódu pre funkciu náhodného vektora. Nech vektor štatistík $\hat{\boldsymbol{\theta}}_n = (\hat{\theta}_{n,1}, \dots, \hat{\theta}_{n,m})'$ má asymptoticky m -rozmerné normálne rozdelenie $\mathbf{N}_m(\boldsymbol{\theta}, \boldsymbol{\Sigma}/n)$. Nech ďalej funkcia $g(t_1, \dots, t_m)$ má nenulový diferenciál $\boldsymbol{\phi} = (\phi_1, \dots, \phi_m)'$ v bode $\boldsymbol{\theta}$, kde

$$\phi_i = \left. \frac{\partial g}{\partial t_i} \right|_{\mathbf{t}=\boldsymbol{\theta}} \quad \text{pre } i = 1, 2, \dots, m.$$

Potom

$$\sqrt{n} \left(g(\hat{\boldsymbol{\theta}}_n) - g(\boldsymbol{\theta}) \right) \text{ konverguje v distribúcii k } \mathbf{N}(0, \boldsymbol{\phi}' \boldsymbol{\Sigma} \boldsymbol{\phi}). \quad (1.5)$$

Teda pre veľký výber n je rozdelenie $g(\hat{\boldsymbol{\theta}}_n)$ blízke normálnemu so strednou hodnotou $g(\boldsymbol{\theta})$ a rozptylom $\boldsymbol{\phi}' \boldsymbol{\Sigma} \boldsymbol{\phi} / n$.

1.2.1. Prehľad výsledkov z teórie matic

V práci bude predpokladaná znalosť základných pojmov z lineárnej algebry, prevažne teórie matic a vektorov. Ďalej budú uvedené potrebné výsledky z tejto teórie. Matematická teória v celej práci bude uvažovaná na poli reálnach čísel \mathbb{R} . Matice budú vždy označované veľkými písmenami latinky tučne. Potom reálne vektory budú označené malými písmenami latinky tučne. Matica bude označovaná s príslušným typom matice takto $\mathbf{A}_{m \times n} = (a_{ij})_{i=1, \dots, m, j=1, \dots, n}$, kde m vyjadruje počet riadkov a n počet stĺpcov matice $\mathbf{A}_{m \times n}$. Vektory budú vždy uvažované v stĺpcovom tvare a budú označené $\mathbf{a} = (a_i)_{i=1, \dots, m}$ bez typu, ktorý bude vopred uvedený. Avšak ak bude zo zadania zrejmé o aký typ matice/vektora sa jedná alebo tento typ nebude dôležitý, matica bude pre jednoduchosť označené bez typu \mathbf{A} a pri vektore \mathbf{a} už nebude uvádzaný typ. Pre j -tý stĺpec matice \mathbf{A} bude označený $\mathbf{a}_{\bullet j}$. V prípade štvorcovej matice \mathbf{A} bude v indexe uvedený rád tejto matice takto \mathbf{A}_n . Jednotková matica bude označená \mathbf{I} a nulová matica bude označená \mathbf{O} . Vektor nul bude označený $\mathbf{0}$. Hodnosť matice \mathbf{A} bude označená $h(\mathbf{A})$. Determinant štvorcovej matice \mathbf{A} bude označený $\det(\mathbf{A})$. Transponovaná matica k matici \mathbf{A} bude označená \mathbf{A}' . Inverzná matica k matici \mathbf{A} bude označená \mathbf{A}^{-1} . Stopa matice \mathbf{A} bude označená $\text{Tr}(\mathbf{A})$. Pre nejakú štvorcovú maticu \mathbf{A} a ľubovoľný vektor \mathbf{c} príslušného typu platí užitočný vzťah

$$\mathbf{c}' \mathbf{A} \mathbf{c} = \text{Tr}(\mathbf{A} \mathbf{c} \mathbf{c}'), \quad (1.6)$$

1.2. ZÁKLADNÉ POJMY A OZNAČENIA

ktorý je možné jednoduchým rozpísaním pravej a ľavej strany potvrdiť. Ďalej bude zavedené označenie $\mathcal{M}(\mathbf{A}_{m \times n}) = \mathcal{M}(\mathbf{a}_{\bullet 1}, \dots, \mathbf{a}_{\bullet n})$, ktoré bude vyjadrovať vektorový priestor vytvorený lineárnou kombináciou stĺpcov matice \mathbf{A} (resp. generovaný $\mathbf{a}_{\bullet 1}, \dots, \mathbf{a}_{\bullet n}$). Vektorový priestor bude označovaný \mathcal{V} . Bude uvažovaná euklidovská norma označená $\|\cdot\|$. Matica \mathbf{A} spĺňajúca $\mathbf{A}\mathbf{A}' = \mathbf{A}'\mathbf{A} = \mathbf{I}$, bude nazývaná *ortogonálnou maticou*. Nasledujúca veta bude užitočná pri dokazovaní mnohých matematických viet v tejto práci.

Veta 1.2. *Nech je uvažovaná matica \mathbf{A} s n stĺpcami. Potom platí pre vektorové priestory $\mathcal{M}(\mathbf{A}')$ a $\mathcal{M}(\mathbf{A}'\mathbf{A})$ v n -rozmernom vektorovom priestore \mathcal{V} nasledujúce*

$$\mathcal{M}(\mathbf{A}') = \mathcal{M}(\mathbf{A}'\mathbf{A}). \quad (1.7)$$

Dôkaz. Zrejme priestory $\mathcal{M}(\mathbf{A}')$, $\mathcal{M}(\mathbf{A}'\mathbf{A})$ sú podpriestormi \mathcal{V} . Ak sú oba podpriestory identické v \mathcal{V} , majú identické aj svoje komplementy. A platí to aj opačne. Nech $\mathbf{v} \in \mathcal{V}$ a zároveň platia obe nasledujúce tvrdenia:

1. pre $\forall \mathbf{v} \notin \mathcal{M}(\mathbf{A}')$ platí $\mathbf{v}'\mathbf{A}' = \mathbf{0}$ a úpravou vyplýva platnosť $\mathbf{v}'\mathbf{A}'\mathbf{A} = \mathbf{0}$, a teda záver $\mathbf{v} \notin \mathcal{M}(\mathbf{A}'\mathbf{A})$
2. pre $\forall \mathbf{v} \notin \mathcal{M}(\mathbf{A}'\mathbf{A})$ platí $\mathbf{v}'\mathbf{A}'\mathbf{A} = \mathbf{0}$ a úpravou vyplýva platnosť taktiež $\mathbf{v}'\mathbf{A}'\mathbf{A}\mathbf{v} = 0$, z čoho jasne vyplýva $\mathbf{v}'\mathbf{A}' = \mathbf{0}$ a teda záver $\mathbf{v} \notin \mathcal{M}(\mathbf{A}')$

Potom sú ich komplementy totožné a taktiež oba podpriestory $\mathcal{M}(\mathbf{A}')$, $\mathcal{M}(\mathbf{A}'\mathbf{A})$. \square

Nepriamym dôsledkom predchádzajúcej vety je platnosť vzťahu $h(\mathbf{A}') = h(\mathbf{A}'\mathbf{A})$.

Veta 1.3. (o skeletnom rozklade) *Nech $h(\mathbf{A}_{m \times n}) = r \geq 1$. Potom existujú matice $\mathbf{B}_{m \times r}$ a $\mathbf{C}_{r \times n}$ také, že platí*

$$\mathbf{A} = \mathbf{B}\mathbf{C} \text{ a } h(\mathbf{B}) = h(\mathbf{C}) = r. \quad (1.8)$$

Dôkaz. Dôkaz existencie matic \mathbf{B} a \mathbf{C} bude založený na ich konštrukcii. Z $h(\mathbf{A}) = r$ vyplýva, že matica \mathbf{A} má r nezávislých stĺpcov. Nech matica \mathbf{B} je tvorená z r nezávislých stĺpcov matice \mathbf{A} , v takom poradí ako boli v matici \mathbf{A} . Zrejme $h(\mathbf{B}) = r$. Nech symbol \mathcal{I} značí množinu čísiel, ktoré vyjadrujú poradie stĺpcov v matici \mathbf{A} vybraných do matice \mathbf{B} . Potom matica \mathbf{C} má stĺpce $\mathbf{c}_{\bullet i}$ na pozícii $i \in \mathcal{I}$ tvorené jednou jedničkou a $r - 1$ nulami, tak aby platil súčin $\mathbf{B}\mathbf{c}_{\bullet i} = \mathbf{a}_{\bullet i}$. Stĺpce $\mathbf{c}_{\bullet j}$ na pozíciách $j \in \{1, \dots, n\} - \mathcal{I}$ sú dané koeficientami lineárnej kombinácie stĺpcov matice \mathbf{B} , tak aby platil súčin $\mathbf{B}\mathbf{c}_{\bullet j} = \mathbf{a}_{\bullet j}$. A platí $h(\mathbf{C}) = r$, pretože matica \mathbf{C} obsahuje r nezávislých stĺpcov na všetkých r pozíciách z množiny \mathcal{I} . A keďže bol zavedený súčin $\mathbf{B}\mathbf{c}_{\bullet i} = \mathbf{a}_{\bullet i}$ pre $\forall i \in \{1, \dots, n\}$, platí taktiež súčin $\mathbf{A} = \mathbf{B}\mathbf{C}$. \square

Predchádzajúci dôkaz odpovedal taktiež na jednoznačnosť. Matice $\mathbf{B}_{m \times r}$ a $\mathbf{C}_{r \times n}$ nie sú jednoznačne maticou $\mathbf{A}_{m \times n}$ určené, ale v prípade výberu matice \mathbf{B} resp. \mathbf{C} je jednoznačne daná matica \mathbf{C} resp. \mathbf{B} .

Nech je teraz uvažovaná matica \mathbf{A} štvorcová. Rovnica $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$ pre neznáme číslo λ , je nazývaná *charakteristickou rovnicou* matice \mathbf{A} . Korene charakteristickej rovnice matice \mathbf{A} sú nazývané *vlastnými číslami* matice \mathbf{A} . Nenulový vektor \mathbf{v} , spĺňujúci rovnicu $(\mathbf{A} - \lambda\mathbf{I})\mathbf{v} = \mathbf{0}$ pre nejaké vlastné číslo λ matice \mathbf{A} , je nazývaný *vlastným vektorom* príslušného vlastného čísla λ matice \mathbf{A} . Bez újmy na všeobecnosti bude ďalšom texte vždy uvažovaný vlastný vektor jednotkový (teda $\|\mathbf{v}\| = 1$).

Ak je matica \mathbf{A} reálna, štvorcová a symetrická a pre každý reálny nenulový vektor \mathbf{c} príslušného typu platí nerovnosť

$$\mathbf{c}'\mathbf{A}\mathbf{c} \geq 0, \quad (1.9)$$

potom matica \mathbf{A} je nazývaná *pozitívne semidefinitnou* maticou. Ak vyššie definovaná matica \mathbf{A} splňa pre každý nenulový vektor \mathbf{c} ostrú nerovnosť vo vzťahu 1.9, potom je matica \mathbf{A} nazývaná *pozitívne definitnou* maticou.

Veta 1.4. *Nech \mathbf{A}_m je reálna symetrická matica. Potom existuje ortogonálna matica \mathbf{B} s vlastnosťou $\mathbf{A} = \mathbf{BDB}'$ resp. $\mathbf{D} = \mathbf{B}'\mathbf{A}\mathbf{B}$, kde \mathbf{D} je diagonálna matica.*

Dôkaz. Vid' Rao ([9], str. 62). □

Následne bude ukázaný takýto rozklad $\mathbf{A} = \mathbf{BDB}'$ pre maticu \mathbf{A}_m , uvedený vo veta 1.4. Nech $\lambda_1, \dots, \lambda_m$ sú vlastné čísla matice \mathbf{A}_m , uvedené toľko krát ako je ich násobnosť (ako koreňa charakteristickej rovnice). Vlastné čísla $\lambda_1, \dots, \lambda_m$ sa po poradí priradia na diagonálu matice \mathbf{D} . Potom do každého stĺpca $\mathbf{b}_{\bullet i}$, $i = 1, \dots, m$ matice \mathbf{B} sa priradí vlastný vektor \mathbf{v}_i príslušný vlastnému číslu λ_i . V prípade k -nasobnosti ($k \leq m$) vlastného čísla λ_i , bude v matici \mathbf{B} uvedených všetkých n rôznych vlastných vektorov príslušných k vlastnému číslu λ_i . Vyššie uvedenému rozkladu symetrickej matice

$$\mathbf{A} = \mathbf{BDB}' = \sum_{j=1}^m d_{jj} \mathbf{b}_{\bullet j} \mathbf{b}'_{\bullet j} = \sum_{j=1}^m \lambda_j \mathbf{v}_j \mathbf{v}'_j, \quad (1.10)$$

sa hovorí *spektrálny rozklad*. Je zrejmé, že každá matica $\mathbf{v}_j \mathbf{v}'_j$, $j = 1, \dots, m$ je symetrická matica a preto aj ich lineárna kombinácia (určená vlastnými číslami) bude opäť tvoriť symetrickú maticu \mathbf{A} . Keďže v spektrálnom rozklade 1.10 je vždy uvažovaná matica \mathbf{B} regulárna, hodnosť matice \mathbf{A} je rovná počtu nenulových prvkov na diagonále matice \mathbf{D} .

Bez dôkazu bude uvedené nasledujúce tvrdenie. Reálna symetrická matica má reálne vlastné čísla a vlastné vektory. Pozitívne semidefinitná matica \mathbf{A}_m má vždy nezáporné vlastné čísla, pretože v definícii zo vzťahu 1.9 a vety 1.4. vyplýva

$$\mathbf{c}'\mathbf{A}\mathbf{c} = \mathbf{c}'\mathbf{BDB}'\mathbf{c} = \mathbf{e}'\mathbf{D}\mathbf{e} = \sum_{j=1}^m \lambda_j e_j^2 \geq 0, \quad (1.11)$$

kde $\mathbf{e} = \mathbf{B}'\mathbf{c}$. V posledné vyjadrenie vo vzťahu 1.11 platí pre ľubovoľný vektor \mathbf{c} resp. \mathbf{e} , len ak všetky vlastné čísla λ_j sú nezáporné. Pozitívne definitná matica \mathbf{A}_m má všetky vlastné čísla kladné. A to sa ukáže rovnako.

Ďalej bude tento odstavec venovaný špeciálnym maticiam, ktoré majú významné uplatnenie v teórii lineárnych regresných modelov. Dobrým príkladom je idempotentná alebo pseudoinverzná matica.

Idempotentná matica

Štvorcová matica \mathbf{A} bude nazývaná *idempotentnou maticou*, ak platí $\mathbf{AA} = \mathbf{A}$.

Veta 1.5. *Ak \mathbf{A} je matica idempotentná, potom platí $h(\mathbf{A}) = \text{Tr}(\mathbf{A})$.*

Dôkaz. Dôkaz je technického charakteru. Vid' Anděl ([2], str. 66). □

1.2. ZÁKLADNÉ POJMY A OZNAČENIA

Pseudoinverzná matica

Nech je daná matica $\mathbf{A}_{m \times n}$. Potom každá matica $\mathbf{B}_{n \times m}$ spĺňajúca rovnosť

$$\mathbf{ABA} = \mathbf{A}, \quad (1.12)$$

bude nazývaná *pseudoinverznou maticou* k matici \mathbf{A} . Pseudoinverzná matica k matici \mathbf{A} bude v ďalšom texte vyznačená horným indexom nasledovne \mathbf{A}^- . Len v prípade regulárnosti matice \mathbf{A} je pseudoinverzná matica určená jednoznačne. Pre regulárnu maticu \mathbf{A} , platí rovnica $\mathbf{ABA} = \mathbf{A}$ ekvivalentne s rovnicami $\mathbf{AB} = \mathbf{I} = \mathbf{BA}$. A tieto rovnice platia súčasne len pre maticu $\mathbf{B} = \mathbf{A}^{-1}$ inverznú k matici \mathbf{A} .

Nech je daná sústava lineárnych rovníc

$$\mathbf{Ab} = \mathbf{y}, \text{ kde } \mathbf{y} \in \mathcal{M}(\mathbf{A}). \quad (1.13)$$

Ďalej nech je uvažovaný vzťah $\mathbf{AA}^- \mathbf{Ab} = \mathbf{Ab}$, z čoho vyplýva $\mathbf{AA}^- \mathbf{y} = \mathbf{y}$ pri platnosti 1.13. Je zrejmé, že sústavy $\mathbf{AA}^- \mathbf{y} = \mathbf{y}$ a 1.13 súčasne platia, ak $\mathbf{b} = \mathbf{A}^- \mathbf{y}$. Záver je nasledovný, pre ľubovoľnú sústavu lineárnych rovníc 1.13 je riešením $\mathbf{b} = \mathbf{A}^- \mathbf{y}$ (matica sústavy \mathbf{A} môže byť singulárna i neštvorcová). Pseudoinverzná matica je teda zovšeobecnením pojmu inverznej matice pri riešení sústavy lineárnych rovníc.

Operátor ortogonálnej projekcie

Nech je uvažovaný vektorový priestor \mathcal{V} konečnej dimenzie na \mathbb{R} . Nech je ďalej uvažovaný podpriestor \mathcal{V}_1 vo \mathcal{V} a jeho ortogonálny komplement \mathcal{V}_1^\perp vo \mathcal{V} (platí $\mathcal{V}_1 \cap \mathcal{V}_1^\perp = \{\mathbf{0}\}$). Potom každý vektor $\mathbf{x} \in \mathcal{V}$ je možné jednoznačne vyjadriť dvomi vektormi $\mathbf{x}_1, \mathbf{x}_2$ nasledovne $\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2$, kde $\mathbf{x}_1 \in \mathcal{V}_1$ a $\mathbf{x}_2 \in \mathcal{V}_1^\perp$. Zobrazenie $\mathbf{P} : \mathcal{V} \rightarrow \mathcal{V}_1$ definované vzťahom $\mathbf{Px} = \mathbf{x}_1$ sa nazýva *ortogonálnym operátorom projekcie* na priestor \mathcal{V}_1 . Ďalej zrejme platí $\mathbf{P}(a\mathbf{x} + b\mathbf{y}) = a\mathbf{Px} + b\mathbf{Py}$ pre $a, b \in \mathbb{R}$ a teda \mathbf{P} je lineárnou transformáciou, ktorú je možné vyjadriť maticou. Ďalej bude označená takáto matica symbolom \mathbf{P} a bude nazývaná *projekčnou maticou*.

Veta 1.6. *Projekčná matica \mathbf{P} je idempotentná matica.*

Dôkaz. Zrejme platí vyššie uvedený vzťah $\mathbf{Px} = \mathbf{x}_1$ a vzťah $\mathbf{Px}_1 = \mathbf{x}_1$. Dosadením vyjadrenia \mathbf{x}_1 z prvého vzťahu do druhého je možné dostať vzťah $\mathbf{PPx} = \mathbf{Px}$, z čoho vyplýva $\mathbf{PP} = \mathbf{P}$. \square

Veta 1.7. *Ak \mathbf{P} je projekčná matica na podpriestor \mathcal{V}_1 , potom $\mathbf{I} - \mathbf{P}$ je projekčná matica na podpriestor \mathcal{V}_1^\perp .*

Dôkaz. Nech je uvažovaný ľubovoľný vektor $\mathbf{x} \in \mathcal{V}$ a príslušný vektor $\mathbf{x}_1 = \mathbf{Px}$, kde zrejme $\mathbf{x}_1 \in \mathcal{V}_1$. V prípade nenulovosti vektora $\mathbf{x} - \mathbf{x}_1 = \mathbf{x}_2$, platí $\mathbf{x}_2 \notin \mathcal{V}_1$ ale $\mathbf{x}_2 \in \mathcal{V}_1^\perp$. Potom $\mathbf{x}_2 = \mathbf{x} - \mathbf{x}_1 = \mathbf{Ix} - \mathbf{Px} = (\mathbf{I} - \mathbf{P})\mathbf{x}$ a teda $\mathbf{I} - \mathbf{P}$ je projekčnou maticou z priestoru \mathcal{V} na podpriestor \mathcal{V}_1^\perp . \square

Nech je uvažovaný n -rozmerný vektorový priestor \mathcal{V} so skalárnym súčinom zadaným pozitívne definitnou maticou \mathbf{W}_n nasledovne $\langle \cdot, \cdot \rangle_{\mathbf{W}} : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$, kde dvom vektormi $\mathbf{x}, \mathbf{y} \in \mathcal{V}$ priradí toto zobrazenie $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{W}}$ hodnotu $\mathbf{x}'\mathbf{W}\mathbf{y}$. V prípade, uvedenia skalárneho súčinu bez maticového indexu takto $\langle \cdot, \cdot \rangle$, bude uvažovaná matica \mathbf{W} jednotková \mathbf{I} .

Veta 1.8. *Nech je daný vektorový priestor \mathcal{V} , kde je definovaný skalárny súčin pozitívne definitnou maticou Σ . Potom matica \mathbf{P} je operátorom ortogonálnej projekcie, práve vtedy, keď matica \mathbf{P} je idempotentná a zároveň $\Sigma\mathbf{P}$ je symetrická matica.*

Dôkaz. Z vety 1.7. vyplýva, že každá projekčná matica musí mať vlastnosť idempotentnej matice. Nech sú uvažované dva ľubovoľné vektory $\mathbf{x}, \mathbf{y} \in \mathcal{V}$, potom $\mathbf{P}\mathbf{x} \in \mathcal{V}_1$, $(\mathbf{I} - \mathbf{P})\mathbf{y} \in \mathcal{V}_1^\perp$ a aby boli na seba ortogonálne, musí pre ne vždy platiť

$$0 = \langle \mathbf{P}\mathbf{x}, (\mathbf{I} - \mathbf{P})\mathbf{y} \rangle_\Sigma = \mathbf{x}'\mathbf{P}'\Sigma(\mathbf{I} - \mathbf{P})\mathbf{y}.$$

Predchádzajúci vzťah ekvivalentne platí, pre ľubovoľné vektory \mathbf{x} a \mathbf{y} , len ak platí $\mathbf{P}'\Sigma(\mathbf{I} - \mathbf{P}) = \mathbf{O}$ a teda $\mathbf{P}'\Sigma\mathbf{P} = \mathbf{P}'\Sigma$. Z tejto druhej rovnice a jej transponovanej verzii $\mathbf{P}'\Sigma\mathbf{P} = \Sigma\mathbf{P}$ vyplýva platnosť $\mathbf{P}'\Sigma = (\Sigma\mathbf{P})' = \Sigma\mathbf{P}$ a teda symetrickosť matice $\Sigma\mathbf{P}$. \square

V prípade skalárneho súčinu definovaného jednotkovou maticou \mathbf{I} je matica \mathbf{P} operátorom ortogonálnej projekcie práve vtedy, keď \mathbf{P} je idempotentná a symetrická matica.

2. Lineárny regresný model úplnej hodnoti

Nech $\mathbf{Y} = (Y_1, \dots, Y_n)'$ je náhodný vektor a $\mathbf{X}_{n \times k}$ je matica daných reálnych čísiel, kde zároveň platí $h(\mathbf{X}) = k$ a $k < n$. Ďalej nech pre náhodný vektor \mathbf{Y} platí

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad (2.1)$$

kde $\boldsymbol{\beta} = (\beta_1, \dots, \beta_k)'$ je vektor neznámych parametrov a $\mathbf{e} = (e_1, \dots, e_n)'$ je náhodný vektor spĺňajúci podmienky

$$\mathbf{E}\mathbf{e} = \mathbf{0}, \quad \text{Var}(\mathbf{e}) = \sigma^2 \mathbf{I}$$

s neznámym parametrom $\sigma^2 > 0$. V takomto prípade sa tiež povie, že \mathbf{Y} sa riadi lineárnym modelom a bude nazývaný *regresným*. Pretože vektor \mathbf{Y} závisí lineárne na prvkoch vektora $\boldsymbol{\beta}$, bude model 2.1 nazývaný *lineárnym regresným modelom*. Matica \mathbf{X} bude nazývaná *regresnou maticou* a jej stĺpce $\mathbf{x}_{\bullet 1}, \dots, \mathbf{x}_{\bullet k}$ budú nazývané *regresormi*. Vektor $\boldsymbol{\beta}$ bude nazývaný *vektorom regresných koeficientov*. Keďže regresná matica má rovnaký počet stĺpcov ako je jej hodnosť, bude model 2.1 nazývaný *model s úplnou stĺpcovou hodnotou* alebo skrátene *model úplnej hodnoti*.

Podmienka $h(\mathbf{X}) = k$ matice $\mathbf{X}_{n \times k}$ v modeli 2.1 určuje, že sa jedná o model úplnej hodnoti a teda stĺpce regresnej matice sú lineárne nezávislé. Zároveň ak matica \mathbf{X} je úplnej hodnoti, potom matica $\mathbf{X}'\mathbf{X}$ je regulárna. Predchádzajúce tvrdenie je tiež nepriamo dôsledkom vety 1.2. Priestor $\mathcal{M}(\mathbf{X})$ bude nazývaný *regresným priestorom*. Ďalej nech je zavedený pojem *reziduálneho priestoru*, ktorý bude vyjadrovať priestor kolmý na regresný priestor a bude označený $\mathcal{M}(\mathbf{X})^\perp$ (resp. $\mathcal{M}(\mathbf{X})^\perp$ je komplementárnym podpriestorom $\mathcal{M}(\mathbf{X})$ v \mathbb{R}^n).

Ďalej nech je uvažovaný všeobecnejší model úplnej hodnoti pre náhodný vektor \mathbf{Z} tvaru

$$\mathbf{Z} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad (2.2)$$

s neznámym vektorom $\boldsymbol{\beta}$ a neznámym náhodným vektorom \mathbf{e}

$$\mathbf{E}\mathbf{e} = \mathbf{0}, \quad \text{Var}(\mathbf{e}) = \sigma^2 \mathbf{G},$$

kde \mathbf{G} je známa pozitívne definitná matica. V prípade že matica \mathbf{G} nie je diagonálna, budú medzi pozorovaniami korelácie. Hodnoty \mathbf{Z} je možné lineárne transformovať nasledovne $\mathbf{Y} = \mathbf{G}_1^{-1}\mathbf{Z}$ a získať pôvodný model 2.1 pre \mathbf{Y} s $\text{Var}(\mathbf{e}) = \sigma^2 \mathbf{I}$, kde matica \mathbf{G}_1 je zavedená vzťahom $\mathbf{G} = \mathbf{G}_1\mathbf{G}_1'$. Pretože matica \mathbf{G} odpovedá variančnej matici, má všetky vlastné čísla kladné a platí, že existuje rozklad podľa vety 1.4, potom existuje aj vyššie zadaná matica \mathbf{G}_1 . Takto je možné previesť všeobecnejší model \mathbf{Z} , s ľubovoľnou známou pozitívne definitnou maticou \mathbf{G} , na pôvodný prípad modelu 2.1.

Regresné koeficienty β_1, \dots, β_k sa odhadujú *metódou najmenších štvorcov*, tj. minimalizáciou funkcie $S(\boldsymbol{\beta}) = (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})$. Tieto odhady budú spoločne označené vektorom $\mathbf{b} = (b_1, \dots, b_k)'$. Ilustratívny môže byť obrázok 3a, v 3. kapitole, kde je prípad uvažovaný na \mathbb{R}^2 .

Veta 2.1. Ak má matica \mathbf{X} úplnú stĺpcovú hodnotu, potom odhady metódou najmenších štvorcov sú $\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$.

2. LINEÁRNY REGRESNÝ MODEL ÚPLNEJ HODNOSTI

Dôkaz. Nech je uvažovaný taký vektor \mathbf{b} pre ktorý buď platí $\mathbf{Y} - \mathbf{Xb} \notin \mathcal{M}(\mathbf{X})$ alebo $\mathbf{Y} - \mathbf{Xb} = \mathbf{0}$. Potom samozrejme platí vzťah

$$\mathbf{X}'(\mathbf{Y} - \mathbf{Xb}) = \mathbf{0}, \quad (2.3)$$

ktorý bude využitý v nasledujúcom odvodzovaní.

$$\begin{aligned} (\mathbf{Y} - \mathbf{Xb})'(\mathbf{Y} - \mathbf{Xb}) &= [(\mathbf{Y} - \mathbf{Xb}) + (\mathbf{Xb} - \mathbf{Xb})]'[(\mathbf{Y} - \mathbf{Xb}) + (\mathbf{Xb} - \mathbf{Xb})] \\ &= (\mathbf{Y} - \mathbf{Xb})'(\mathbf{Y} - \mathbf{Xb}) + (\mathbf{b} - \mathbf{b})'\mathbf{X}'\mathbf{X}(\mathbf{b} - \mathbf{b}) \\ &\geq (\mathbf{Y} - \mathbf{Xb})'(\mathbf{Y} - \mathbf{Xb}) \end{aligned}$$

Pretože matica $\mathbf{X}'\mathbf{X}$ je regulárna, tak je i pozitívne definitná. A práve pre túto vlastnosť, v prípade $(\mathbf{b} - \mathbf{b}) = \mathbf{0}$, nadobúda súčet štvorcov $(\mathbf{Y} - \mathbf{Xb})'(\mathbf{Y} - \mathbf{Xb})$ svoje minimum $(\mathbf{Y} - \mathbf{Xb})'(\mathbf{Y} - \mathbf{Xb})$. Sústavu lineárnych rovníc 2.3 je možné prepísať na $\mathbf{X}'\mathbf{Xb} = \mathbf{X}'\mathbf{Y}$ a kvôli úplnej hodnosti \mathbf{X} , priamo dostať výsledok pre odhady \mathbf{b} . \square

Je vhodné si uvedomiť, že súčin 2.3 vyjadruje podmienku: vektor $\mathbf{Y} - \mathbf{Xb}$ je ortogonálny na každý stĺpec matice \mathbf{X} a teda zároveň na $\mathcal{M}(\mathbf{X})$. Tu môže byť vhodná geometrická predstava vid' obrázok 1.

Je možné získať aj iný pohľad na metódu najmenších štvorcov (minimalizáciu funkcie $S(\boldsymbol{\beta}) = (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})$) s ekvivalentným záverom, ako v dôkaze vety 2.1. Keďže funkcia $S(\boldsymbol{\beta})$ je konvexná na \mathbb{R}^k , pre neobmedzené odhady \mathbf{b} musí platiť

$$\frac{\partial}{\partial b_i}(\mathbf{Y} - \mathbf{Xb})'(\mathbf{Y} - \mathbf{Xb}) = 0, \quad i = 1, 2, \dots, k,$$

a zjednodušene platí

$$\frac{\partial}{\partial b_i}(\mathbf{Y} - \mathbf{x}_{\bullet i}b_i)'(\mathbf{Y} - \mathbf{x}_{\bullet i}b_i) = 0, \quad i = 1, 2, \dots, k, \quad (2.4)$$

kde $\mathbf{x}_{\bullet i}$ značí i -tý stĺpec matice \mathbf{X} . Každú z parciálnych derivácií zo vzťahu 2.4 je možné rozpísať nasledovne

$$\frac{\partial}{\partial b_i}(\mathbf{Y} - \mathbf{x}_{\bullet i}b_i)'(\mathbf{Y} - \mathbf{x}_{\bullet i}b_i) = -\mathbf{x}'_{\bullet i}(\mathbf{Y} - \mathbf{x}_{\bullet i}b_i) - (\mathbf{Y} - \mathbf{x}_{\bullet i}b_i)'\mathbf{x}_{\bullet i} = -2\mathbf{x}'_{\bullet i}(\mathbf{Y} - \mathbf{x}_{\bullet i}b_i).$$

Potom zo vzťahu 2.4 vyplýva platnosť sústavy lineárnych rovníc $\mathbf{x}'_{\bullet i}(\mathbf{Y} - \mathbf{x}_{\bullet i}b_i) = 0$, pre $i = 1, 2, \dots, k$, ktorú je možné zapísať maticovo

$$\mathbf{X}'(\mathbf{Y} - \mathbf{Xb}) = \mathbf{0} \quad \text{resp.} \quad \mathbf{X}'\mathbf{Xb} = \mathbf{X}'\mathbf{Y}. \quad (2.5)$$

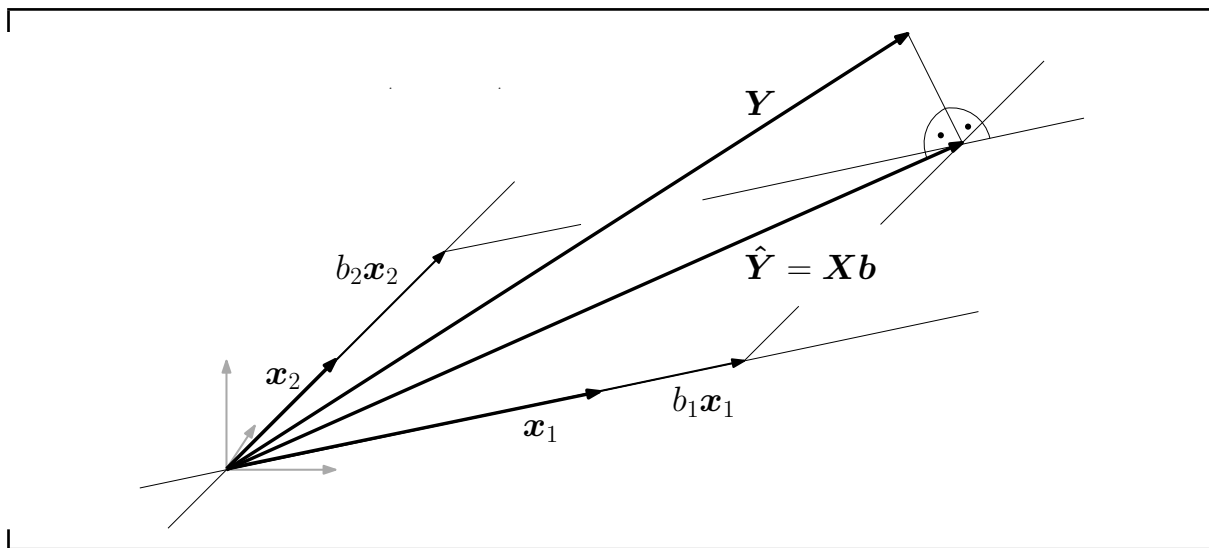
Sústava rovníc $\mathbf{X}'\mathbf{Xb} = \mathbf{X}'\mathbf{Y}$, vyplývajúca taktiež z dôkazu vety 2.1., určuje teda riešenia pre odhady \mathbf{b} regresných koeficientov $\boldsymbol{\beta}$ získaných metódou najmenších štvorcov. Je určite dobré si položiť otázku existencie riešenia, či vôbec pre takúto sústavu rovníc musí existovať riešenie. Z platnosti vzťahu $\mathbf{X}'\mathbf{Y} \in \mathcal{M}(\mathbf{X}')$ a vety 1.7. vyplýva platnosť $\mathbf{X}'\mathbf{Y} \in \mathcal{M}(\mathbf{X}'\mathbf{X})$. Keďže pravá strana sústavy $\mathbf{X}'\mathbf{Xb} = \mathbf{X}'\mathbf{Y}$ je podľa vety 1.2. taktiež z priestoru $\mathcal{M}(\mathbf{X}'\mathbf{X})$, má táto sústava vždy riešenie. V prípade úplnej hodnosti matice \mathbf{X} je riešenie tejto lineárnej sústavy vždy jednoznačné. Sústava lineárnych rovníc

$$\mathbf{X}'\mathbf{Xb} = \mathbf{X}'\mathbf{Y} \quad (2.6)$$

bude nazývaná *sústava normálnych rovníc* alebo skrátene *normálne rovnice* (*normálna rovnica*). Vektor

$$\hat{Y} = Xb = X(X'X)^{-1}X'Y \quad (2.7)$$

je najlepšou aproximáciou vektora Y , akú je možné vytvoriť v priestore $\mathcal{M}(X)$. Z geometrickej predstavy riešenia normálnych rovníc vyplýva, že sa jedná o kolmý priemet vektora Y na $\hat{Y} \in \mathcal{M}(X)$ a ten je určený jednoznačne. Táto situácia je zobrazená na obrázku 1.



Obrázok 1: Zobrazený je kolmý priemet vektora Y na vektor $\hat{Y} = Xb = b_1x_1 + b_2x_2$ do regresného priestoru $\mathcal{M}(X) = \mathcal{M}(x_1, x_2)$. Vektor b je odhadom vektora β metódou najmenších štvorcov.

Vektor Y je lineárne transformovaný do priestoru $\mathcal{M}(X)$ na vektor \hat{Y} pomocou vzťahu 2.7. V ďalšom texte bude H označovať maticu tejto lineárnej transformácie

$$H = X(X'X)^{-1}X'$$

Je zrejmé, že platí vzťah $HX = X$, totiž $X(X'X)^{-1}X'X = X$. Preto vektor z priestoru $\mathcal{M}(X)$ sa už touto ortogonálnou projekciou H nezmení, a teda každý stĺpec matice X sa transformuje na seba $HX = X$ z čoho vyplýva, že matica H je projekčnou maticou na priestor $\mathcal{M}(X)$. Zrejme ďalej platí tiež $(I - X(X'X)^{-1}X')X = O$, teda že stĺpce projekčnej matice $(I - X(X'X)^{-1}X')$ (z vety 1.7.) sú ortogonálne na regresný priestor $\mathcal{M}(X)$. Ďalej nech M vyjadruje $M = I - X(X'X)^{-1}X' = I - H$.

Zhrnutím platia vzťahy

$$HX = X, \quad M = I - H \quad \text{a} \quad MX = O. \quad (2.8)$$

Dá sa ľahko ukázať, že matice H a M sú okrem symetrickosti aj idempotentné, keďže platí

$$HH = X(X'X)^{-1}X'X(X'X)^{-1}X' = X(X'X)^{-1}X' = H,$$

$$MM = (I - H)(I - H) = I - H - H + HH = I - H = M.$$

Alebo inak, každý stĺpec z projekčnej matice H resp. M sa transformuje na seba do priestoru $\mathcal{M}(H)$ resp. $\mathcal{M}(M)$, teda platí $HH = H$ a $MM = M$. A tým sa ukázalo, že

2. LINEÁRNY REGRESNÝ MODEL ÚPLNEJ HODNOSTI

matice \mathbf{H} a \mathbf{M} sú symetrické a idempotentné. Podľa vety 1.8. sú obe matice operátormi ortogonálnej projekcie (na regresný a reziduálny priestor) a tým boli potvrdené doterajšie intuitívne predstavy.

Pre idempotentné matice, podľa vety 1.5., platia vzťahy

$$h(\mathbf{H}) = \text{Tr}(\mathbf{H}) = k, \quad h(\mathbf{M}) = \text{Tr}(\mathbf{M}) = n - k.$$

Matica \mathbf{H} je symetrická a preto môže byť uvažovaný spektrálny rozklad 1.10. Keďže $h(\mathbf{H}) = k$ platí, že matice \mathbf{H} má práve k nenulových vlastných čísiel (i s násobnosťami) s príslušnými vlastnými vektormi. Nech \mathbf{v} je nejaký vlastný vektor matice \mathbf{H} príslušný k nenulovému vlastnému číslu λ . Potom platí $\mathbf{H}\mathbf{v} = \lambda\mathbf{v}$. Keďže $\mathbf{v} \in \mathcal{M}(\mathbf{H})$ a \mathbf{H} je projekčná matice na $\mathcal{M}(\mathbf{H})$ musí platiť $\lambda = 1$. Preto nech pre \mathbf{H} sú označené vlastné čísla nasledovne $\lambda_i = 1, i = 1, \dots, k$ a $\lambda_i = 0, i = k+1, \dots, n$. Potom podľa spektrálneho rozkladu 1.10 platí pre matice \mathbf{H}

$$\mathbf{H} = \sum_{j=1}^m \lambda_j \mathbf{v}_j \mathbf{v}_j' = \sum_{j=1}^k \mathbf{v}_j \mathbf{v}_j' = \mathbf{Q}\mathbf{Q}', \quad (2.9)$$

kde matice $\mathbf{Q}_{n \times k}$ má stĺpce vyplnené vlastnými vektormi takto $\mathbf{Q} = (\mathbf{v}_1, \dots, \mathbf{v}_k)$. Vektory $\mathbf{v}_1, \dots, \mathbf{v}_k$ sú ortonormálne bázové vektory priestoru $\mathcal{M}(\mathbf{H})$ (a teda i regresného priestoru $\mathcal{M}(\mathbf{X})$) a preto platí $\mathbf{Q}'\mathbf{Q} = \mathbf{I}_k$. Obdobne pre matice \mathbf{M} existuje $n-k$ vlastných vektorov matice \mathbf{M} generujúcich priestor $\mathcal{M}(\mathbf{M})$ (a teda i reziduálneho priestoru $\mathcal{M}(\mathbf{X})^\perp$). Keďže sú priestory $\mathcal{M}(\mathbf{H})$ a $\mathcal{M}(\mathbf{M})$ na seba ortogonálne, ortonormálnou bázou sú práve vektory $\mathbf{v}_{k+1}, \dots, \mathbf{v}_n$ zo spektrálneho rozkladu 2.9 matice \mathbf{H} . Nech je teda, rovnako ako matice \mathbf{Q} , vytvorená matice $\mathbf{N}_{n \times n-k} = (\mathbf{v}_{k+1}, \dots, \mathbf{v}_n)$. Potom platí okrem vzťahu $\mathbf{M} = \mathbf{N}\mathbf{N}'$ aj vzťah $\mathbf{N}'\mathbf{N} = \mathbf{I}_{n-k}$. Pomocou platnosti $\mathbf{H} + \mathbf{M} = \mathbf{I}$ (zo vzťahu 2.8) je možné ukázať

$$\mathbf{I} = \mathbf{H} + \mathbf{M} = \sum_{j=1}^k \mathbf{v}_j \mathbf{v}_j' + \sum_{j=k+1}^n \mathbf{v}_j \mathbf{v}_j' = (\mathbf{Q}, \mathbf{N})(\mathbf{Q}, \mathbf{N})'. \quad (2.10)$$

A teda, že matice (\mathbf{Q}, \mathbf{N}) (v blokovom tvare) je skutočne ortogonálna, a splňa podmienky vety 1.4. Totiž len rovnosť $(\mathbf{Q}, \mathbf{N})'(\mathbf{Q}, \mathbf{N}) = \mathbf{I}$ je hneď zrejmé.

Náhodný vektor \mathbf{b}_L s k -zložkami bude nazývaný *lineárnym* odhadom parametra β (v lineárnom regresnom modeli 2.1), ak existuje matice $\mathbf{B}_{k \times n}$, že

$$\mathbf{b}_L = \mathbf{B}\mathbf{Y}. \quad (2.11)$$

Lineárny odhad \mathbf{b}_L bude nazývaný *nestranným* odhadom, ak platí vzťah $\mathbf{E}(\mathbf{b}_L) = \beta$ pre každé $\beta \in \mathbb{R}^k$. Nestranný lineárny odhad \mathbf{b}_L bude nazývaný *najlepším* nestranným odhadom, ak platí pre ľubovoľný iný nestranný lineárny odhad \mathbf{b}_L^* , že matice

$$(\text{Var}(\mathbf{b}_L^*) - \text{Var}(\mathbf{b}_L)) \quad (2.12)$$

je pozitívne semidefinitná. Vysvetlenie tejto definície môže byť nasledujúce. Odhad \mathbf{b}_L je najlepší, pretože sa rozptyl $\text{D}(\mathbf{c}'\mathbf{b}_L)$ ľubovoľnej lineárnej kombinácie odhadov $\mathbf{c}'\mathbf{b}_L$ zmení oproti rozptylu $\text{D}(\mathbf{c}'\mathbf{b}_L^*)$ o hodnotu $\mathbf{c}'(\text{Var}(\mathbf{b}_L^*) - \text{Var}(\mathbf{b}_L))\mathbf{c}$, ktorá je nezáporná (pretože matice 2.12 je pozitívne semidefinitnou). Teda odhad \mathbf{b}_L je najlepší, pretože pre ľubovoľný vektor \mathbf{c} minimalizuje rozptyl na $\text{D}(\mathbf{c}'\mathbf{b}_L)$ oproti ostatným nestranným lineárnym odhadom.

Veta 2.2. Pre odhad \mathbf{b} metódou najmenších štvorcov platí $E(\mathbf{b}) = \boldsymbol{\beta}$, $\text{Var}(\mathbf{b}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$.

Dôkaz.

$$\begin{aligned} E(\mathbf{b}) &= E(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E\mathbf{Y} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \boldsymbol{\beta}, \\ \text{Var}(\mathbf{b}) &= \text{Var}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}] = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\text{Var}\mathbf{Y})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\sigma^2\mathbf{I})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}. \end{aligned} \tag{2.13}$$

□

Teda bolo ukázané, že odhad metódou najmenších štvorcov $\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ uvedený vo vete 2.1. je nestranným lineárnym odhadom vektora $\boldsymbol{\beta}$.

Veta 2.3. (Gaussova-Markovova) V modeli \mathbf{Y} , pre ktorý platí $E\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$ a $\text{Var}\mathbf{Y} = \sigma^2\mathbf{I}$, je $\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$ najlepším nestranným lineárnym odhadom vektora $\mathbf{X}\boldsymbol{\beta}$. Zároveň platí $\text{Var}\hat{\mathbf{Y}} = \sigma^2\mathbf{H}$.

Dôkaz. Odhad $\hat{\mathbf{Y}}$ je lineárny, pretože $\hat{\mathbf{Y}}$ je lineárnou funkciou \mathbf{Y} ($\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$). Keďže platí

$$E\hat{\mathbf{Y}} = E\mathbf{H}\mathbf{Y} = \mathbf{H}\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta},$$

je odhad $\hat{\mathbf{Y}}$ nestranný. Ďalej nech je uvažovaný ľubovoľný iný lineárny odhad, teda napríklad $\tilde{\mathbf{Y}} = \mathbf{a} + \mathbf{B}\mathbf{Y}$. Aby bol lineárny odhad $\tilde{\mathbf{Y}}$ nestranný musí platiť

$$E\tilde{\mathbf{Y}} = E(\mathbf{a} + \mathbf{B}\mathbf{Y}) = \mathbf{a} + \mathbf{B}\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta} \quad \text{pre každé } \boldsymbol{\beta} \in \mathbb{R}^k.$$

Z toho vyplýva $\mathbf{a} = \mathbf{0}$ a $\mathbf{B}\mathbf{X} = \mathbf{X}$. Matica \mathbf{B} teda transformuje stĺpce matice \mathbf{X} na samé seba. Preto matica \mathbf{B} taktiež lineárne transformuje lineárnu kombináciu stĺpcov \mathbf{X} na seba. Preto $\mathbf{B}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ z čoho vyplýva rovnosť $\mathbf{B}\mathbf{H} = \mathbf{H}$. Posledná rovnosť sa dá prepísať na vzťah $(\mathbf{B} - \mathbf{H})\mathbf{H} = \mathbf{O}$, využitím idemotentnosti matice \mathbf{H} . Tento vzťah bude využitý nižšie. Pre variančnú maticu odhadu $\hat{\mathbf{Y}}$ platí

$$\text{Var}(\hat{\mathbf{Y}}) = \text{Var}(\mathbf{H}\mathbf{Y}) = \mathbf{H}\text{Var}(\mathbf{Y})\mathbf{H}' = \sigma^2\mathbf{H}.$$

Pre variančnú maticu lineárneho nestranného odhadu $\tilde{\mathbf{Y}}$ platí

$$\begin{aligned} \text{Var}(\tilde{\mathbf{Y}}) &= \mathbf{B}\sigma^2\mathbf{I}\mathbf{B}' = \sigma^2[(\mathbf{B} - \mathbf{H}) + \mathbf{H}][(\mathbf{B} - \mathbf{H}) + \mathbf{H}]' = \\ &= \sigma^2(\mathbf{B} - \mathbf{H})(\mathbf{B} - \mathbf{H})' + \sigma^2\mathbf{H}\mathbf{H}' = \sigma^2(\mathbf{B} - \mathbf{H})(\mathbf{B} - \mathbf{H})' + \text{Var}(\hat{\mathbf{Y}}), \end{aligned}$$

pretože $(\mathbf{B} - \mathbf{H})\mathbf{H} = \mathbf{O}$. Aby vektor $\hat{\mathbf{Y}}$ bol najlepší nestranný odhad, musí platiť, že matica $(\text{Var}(\tilde{\mathbf{Y}}) - \text{Var}(\hat{\mathbf{Y}}))$ je pozitívne semidefinitná matica podľa vzťahu 2.12. Podľa tvaru matice $(\mathbf{B} - \mathbf{H})(\mathbf{B} - \mathbf{H})'$, sa jedná o pozitívne semidefinitnú maticu, pretože je symetrická a má nezáporné vlastné čísla. Platí totiž ekvivalentne $(\mathbf{B} - \mathbf{H})\mathbf{I}(\mathbf{B} - \mathbf{H})'$, kde matica \mathbf{I} určuje nezápornosť vlastných čísiel. □

Dôležitým dôsledkom predchádzajúcej Gauss-Markovovej vety je pre odhad \mathbf{b} vektora $\boldsymbol{\beta}$, že tento odhad je dokonca najlepším nestranným lineárnym odhadom vektora $\boldsymbol{\beta}$. Totiž odhad \mathbf{b} je lineárnou funkciou $\hat{\mathbf{Y}}$ a preto je tiež najlepším nestranným lineárnym odhadom svojej strednej hodnoty, teda $\boldsymbol{\beta}$.

Minimálna hodnota funkcie $S(\boldsymbol{\beta})$ bude označená

$$RSS = S(\mathbf{b}) = (\mathbf{Y} - \hat{\mathbf{Y}})'(\mathbf{Y} - \hat{\mathbf{Y}}),$$

a nazývaná *reziduálnym súčtom štvorcov* (z angl. Residual sum of squares). Iné vyjádrenie pre štatistiku RSS popisuje nasledujúca veta.

2. LINEÁRNY REGRESNÝ MODEL ÚPLNEJ HODNOSTI

Veta 2.4. Pre reziduálny súčet štvorcov platí $RSS = \mathbf{Y}'\mathbf{M}\mathbf{Y}$ a $RSS = \mathbf{Y}'\mathbf{Y} - \mathbf{b}'\mathbf{X}'\mathbf{Y}$.

Dôkaz.

$$\begin{aligned} RSS &= (\mathbf{Y} - \hat{\mathbf{Y}})'(\mathbf{Y} - \hat{\mathbf{Y}}) = (\mathbf{Y} - \mathbf{H}\mathbf{Y})'(\mathbf{Y} - \mathbf{H}\mathbf{Y}) = \\ &= (\mathbf{M}\mathbf{Y})'(\mathbf{M}\mathbf{Y}) = \mathbf{Y}'\mathbf{M}'\mathbf{M}\mathbf{Y} = \mathbf{Y}'\mathbf{M}\mathbf{Y} = \\ &= \mathbf{Y}'(\mathbf{I} - \mathbf{H})\mathbf{Y} = \mathbf{Y}'\mathbf{Y} - \mathbf{Y}'\mathbf{H}\mathbf{Y} = \mathbf{Y}'\mathbf{Y} - (\mathbf{H}'\mathbf{Y})'\mathbf{Y} = \\ &= \mathbf{Y}'\mathbf{Y} - \hat{\mathbf{Y}}'\mathbf{Y} = \mathbf{Y}'\mathbf{Y} - \mathbf{b}'\mathbf{X}'\mathbf{Y}. \end{aligned}$$

□

Veta 2.5. Náhodná veličina $s^2 = RSS/(n - k)$ je nestranným odhadom parametra σ^2 .

Dôkaz. Pretože platí $RSS = \mathbf{Y}'\mathbf{M}\mathbf{Y}$, je možné vyvodiť

$$E(RSS) = E(\mathbf{Y}'\mathbf{M}\mathbf{Y}) = E[\text{Tr}(\mathbf{M}\mathbf{Y}\mathbf{Y}')] = \text{Tr}[\mathbf{M}E(\mathbf{Y}\mathbf{Y}')], \quad (2.14)$$

použitím vzťahu 1.6. Pre variančnú maticu platí vzťah

$$\text{Var}(\mathbf{Y}) = E(\mathbf{Y}\mathbf{Y}') - (E\mathbf{Y})(E\mathbf{Y})' = E(\mathbf{Y}\mathbf{Y}') - \mathbf{X}\boldsymbol{\beta}(\mathbf{X}\boldsymbol{\beta})'. \quad (2.15)$$

Zo vzťahu 2.15 je možné vyjadriť $E(\mathbf{Y}\mathbf{Y}') = \text{Var}(\mathbf{Y}) + \mathbf{X}\boldsymbol{\beta}(\mathbf{X}\boldsymbol{\beta})'$. Použitím predchádzajúceho vzťahu $E(\mathbf{Y}\mathbf{Y}')$ v odvodzovaní $E(RSS)$ vo vzťahu 2.14 sa pokračuje nasledovne

$$\begin{aligned} \text{Tr}[\mathbf{M}E(\mathbf{Y}\mathbf{Y}')] &= \text{Tr}[\mathbf{M}(\text{Var}(\mathbf{Y}) + \mathbf{X}\boldsymbol{\beta}(\mathbf{X}\boldsymbol{\beta})')] = \text{Tr}[\mathbf{M}(\mathbf{I}\sigma^2 + \mathbf{X}\boldsymbol{\beta}(\mathbf{X}\boldsymbol{\beta})')] \\ &= \sigma^2\text{Tr}(\mathbf{M}) + \text{Tr}[\mathbf{M}\mathbf{X}\boldsymbol{\beta}(\mathbf{X}\boldsymbol{\beta})'] = (n - k)\sigma^2 + \text{Tr}[\mathbf{O}] \\ &= (n - k)\sigma^2, \end{aligned}$$

kde bol využitý vyťah $\mathbf{M}\mathbf{X} = \mathbf{O}$.

□

Bude zavedené nové označenie pre vektor odhadov $\mathbf{u} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{M}\mathbf{Y} = \mathbf{M}\mathbf{e}$ náhodných zložiek \mathbf{e} , ktorý bude nazývaný *vektorom reziduí*. Potom platí tiež iné vyjadrenia pre reziduálny súčet štvorcov, napríklad $RSS = \mathbf{u}'\mathbf{u} = \|\mathbf{u}\|^2$.

Veta 2.6. Pre vektor reziduí platia vzťahy $E(\mathbf{u}) = \mathbf{0}$ a $\text{Var}(\mathbf{u}) = \sigma^2\mathbf{M}$.

Dôkaz. Na platnosť tejto vety sa dá ľahko nahliadnuť podľa nasledujúcich odvození

$$\begin{aligned} E(\mathbf{u}) &= E(\mathbf{M}\mathbf{Y}) = \mathbf{M}E\mathbf{Y} = \mathbf{M}\mathbf{X}\boldsymbol{\beta} = \mathbf{0}, \\ \text{Var}(\mathbf{u}) &= \text{Var}(\mathbf{M}\mathbf{Y}) = \mathbf{M}\text{Var}(\mathbf{Y})\mathbf{M}' = \sigma^2\mathbf{M}, \end{aligned}$$

kde bol použitý vyťah $\mathbf{M}\mathbf{X} = \mathbf{O}$.

□

V regresnej matici je v prvom stĺpci často uvažovaný vektor jednotiek. Potom sa formuluje matica $\mathbf{X} = (\mathbf{1}, \mathbf{X}_1)$, kde matica \mathbf{X}_1 je typu $n \times r$ ($r = k - 1$.) Potom sú preformulované koeficienty $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_r)'$ a odhady $\mathbf{b} = (b_0, b_1, \dots, b_r)'$.

V prípade, že bude uvažovaný vektor $\mathbf{1} \in \mathcal{M}(\mathbf{X})$, zo vzťahu 2.3 z ktorého boli odvodené normálne rovnice (teda metódu najmenších štvorcov) vyplýva

$$\mathbf{1}'(\mathbf{Y} - \mathbf{X}\mathbf{b}) = \mathbf{1}'(\mathbf{Y} - \hat{\mathbf{Y}}) = 0 \quad \text{a z toho vyplýva} \quad \frac{1}{n} \sum Y_i = \bar{Y} = \frac{1}{n} \sum \hat{Y}_i.$$

Potom bude výberový korelačný koeficient $r_{Y,\hat{Y}}$ medzi vektormi \mathbf{Y} a $\hat{\mathbf{Y}}$ zavedený nasledovne

$$r_{Y,\hat{Y}} = \frac{S_{Y,\hat{Y}}}{\sqrt{S_Y^2 S_{\hat{Y}}^2}} = \frac{\sum(Y_i - \bar{Y})(\hat{Y}_i - \bar{Y})/(n-1)}{\sqrt{\sum(Y_i - \bar{Y})^2/(n-1) \sum(\hat{Y}_i - \bar{Y})^2/(n-1)}},$$

kde $S_{Y,\hat{Y}}$ je výberová kovariancia a S_Y^2 , $S_{\hat{Y}}^2$ sú výberové rozptyly, podľa zavedenia v úvode.

Korelačný koeficient je dobrým ukazovateľom lineárnej závislosti. Preto pri "dobrej" aproximácii vektora \mathbf{Y} vektorom $\hat{\mathbf{Y}}$ by bol očakávaný koeficient korelácie blízky 1. K vyhodnoteniu kvality preloženia lineárnym modelom sa najčastejšie používa druhá mocnina tohto korelačného koeficientu. Túto štatistiku je možné vyjadriť pomocou vektorov nasledovne

$$r_{Y,\hat{Y}}^2 = \frac{S_{Y,\hat{Y}}}{\sqrt{S_Y^2 S_{\hat{Y}}^2}} = \frac{(\sum(Y_i - \bar{Y})(\hat{Y}_i - \bar{Y}))^2}{\sum(Y_i - \bar{Y})^2 \sum(\hat{Y}_i - \bar{Y})^2} = \frac{((\mathbf{Y} - \bar{\mathbf{Y}})'(\hat{\mathbf{Y}} - \bar{\mathbf{Y}}))^2}{\|\mathbf{Y} - \bar{\mathbf{Y}}\|^2 \|\hat{\mathbf{Y}} - \bar{\mathbf{Y}}\|^2}, \quad (2.16)$$

kde bol označený vektor $\bar{\mathbf{Y}} = \bar{Y}\mathbf{1}$. Vektor $\bar{\mathbf{Y}}$ odpovedá odhadu parametra modelu s regresnou maticou tvorenou len stĺpcom jednotiek. Vtedy je samozrejme najlepším odhadom stredná hodnota. Pre ďalšie odvodzovanie bude označený nový vektor $\mathbf{d} = \hat{\mathbf{Y}} - \bar{\mathbf{Y}}$, ktorý bude vyjadrovať jednotlivé vzdialenosti týchto odhadov. Viac o tejto problematike (vzdialenosti rôznych odhadov) bude pojednávať kapitola 4. V odvodzovaní $r_{Y,\hat{Y}}^2$ zo vzťahu 2.16 je možné pokračovať takto

$$r_{Y,\hat{Y}}^2 = \frac{((\mathbf{d} + \mathbf{u})'\mathbf{d})^2}{\|\mathbf{d} + \mathbf{u}\|^2 \|\mathbf{d}\|^2} = \frac{(\mathbf{d}'\mathbf{d} + \mathbf{u}'\mathbf{d})^2}{\|\mathbf{d} + \mathbf{u}\|^2 \|\mathbf{d}\|^2} = \frac{(\|\mathbf{d}\|^2 + 0)^2}{\|\mathbf{d} + \mathbf{u}\|^2 \|\mathbf{d}\|^2} = \frac{\|\mathbf{d}\|^2}{\|\mathbf{d} + \mathbf{u}\|^2}, \quad (2.17)$$

kde $\mathbf{u}'\mathbf{d} = \mathbf{Y}'\mathbf{M}'\mathbf{d} = \mathbf{u}'\mathbf{M}\mathbf{d} = 0$ z platnosti $\hat{\mathbf{Y}}, \bar{\mathbf{Y}} \in \mathcal{M}(\mathbf{X})$ vyplýva $\mathbf{d} \in \mathcal{M}(\mathbf{X})$ a teda $\mathbf{M}\mathbf{d} = \mathbf{0}$. Vektory \mathbf{u} a \mathbf{d} sú na seba kolmé, pretože sa oba nachádzajú v navzájom kolmých podpriestoroch a preto platí vzťah $\langle \mathbf{u}, \mathbf{d} \rangle = \mathbf{u}'\mathbf{d} = 0$.

Nech je uvažovaný výraz $\sum_{i=1}^n (Y_i - \bar{Y})^2$, ktorý určuje súčet štvorcov odchýlok v modeli s regresnou maticou $\mathbf{X} = \mathbf{1}\bar{Y}$. Táto hodnota bude označená TSS (z angl. Total sum of squares) a bude vyjadrená vzťahmi

$$TSS = (\mathbf{Y} - \bar{\mathbf{Y}})'(\mathbf{Y} - \bar{\mathbf{Y}}) = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - n\bar{Y}^2.$$

Je teda zrejmé, že v prípade regresnej matice tvorenej len stĺcom jednotiek a odhadu $\hat{\mathbf{Y}}$ získaného metódou najmenších štvorcov platí $RSS = TSS$. Toto nové označenie štatistiky TSS je možné využiť v ďalšom odvodzovaní štatistiky $r_{Y,\hat{Y}}^2$ zo vzťahu 2.17

$$r_{Y,\hat{Y}}^2 = \frac{\|\mathbf{d}\|^2}{\|\mathbf{d} + \mathbf{u}\|^2} = \frac{\|\mathbf{d} + \mathbf{u} - \mathbf{u}\|^2}{\|\mathbf{d} + \mathbf{u}\|^2} = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}. \quad (2.18)$$

Ako uvádza Anděl v literatúre ([3], str. 82), k popisu presnosti regresného modelu sa používa koeficient determinácie R^2 a *adjustovaný (korigovaný) koeficient determinácie* R_{adj}^2 , ktoré sú dané vzorcami

$$R^2 = 1 - \frac{RSS}{TSS}, \quad R_{adj}^2 = 1 - \frac{(n-1)RSS}{(n-r-1)TSS}. \quad (2.19)$$

2. LINEÁRNY REGRESNÝ MODEL ÚPLNEJ HODNOSTI

Vo vzťahu 2.18 je teda ukázané, že koeficient determinácie R^2 je priamo odvodený z druhej mocniny korelačného koeficientu $r_{Y,\hat{Y}}^2$.

Menej matematický pohľad na vzorec 2.19 koeficientu determinácie: Čím je hodnota R^2 bližšia k 1 (resp. hodnota RSS je výrazne menšia ako hodnota TSS), tým lepšie model zachytáva pozorované data. V hodnote adjustovaného koeficientu determinácie sú zohľadňované jednotlivé stupne voľnosti štatistik RSS a TSS .

2.1. Normálny lineárny regresný model

Pre predchádzajúce tvrdenia nebola vyžadovaná podmienka normálneho rozdelenia vektora \mathbf{Y} . Avšak pre nasledujúce tvrdenia, už bude treba uvažovať normálne rozdelenie, keďže je testovanie hypotéz založené na znalosti rozdelenia dat. Nech je teda predpokladané, že náhodná veličina \mathbf{e} má n -rozmerné normálne rozdelenie $\mathbf{e} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{I})$, z čoho vyplýva $\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$. Tento model bude ďalej nazývaný *normálny lineárny regresný model*.

Veta 2.7. *Pre normálny lineárny regresný model úplnej hodnosti platí*

$$\mathbf{b} \sim \mathbf{N}_k(\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}).$$

Dôkaz. Vzťah $\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ je známy. Matica $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ je maticu čísel typu $k \times n$, ktorá určuje lineárnu transformáciu vektora z normálneho rozdelenia $\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$. Parametre pre tento vektor boli určené vo vete 2.2. a preto podľa vzťahu 1.3 platí $\mathbf{b} \sim \mathbf{N}_k(\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1})$. □

Veta 2.8. *Pre odhady platí*

$$\hat{\mathbf{Y}} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{H}), \quad \mathbf{u} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{M}).$$

Dôkaz. Pre vektor $\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$ platí, že matica čísel \mathbf{H} určuje jeho lineárnu transformáciu z vektora $\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$. Preto aj vektor $\hat{\mathbf{Y}}$ má normálne rozdelenie podľa vzťahu 1.3. Z výpočtu parametrov

$$\begin{aligned} \mathbf{E}(\hat{\mathbf{Y}}) &= \mathbf{E}(\mathbf{H}\mathbf{Y}) = \mathbf{H}\mathbf{E}\mathbf{Y} = \mathbf{H}\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta}, \\ \text{Var}(\hat{\mathbf{Y}}) &= \text{Var}(\mathbf{H}\mathbf{Y}) = \mathbf{H}\text{Var}(\mathbf{Y})\mathbf{H}' = \sigma^2 \mathbf{H}, \end{aligned}$$

sa určí normálne rozdelenie $\hat{\mathbf{Y}} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{H})$. Podobne sa určí rozdelenie aj pre reziduá \mathbf{u} . Pre reziduá \mathbf{u} platí vzťah $\mathbf{u} = \mathbf{M}\mathbf{e}$, kde \mathbf{M} je matica čísel a vektor chýb \mathbf{e} má rozdelenie $\mathbf{e} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{I})$. Parametre odhadov \mathbf{u} boli určené vo vete 2.6. a z toho vyplýva, že odhady \mathbf{u} majú rozdelenie $\mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{M})$. □

Veta 2.9. *Platí $RSS/\sigma^2 \sim \chi^2(n - k)$.*

Dôkaz. Matica $\mathbf{N}_{n \times n-k}$ bola spomenutá vyššie, pri spektrálnom rozklade matice \mathbf{M} . Jej stĺpce boli ortonormálne bázové vektory reziduálneho priestoru $\mathcal{M}(\mathbf{X})^\perp$ a spĺňala vzťah $\mathbf{N}\mathbf{N}' = \mathbf{M}$. Nech je uvažovaný náhodný vektor $\mathbf{Z} \sim \mathbf{N}_{n-k}(\mathbf{0}, \mathbf{I}_{n-k})$, potom platí vzťah

2.1. NORMÁLNÝ LINEÁRNY REGRESNÝ MODEL

$\sigma \mathbf{NZ} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{M})$. Je zřejmé (vid' ostavec 1.2), že vektor reziduí \mathbf{u} a náhodný vektor $\sigma \mathbf{NZ}$ majú rovnaké rozdelenie $\mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{M})$. Pretože platí vzťah

$$\sigma^2 \mathbf{Z}' \mathbf{N}' \mathbf{N} \mathbf{Z} = \sigma^2 \mathbf{Z}' \mathbf{Z} \sim \sigma^2 \chi^2(n - k), \text{ vyplýva z toho } \frac{RSS}{\sigma^2} = \frac{\mathbf{u}' \mathbf{u}}{\sigma^2} \sim \chi^2(n - k).$$

□

Veta 2.10. Vektor \mathbf{b} a veličina s^2 sú nezávislé.

Dôkaz. Vzťahy $\mathbf{b} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}$ a $RSS = \mathbf{Y}' \mathbf{M} \mathbf{Y}$ sú známe. Anděl vo vete 4.19. literatúry ([3], str. 69), formuluje kedy kvadratická forma $\mathbf{Y}' \mathbf{M} \mathbf{Y}$ a náhodný vektor $(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}$ sú nezávislé. V tomto prípade, predpoklady uvedenej Andělovej vety sú: matica \mathbf{M} je pozitívne semidefinitná matica, $\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X} \boldsymbol{\beta}, \sigma^2 \mathbf{I})$ a zároveň platí

$$(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' (\sigma^2 \mathbf{I}) \mathbf{M} = \mathbf{O}. \quad (2.20)$$

Bolo už ukázané vyššie, pri výpisu spektrálneho rozkladu matice \mathbf{M} , že matica \mathbf{M} je symetrická a vlastné čísla matice \mathbf{M} sú rovné buď 0 alebo 1. Preto matica \mathbf{M} je pozitívne semidefinitná. Keďže \mathbf{M} je symetrická matica a platí $\mathbf{M} \mathbf{X} = \mathbf{O}$, platí tiež vzťah 2.20. Vektor \mathbf{b} a hodnota RSS sú teda nezávislé. Z toho vyplýva tiež nezávislosť medzi vektorom \mathbf{b} a veličinou s^2 . □

2.1.1. Testy hypotéz a intervaly spoľahlivosti

Veta 2.11. Nech v_{ij} značia prvky matice $(\mathbf{X}' \mathbf{X})^{-1}$. Potom platí pre náhodnú veličinu

$$\frac{b_i - \beta_i}{\sqrt{s^2 v_{ii}}} \sim t(n - k),$$

pre každé $i = 1, \dots, k$.

Dôkaz. Z vety 2.7 vyplýva $(b_i - \beta_i) / \sqrt{\sigma^2 v_{ii}} \sim \mathbf{N}(0, 1)$. Keďže $RSS / \sigma^2 \sim \chi^2(n - k)$ a vektor \mathbf{b} s odhadom s^2 sú nezávislé, je možné urobiť záver pre náhodnú veličinu

$$\frac{(b_i - \beta_i) / \sqrt{\sigma^2 v_{ii}}}{\sqrt{\frac{RSS / \sigma^2}{n - k}}} = \frac{(b_i - \beta_i)}{\sqrt{s^2 v_{ii}}} = T_i \sim t(n - k).$$

A teda, vyššie uvedená štatistika má Studentovo rozdelenie s $n - k$ stupňami voľnosti. □

Pomocou vety 2.11 sa testuje hypotéza $H_0 : \beta_i = \beta_i^0$ proti alternatívnej hypotéze $H_1 : \beta_i \neq \beta_i^0$, pre práve jedno $i = 1, \dots, k$. Hypotéza H_0 je zamietnutá na hladine α , práve keď, platí

$$\frac{|(b_i - \beta_i^0)|}{\sqrt{s^2 v_{ii}}} \geq t_{1 - \frac{\alpha}{2}}(n - k).$$

Potom obojstranný interval spoľahlivosti pre odhad parametra β_i so spoľahlivosťou $1 - \alpha$ je interval

$$\left(b_i - s \sqrt{v_{ii}} t_{1 - \frac{\alpha}{2}}(n - k), b_i + s \sqrt{v_{ii}} t_{1 - \frac{\alpha}{2}}(n - k) \right). \quad (2.21)$$

Ekvivalentne sa zamieta nulová hypotéza $H_0 : \beta_i = \beta_i^0$ na hladine významnosti α , ak hodnota β_i^0 neleží v intervale spoľahlivosti 2.21 so spoľahlivosťou $1 - \alpha$.

2. LINEÁRNY REGRESNÝ MODEL ÚPLNEJ HODNOSTI

Najčastejšie sa testuje nulová hypotéza $H_0 : \beta_i = 0$ a teda sa testuje, či je koeficient β_i významne odlišný od 0. Ak sa nulovú hypotézu H_0 zamietne, znamená to, že \mathbf{Y} významne (na hladine α) závisí na i -tom stĺpci regresnej matice \mathbf{X} . Avšak niekedy je potrebné testovať viacero regresných koeficientov súčasne. Preto bude uvedená nasledujúca veta, pre ktorú je potrebné rozdeliť vektor regresných koeficientov $\boldsymbol{\beta}$ na vektor $\boldsymbol{\beta}_1$ s p zložkami a vektor $\boldsymbol{\beta}_2$ s q zložkami s odpovedajúcou maticou $(\mathbf{X}'\mathbf{X})^{-1}$ takto

$$\boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{pmatrix}, \quad (\mathbf{X}'\mathbf{X})^{-1} = \begin{pmatrix} \mathbf{V}_{p \times p} & \mathbf{U} \\ \mathbf{U}' & \mathbf{W}_{q \times q} \end{pmatrix}.$$

Veta 2.12. *Platí*

$$\frac{1}{qs^2}(\mathbf{b}_2 - \boldsymbol{\beta}_2)' \mathbf{W}^{-1}(\mathbf{b}_2 - \boldsymbol{\beta}_2) \sim F(q, n - k).$$

Dôkaz. Pre marginálne rozdelenie platí $\mathbf{b}_2 \sim N_q(\boldsymbol{\beta}_2, \sigma^2 \mathbf{W})$. Keďže matica $\mathbf{X}'\mathbf{X}$ je pozitívne definitná, potom aj matica $(\mathbf{X}'\mathbf{X})^{-1}$ je pozitívne definitná (všetky vlastné čísla sú prevrátené hodnoty vlastných čísel matice $\mathbf{X}'\mathbf{X}$ a teda opäť kladné). Potom submatica \mathbf{W} je tiež pozitívne definitná a preto existuje rozklad $\mathbf{B}\mathbf{B}' = \mathbf{W}$, s regulárnou maticou \mathbf{B} . Nech je uvažovaný náhodný vektor $\boldsymbol{\beta}_2 + \mathbf{B}\mathbf{Z} \sim N_q(\boldsymbol{\beta}_2, \sigma^2 \mathbf{W})$, kde $\mathbf{Z} \sim N_q(\mathbf{0}, \sigma^2 \mathbf{I})$. Potom je možné kvadratickú formu

$$\sigma^{-2}(\mathbf{b}_2 - \boldsymbol{\beta}_2)' \mathbf{W}^{-1}(\mathbf{b}_2 - \boldsymbol{\beta}_2),$$

nahradiť za nasledujúci vzťah bez toho aby sa zmenilo rozdelenie tejto kvadratickej formy, takto

$$\sigma^{-2}(\mathbf{B}\mathbf{Z})'(\mathbf{B}\mathbf{B}')^{-1}(\mathbf{B}\mathbf{Z}) = \sigma^{-2}\mathbf{Z}'\mathbf{B}'(\mathbf{B}\mathbf{B}')^{-1}\mathbf{B}\mathbf{Z}.$$

Ďalej sa ukáže, že platí $\mathbf{B}'(\mathbf{B}\mathbf{B}')^{-1}\mathbf{B} = \mathbf{I}$. Pri vynásobení zľava maticou \mathbf{B} , nastáva rovnosť $\mathbf{B} = \mathbf{B}$. Týmto tvrdením je rovnosť $\mathbf{B}'(\mathbf{B}\mathbf{B}')^{-1}\mathbf{B} = \mathbf{I}$ dokázaná, len ak matica \mathbf{B} je regulárna. Následne platí

$$\sigma^{-2}\mathbf{Z}'\mathbf{B}'(\mathbf{B}\mathbf{B}')^{-1}\mathbf{B}\mathbf{Z} = \sigma^{-2}\mathbf{Z}'\mathbf{Z} \sim \chi^2(q).$$

Už bola ukázaná platnosť $RSS/\sigma^2 \sim \chi^2(n - k)$ a nezávislosť vektora \mathbf{b} s hodnotou RSS , potom platí

$$\frac{\sigma^{-2}(\mathbf{b}_2 - \boldsymbol{\beta}_2)' \mathbf{W}^{-1}(\mathbf{b}_2 - \boldsymbol{\beta}_2)}{RSS/\sigma^2} \frac{n - k}{q} = \frac{1}{qs^2}(\mathbf{b}_2 - \boldsymbol{\beta}_2)' \mathbf{W}^{-1}(\mathbf{b}_2 - \boldsymbol{\beta}_2) \sim F(q, n - k).$$

□

Dôsledok tejto vety má využitie pri testovaní nulovej hypotézu $H_0 : \boldsymbol{\beta}_2 = \boldsymbol{\beta}_2^0$ proti alternatívnej hypotéze $H_1 : \boldsymbol{\beta}_2 \neq \boldsymbol{\beta}_2^0$, kde $\boldsymbol{\beta}_2^0$ je q -prvkový vektor. Nulová hypotéza H_0 bude zamietnutá na hladine α , práve keď, vektor $\boldsymbol{\beta}_2^0$ nepatrí do konfidenčnej množiny pre $\boldsymbol{\beta}_2$ so spoľahlivosťou $1 - \alpha$

$$\left\{ \mathbf{b} \in \mathbb{R}^q \mid (\mathbf{b} - \boldsymbol{\beta}_2)' \mathbf{W}^{-1}(\mathbf{b} - \boldsymbol{\beta}_2) < qs^2 F_{1-\alpha}(q, n - k) \right\}. \quad (2.22)$$

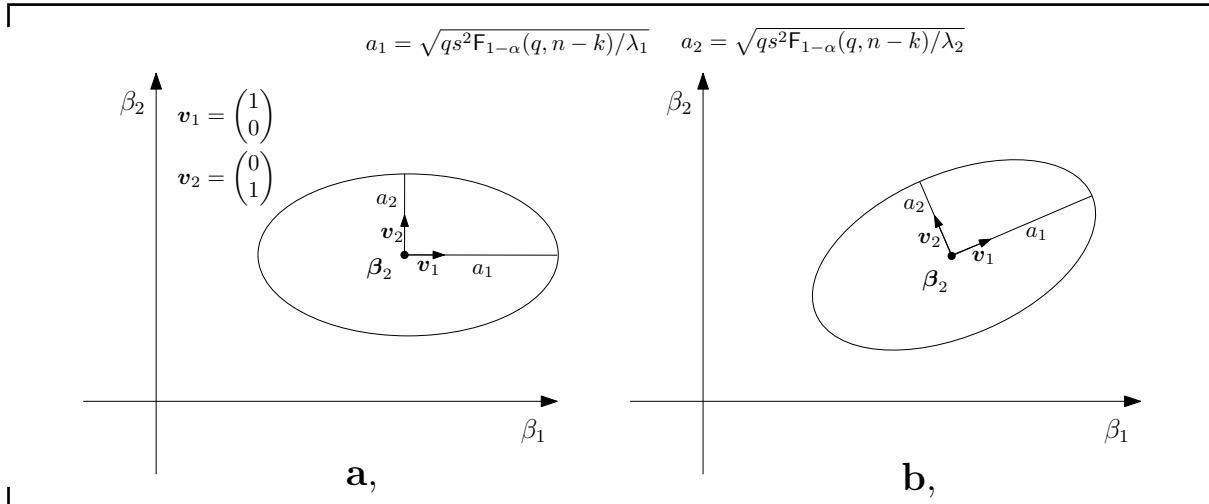
Pretože je matica \mathbf{W}^{-1} pozitívne definitná, konfidenčná množina 2.22 je tvaru elipsoidu¹ so stredom v bode $\boldsymbol{\beta}_2$. Pomocou spektrálneho rozkladu 1.10 matice $\mathbf{W}^{-1} = \sum_{i=1}^q \lambda_i \mathbf{v}_i \mathbf{v}_i'$, je možné prepísať nerovnosť v konfidenčnej množine 2.22 nasledovne

$$(\mathbf{b} - \boldsymbol{\beta}_2)' \left(\sum_{i=1}^q \lambda_i \mathbf{v}_i \mathbf{v}_i' \right) (\mathbf{b} - \boldsymbol{\beta}_2) = \sum_{i=1}^q \lambda_i ((\mathbf{b} - \boldsymbol{\beta}_2)' \mathbf{v}_i)^2 < qs^2 F_{1-\alpha}(q, n - k). \quad (2.23)$$

¹Myslený je všeobecnejší pojem elipsoidu, zahrňujúci dvojrozmerný (elipsa) a viacrozmerný prípad. V euklidovskom priestore dimenzie n , splňujúci $\sum_{j=1}^n x_j^2/r_j^2 < 1$, kde x_j sú súradnice bodu a $r_j \in \mathbb{R}$.

2.1. NORMÁLNÝ LINEÁRNY REGRESNÝ MODEL

Súčin $(\mathbf{b} - \boldsymbol{\beta}_2)' \mathbf{v}_i$ vyjadruje veľkosť premietaného vektora $(\mathbf{b} - \boldsymbol{\beta}_2)$ na priestor generovaný vlastným vektorom \mathbf{v}_i . V tomto prípade, vlastné čísla $\lambda_1, \dots, \lambda_q$ figurujú ako váhy sumy uvedenej vo vzťahu 2.23. Nech najväčším vlastným číslom matice \mathbf{W}^{-1} je λ_j , potom je elipsoid práve v smere vektora \mathbf{v}_j najviac obmedzený. Osami elipsoidu sú vlastné vektory $\mathbf{v}_1, \dots, \mathbf{v}_q$ matice \mathbf{W}^{-1} . Pre lepšiu predstavu je uvedený obrázok 2, ktorý je uvažovaný pre $q = 2$.



Obrázok 2: Konfidenčná množina v **a**, je prípad kedy matice \mathbf{W}^{-1} typu 2×2 je diagonálna a teda vlastné vektory sú rovnobežné s euklidovskou bázou. Konfidenčná množina v **b**, je všeobecnejším prípadom s vlastnými vektormi \mathbf{v}_1 a \mathbf{v}_2 . V oboch prípadoch je vlastné číslo λ_2 väčšie, a_1 vyjadruje dĺžku hlavnej poloosy a a_2 dĺžku vedľajšej poloosy elíps.

Väčšinou je nulová hypotéza nastavená $H_0 : \boldsymbol{\beta}_2 = \mathbf{0}$ proti alternatívnej hypotéze $H_1 : \boldsymbol{\beta}_2 \neq \mathbf{0}$. Teda je testovaná situácia, že posledných q regresných koeficientov v $\boldsymbol{\beta}$ je súčasne nevýznamných pri zvolenej hladine α . V prípade, že je testovaný celý vektor $H_0 : \boldsymbol{\beta} = \mathbf{0}$ proti alternatíve $H_1 : \boldsymbol{\beta} \neq \mathbf{0}$, je pokladané $\mathbf{W} = (\mathbf{X}'\mathbf{X})^{-1}$ a je možné odvodiť vzťah

$$\frac{\mathbf{b}'\mathbf{X}'\mathbf{X}\mathbf{b}}{ks^2} = \frac{\hat{\mathbf{Y}}'\hat{\mathbf{Y}}}{ks^2} = \frac{\mathbf{Y}'\mathbf{Y} - RSS}{ks^2} \sim F_{k,n-k}.$$

Je známy vzťah $\mathbf{b} \sim N_k(\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1})$ z vety 2.7. Ďalej nech je uvažovaný ľubovoľný vektor $\mathbf{c} \in \mathbb{R}^k$. Potom pre odhad lineárnej kombinácie regresných koeficientov $\mathbf{c}'\mathbf{b}$ platí

$$\begin{aligned} E(\mathbf{c}'\mathbf{b}) &= \mathbf{c}'E(\mathbf{b}) = \mathbf{c}'\boldsymbol{\beta}, \\ \text{Var}(\mathbf{c}'\mathbf{b}) &= \mathbf{c}'\text{Var}(\mathbf{b})\mathbf{c} = \sigma^2\mathbf{c}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{c}. \end{aligned}$$

Z definície k -rozmerného normálneho rozdelenia pre náhodný vektor vo vzťahu 1.2, platí $\mathbf{c}'\mathbf{b} \sim N(\mathbf{c}'\boldsymbol{\beta}, \sigma^2\mathbf{c}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{c})$.

Veta 2.13. Pre nenulový vektor $\mathbf{c} = (c_1, \dots, c_n)'$ platí

$$\frac{\mathbf{c}'\mathbf{b} - \mathbf{c}'\boldsymbol{\beta}}{\sqrt{s^2\mathbf{c}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{c}}} \sim t(n-k).$$

2. LINEÁRNY REGRESNÝ MODEL ÚPLNEJ HODNOSTI

Dôkaz. Z vety 2.10. vyplýva, že taktiež hodnoty $\mathbf{c}'\mathbf{b}$ a s^2 sú nezávislé. Potom platí $(\mathbf{c}'\mathbf{b} - \mathbf{c}'\boldsymbol{\beta})/\sqrt{\sigma^2\mathbf{c}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{c}} \sim \mathbf{N}(0, 1)$ a z vety 2.9. platí $RSS/\sigma^2 \sim \chi^2(n - k)$. Preto platí

$$\frac{(\mathbf{c}'\mathbf{b} - \mathbf{c}'\boldsymbol{\beta})/\sqrt{\sigma^2\mathbf{c}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{c}}}{\sqrt{RSS\sigma^{-2}/(n - k)}} = \frac{\mathbf{c}'\mathbf{b} - \mathbf{c}'\boldsymbol{\beta}}{\sqrt{s^2\mathbf{c}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{c}}} \sim \mathbf{t}(n - k).$$

Tvrdenie vyplýva z definovania Studentovho rozdelenia. □

Predošlá veta 2.13. je užitočná pri testovaní hypotéz o lineárnych kombináciach regresných koeficientov. Príkladom môže byť odhad individuálnej hodnoty (funkcie strednej hodnoty) v bode \mathbf{c} . Následne bude uvedený aspoň obojstranný interval spoľahlivosti pre túto hodnotu $\mathbf{c}'\boldsymbol{\beta}$ so spoľahlivosťou $1 - \alpha$

$$\left(\mathbf{c}'\mathbf{b} - t_{1-\frac{\alpha}{2}}(n - k)\sqrt{s^2\mathbf{c}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{c}}, \mathbf{c}'\mathbf{b} - t_{1-\frac{\alpha}{2}}(n - k)\sqrt{s^2\mathbf{c}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{c}} \right).$$

Tvrdenie vety 2.13. už nebude všeobecne platné pre ľubovoľný nenulový vektor \mathbf{c} , v prípade modelu s neúplnou hodnosťou. Tým sa bude zaoberať nasledujúca kapitola.

3. Lineárny regresný model neúplnej hodnoty

Najprv bude zavedený lineárny regresný model neúplnej hodnoty, rovnako ako model plnej hodnoty 2.1.

Nech $\mathbf{Y} = (Y_1, \dots, Y_n)'$ je náhodný vektor a $\mathbf{X}_{n \times k}$ je matica daných reálnych čísel, kde zároveň platí $h(\mathbf{X}) = r$ a $r < k < n$. Ďalej nech sa \mathbf{Y} riadi lineárnym modelom

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad (3.1)$$

keď $\boldsymbol{\beta} = (\beta_1, \dots, \beta_k)'$ je vektor neznámych parametrov a $\mathbf{e} = (e_1, \dots, e_n)'$ je náhodný vektor spĺňajúci podmienky

$$\mathbb{E}\mathbf{e} = \mathbf{0}, \quad \text{Var}(\mathbf{e}) = \sigma^2 \mathbf{I}$$

s neznámym parametrom $\sigma^2 > 0$. Bude uvažované rovnaké pomenovanie pre \mathbf{Y} , \mathbf{X} a $\boldsymbol{\beta}$ ako v modeli 2.1. V tomto prípade, keď regresná matica \mathbf{X} nemá úplnú stĺpcovú hodnotu, bude model 3.1 nazývaný *modelom s neúplnou stĺpcovou hodnotou* alebo skrátene *modelom neúplnej hodnoty*.

Tak ako v predchádzajúcej kapitole, je dôležitou úlohou uvažovať také odhady \mathbf{b} regresných koeficientov $\boldsymbol{\beta}$ aby minimalizovali funkciu $S(\boldsymbol{\beta}) = (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})$. Je dobré pripomenúť, že na minimalizácii tejto kvadratickej funkcie $S(\boldsymbol{\beta})$ je založená celá metóda najmenších štvorcov. Otázkou je, či tak, ako v modeli úplnej hodnoty stačí splnenie normálnych rovníc 2.6

$$\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{Y},$$

pre odhady \mathbf{b} (minimalizujúcich $S(\boldsymbol{\beta})$). Nasledujúca veta o tomto pojednáva.

Veta 3.1. *Ak matica \mathbf{X} nemá úplnú stĺpcovú hodnotu, potom všetky odhady metódou najmenších štvorcov sú tvaru*

$$\mathbf{b} = \mathbf{b}_p + \mathbf{b}_h,$$

kde $\mathbf{b}_p = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}$, $\mathbf{b}_h \in \mathcal{M}(\mathbf{X}'\mathbf{X})^\perp$ a $(\mathbf{X}'\mathbf{X})^{-}$ je ľubovoľná pseudoinverzia k $\mathbf{X}'\mathbf{X}$.

Dôkaz. Tento dôkaz bude obdobný dôkazu vety 2.1. Rovnako je uvažovaný nejaký vektor \mathbf{b} , pre ktorý bude platiť buď $\mathbf{Y} - \mathbf{X}\mathbf{b} \notin \mathcal{M}(\mathbf{X})$ alebo $\mathbf{Y} - \mathbf{X}\mathbf{b} = \mathbf{0}$. Taktiež platia vzťahy $\mathbf{X}'(\mathbf{Y} - \mathbf{X}\mathbf{b}) = \mathbf{0}$ a následné odvodenie

$$\begin{aligned} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) &= (\mathbf{Y} - \mathbf{X}\mathbf{b})'(\mathbf{Y} - \mathbf{X}\mathbf{b}) + (\mathbf{b} - \boldsymbol{\beta})'\mathbf{X}'\mathbf{X}(\mathbf{b} - \boldsymbol{\beta}) \\ &\geq (\mathbf{Y} - \mathbf{X}\mathbf{b})'(\mathbf{Y} - \mathbf{X}\mathbf{b}). \end{aligned}$$

Keďže matica $\mathbf{X}'\mathbf{X}$ je singularná, tak je pozitívne semidefinitná. Riešenie z rovníc $\mathbf{X}'(\mathbf{Y} - \mathbf{X}\mathbf{b}) = \mathbf{0}$ teda minimalizuje funkciu $S(\boldsymbol{\beta})$ na minimum $(\mathbf{Y} - \mathbf{X}\mathbf{b})'(\mathbf{Y} - \mathbf{X}\mathbf{b})$. Predchádzajúcu sústavu lineárnych rovníc je možné upraviť na tvar $\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{Y}$. Avšak kvôli neúplnej hodnote \mathbf{X} , nie je riešenie \mathbf{b} jediné a dá sa rozdeliť na partikulárne $\mathbf{b}_p = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}$ a homogénne riešenie $\mathbf{b}_h \in \mathcal{M}(\mathbf{X}'\mathbf{X})^\perp$. Totiž pre ľubovoľné riešenie partikulárne riešenie $\mathbf{b}_p = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}$ a homogénne riešenie $\mathbf{b}_h \in \mathcal{M}(\mathbf{X}'\mathbf{X})^\perp$ platia vzťahy

$$\mathbf{X}'\mathbf{X}\mathbf{b}_p = \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y} = \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{Y} \quad \text{a} \quad \mathbf{X}'\mathbf{X}\mathbf{b}_h = \mathbf{0}.$$

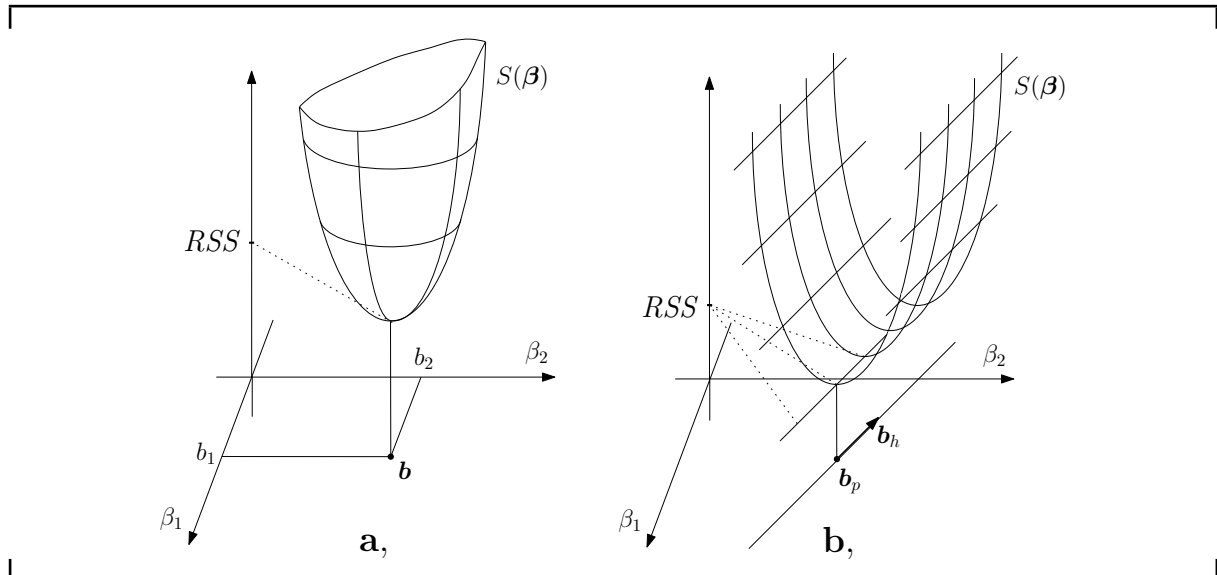
□

Je potreba pár slovami doplniť záverečné myšlienky dôkazu. Ak je matica $\mathbf{X}'\mathbf{X}$ singularná, výpočet inverznej matice k nej nie je možný. Pre riešenie sústavy lineárnych rovníc $\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{Y}$, so singularnou maticou sústavy, bola v teórii matic zavedená pseudoinverzná matica (ako bolo uvedené v úvode v časti o pseudoiverznej matici). V tomto prípade, partikulárnym riešením normálnych rovníc sú

$$\mathbf{b}_p = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}.$$

V úvode bolo poznamenané, že pseudoinverzná matica $(\mathbf{X}'\mathbf{X})^{-}$ nie je jednoznačne daná pre singularnú maticu $\mathbf{X}'\mathbf{X}$. V tomto prípade, je partikulárných riešení \mathbf{b}_p nekonečne veľa, ktoré sa líšia o vektor homogénneho riešenia \mathbf{b}_h . Preto ak matica \mathbf{X} nemá úplnú stĺpcovú hodnotu, majú normálne rovnice vždy nekonečne veľa riešení tvaru $\mathbf{b} = \mathbf{b}_p + \mathbf{b}_h$ ako je znázornené v obrázku 3.

Minimálna hodnota funkcie $S(\mathbf{b})$, pre nejaké riešenie \mathbf{b} normálnych rovníc, bude nazývaná *reziduálny súčet štvorcov* a označená RSS (tak ako tomu bolo v predchádzajúcej kapitole).



Obrázok 3: V obrázku **a**, je zobrazená funkcia $S(\boldsymbol{\beta}) = (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})$ pre rôzne odhady $\boldsymbol{\beta}$, v prípade lineárneho modelu úplnej hodnoty. V obrázku **b**, zobrazená táto funkcia pre model neúplnej hodnoty. Tu \mathbf{b}_p vyjadruje nejaké partikulárne riešenie normálnych rovníc $\mathbf{b}_p = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}$.

Geometrická predstava: bez úplnej hodnoty regresnej matice, sú niektoré stĺpce matice \mathbf{X} lineárne závislé a práve preto nie je možné jednoznačne vyjadriť odhad vektora $\boldsymbol{\beta}$. V takýchto prípadoch uľahčuje prácu s odhadmi, v modeloch neúplnej hodnoty, pojem odhadnuteľnosť parametrickej funkcie.

3.1. Odhadnuteľné parametrické funkcie

V prípade úplnej stĺpcovej hodnoty regresnej matice je odhad metódou najmenších štvorcov parametra $\mathbf{t}'\boldsymbol{\beta}$ rovný $\mathbf{t}'\mathbf{b}$ pre ľubovoľný vektor \mathbf{t} daného typu. Vektor \mathbf{b} je totiž jediným riešením sústavy normálnych rovníc. No otázkou je, ako je to v prípade nekonečného počtu

3.1. ODHADNUTEĽNÉ PARAMETRICKÉ FUNKCIE

riešení vektora koeficientov \mathbf{b} sústavy normálnych rovníc.

Nech je daný nejaký nenáhodný vektor $\mathbf{t} \in \mathbb{R}^k$. Potom lineárna funkcia $\mathbf{t}'\boldsymbol{\beta}$ bude nazývaná *odhadnutel'nou*, ak existuje jej nestranný lineárny odhad.

Veta 3.2. *Lineárna funkcia $\mathbf{t}'\boldsymbol{\beta}$ je odhadnutel'ná práve vtedy, keď platí*

$$\mathbf{t} \in \mathcal{M}(\mathbf{X}') = \mathcal{M}(\mathbf{X}'\mathbf{X}). \quad (3.2)$$

Dôkaz. Aby pre nejaký vektor $\mathbf{s} \in \mathbb{R}^n$ bola lineárna funkcia $\mathbf{s}'\mathbf{Y}$ nestranným odhadom $\mathbf{t}'\boldsymbol{\beta}$, musí pre každé $\boldsymbol{\beta} \in \mathbb{R}^k$ platiť

$$E(\mathbf{s}'\mathbf{Y}) = \mathbf{s}'\mathbf{X}\boldsymbol{\beta} = \mathbf{t}'\boldsymbol{\beta}. \quad (3.3)$$

Odkadiaľ $\mathbf{t} = \mathbf{X}'\mathbf{s}$. Preto ak je lineárna funkcia $\mathbf{t}'\boldsymbol{\beta}$ odhadnutel'nou, potom $\mathbf{t} \in \mathcal{M}(\mathbf{X}')$ (resp. $\mathbf{t} \in \mathcal{M}(\mathbf{X}'\mathbf{X})$). A obrátene, ak $\mathbf{t} \in \mathcal{M}(\mathbf{X}')$ potom musí existovať vektor $\mathbf{s} \in \mathbb{R}^n$ taký, že platí $\mathbf{t} = \mathbf{X}'\mathbf{s}$. Potom platí vzťah 3.3, a preto je lineárna funkcia $\mathbf{t}'\boldsymbol{\beta}$ odhadnutel'ná. \square

Veta 3.3. *Ak je lineárna funkcia $\mathbf{t}'\boldsymbol{\beta}$ odhadnutel'nou, potom $\mathbf{t}'\mathbf{b}$ je nestranným lineárnym odhadom funkcie $\mathbf{t}'\boldsymbol{\beta}$, kde \mathbf{b} je ľubovoľné riešenie normálnych rovníc. Zároveň je hodnota $\mathbf{t}'\mathbf{b}$ rovnaká pre všetky riešenia \mathbf{b} normálnych rovníc.*

Dôkaz. Najprv bude ukázaný dôkaz druhého tvrdenia. Nech $\mathbf{b} = \mathbf{b}_p + \mathbf{b}_h$ a $\mathbf{b}^* = \mathbf{b}_p^* + \mathbf{b}_h^*$ sú rôzne riešenia normálnych rovníc vyjadrené podľa vety 3.1. Aby platilo $\mathbf{t}'\mathbf{b} = \mathbf{t}'\mathbf{b}^*$, musí ekvivalentne platiť $\mathbf{t}'(\mathbf{b} - \mathbf{b}^*) = 0$. Platí vzťah

$$\mathbf{t}'(\mathbf{b} - \mathbf{b}^*) = \mathbf{t}'(\mathbf{b}_p - \mathbf{b}_p^*) + \mathbf{t}'(\mathbf{b}_h - \mathbf{b}_h^*) = \mathbf{t}'(\mathbf{b}_p - \mathbf{b}_p^*), \quad (3.4)$$

pretože $\mathbf{b}_h, \mathbf{b}_h^* \in \mathcal{M}(\mathbf{X}'\mathbf{X})^\perp$ (z čoho vyplýva $(\mathbf{b}_h - \mathbf{b}_h^*) \in \mathcal{M}(\mathbf{X}'\mathbf{X})^\perp$) a zároveň $\mathbf{t} \in \mathcal{M}(\mathbf{X}'\mathbf{X})$. Nech $\mathbf{b}_p = (\mathbf{X}'\mathbf{X})_1^- \mathbf{X}'\mathbf{Y}$ a $\mathbf{b}_p^* = (\mathbf{X}'\mathbf{X})_2^- \mathbf{X}'\mathbf{Y}$ podľa vety 3.1., kde sú uvedené dve rôzne pseudoinverzné matice (označené indexami 1 a 2) k matici $\mathbf{X}'\mathbf{X}$. Anděl vo vete 15.c, v literatúre ([2], str. 69) uviedol, že matica

$$\mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}' \quad (3.5)$$

je nezávislá na zvolenej pseudoinverzii $(\mathbf{X}'\mathbf{X})^-$. Použitím predchádzajúceho tvrdenia a vzťahu $\mathbf{t} = \mathbf{X}'\mathbf{s}$ (z odhadnutel'nosti $\mathbf{t}'\boldsymbol{\beta}$) je možné ďalej pokračovať v úprave vzťahu 3.4 takto

$$\mathbf{t}'(\mathbf{b}_p - \mathbf{b}_p^*) = \mathbf{s}'\mathbf{X}[(\mathbf{X}'\mathbf{X})_1^- \mathbf{X}'\mathbf{Y} - (\mathbf{X}'\mathbf{X})_2^- \mathbf{X}'\mathbf{Y}] = \mathbf{s}'[\mathbf{X}(\mathbf{X}'\mathbf{X})_1^- \mathbf{X}' - \mathbf{X}(\mathbf{X}'\mathbf{X})_2^- \mathbf{X}']\mathbf{Y} = 0.$$

Z toho vyplýva rovnosť $\mathbf{t}'(\mathbf{b} - \mathbf{b}^*) = 0$ resp. $\mathbf{t}'\mathbf{b} - \mathbf{t}'\mathbf{b}^* = 0$, pre ľubovoľné riešenia normálnych rovníc. Nasleduje dôkaz prvého tvrdenia. Použitím vzťahu $\mathbf{t} = \mathbf{X}'\mathbf{X}\mathbf{g}$ (z odhadnutel'nosti $\mathbf{t}'\boldsymbol{\beta}$) je možné vyjadriť nasledujúci vzťah

$$E(\mathbf{t}'\mathbf{b}) = E(\mathbf{g}'\mathbf{X}'\mathbf{X}\mathbf{b}) = E(\mathbf{g}'\mathbf{X}'\mathbf{Y}) = \mathbf{g}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{t}'\boldsymbol{\beta}, \quad (3.6)$$

kde v druhej rovnosti bola využitá platnosť normálnych rovníc $\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{Y}$ (keďže platí pre každé \mathbf{b} spĺňajúce normálne rovnice). Zo vzťahu 3.6 vyplýva, že náhodná veličina $\mathbf{t}'\mathbf{b}$ je nestranným lineárnym odhadom funkcie $\mathbf{t}'\boldsymbol{\beta}$ (lineárnosť podľa vzťahu 2.11). \square

Keďže pri úplnej hodnosti \mathbf{X} je jednoznačne určené riešenie \mathbf{b} normálnej rovnice, potom pre ľubovoľnú lineárnu transformáciu \mathbf{t}' je funkcia $\mathbf{t}'\boldsymbol{\beta}$ odhadnutel'ná.

3. LINEÁRNY REGRESNÝ MODEL NEÚPLNEJ HODNOSTI

Veta 3.4. Ak je parametrická funkcia $\mathbf{t}'\boldsymbol{\beta}$ odhadnuteľná, potom $\mathbf{t}'\mathbf{b}$ je najlepší nestranný lineárny odhad parametrickej funkcie $\mathbf{t}'\boldsymbol{\beta}$.

Dôkaz. Nestrannosť a lineárnosť odhadu $\mathbf{t}'\mathbf{b}$ bola ukazaná vo vete 3.3. Z odhadnuteľnosti funkcie $\mathbf{t}'\boldsymbol{\beta}$, vyplýva existencia $\mathbf{g} \in \mathbb{R}^k$ tak, že platí $\mathbf{t} = \mathbf{X}'\mathbf{X}\mathbf{g}$. Potom platí

$$\mathbf{t}'\mathbf{b} = \mathbf{g}'\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{g}'\mathbf{X}'\mathbf{Y}, \quad (3.7)$$

kde v druhej rovnosti boli použité normálne rovnice (a ich ľubovoľné riešenie \mathbf{b}). Zo vzťahu 3.7 vyplýva, že pre každý lineárny nestranný odhad $\mathbf{v}'\mathbf{Y}$ platí $\mathbf{v}'\mathbf{X}\boldsymbol{\beta} = \mathbb{E}(\mathbf{v}'\mathbf{Y}) = \mathbf{t}'\boldsymbol{\beta}$, z čoho vyplýva $(\mathbf{v}'\mathbf{X} - \mathbf{t}')\boldsymbol{\beta} = 0$ pre $\forall \boldsymbol{\beta}$ (teda $\mathbf{v}'\mathbf{X} = \mathbf{t}'$). Využitím vzťahu 3.7 sa vyjadri rozptyl funkcie $\mathbf{v}'\mathbf{Y}$ nasledovne

$$\begin{aligned} D(\mathbf{v}'\mathbf{Y}) &= D[\mathbf{v}'\mathbf{Y} - \mathbf{t}'\mathbf{b} + \mathbf{t}'\mathbf{b}] = D[(\mathbf{v}'\mathbf{Y} - \mathbf{g}'\mathbf{X}'\mathbf{Y}) + \mathbf{g}'\mathbf{X}'\mathbf{Y}] = \\ &= D(\mathbf{v}'\mathbf{Y} - \mathbf{g}'\mathbf{X}'\mathbf{Y}) + D(\mathbf{g}'\mathbf{X}'\mathbf{Y}) + 2\text{Cov}(\mathbf{v}'\mathbf{Y} - \mathbf{g}'\mathbf{X}'\mathbf{Y}, \mathbf{g}'\mathbf{X}'\mathbf{Y}) = \\ &= D(\mathbf{v}'\mathbf{Y} - \mathbf{g}'\mathbf{X}'\mathbf{Y}) + D(\mathbf{g}'\mathbf{X}'\mathbf{Y}) = D(\mathbf{v}'\mathbf{Y} - \mathbf{g}'\mathbf{X}'\mathbf{Y}) + D(\mathbf{t}'\mathbf{b}), \end{aligned}$$

keďže platí rovnosť

$$\text{Cov}(\mathbf{v}'\mathbf{Y} - \mathbf{g}'\mathbf{X}'\mathbf{Y}, \mathbf{g}'\mathbf{X}'\mathbf{Y}) = (\mathbf{v}' - \mathbf{g}'\mathbf{X}')(\mathbf{I}\sigma^2)\mathbf{X}\mathbf{g} = \sigma^2(\mathbf{v}'\mathbf{X} - \mathbf{g}'\mathbf{X}'\mathbf{X})\mathbf{g} = \sigma^2(\mathbf{t}' - \mathbf{t}')\mathbf{g} = 0.$$

Aby odhad $\mathbf{t}'\mathbf{b}$ bol najlepším nestranným lineárnym odhadom, musí platiť, že hodnota $D(\mathbf{v}'\mathbf{Y}) - D(\mathbf{t}'\mathbf{b}) = D(\mathbf{v}'\mathbf{Y} - \mathbf{g}'\mathbf{X}'\mathbf{Y})$ je nezáporná. Je zrejmé, že platí

$$D(\mathbf{v}'\mathbf{Y} - \mathbf{g}'\mathbf{X}'\mathbf{Y}) = (\mathbf{v}' - \boldsymbol{\xi}'\mathbf{X}')\text{Var}(\mathbf{Y})(\mathbf{v} - \mathbf{X}\mathbf{g}) = \sigma^2(\mathbf{v} - \boldsymbol{\xi}\mathbf{X})'(\mathbf{v} - \mathbf{X}\mathbf{g}),$$

a pretože je táto kvadratická forma tvaru $\sigma^2\mathbf{w}'\mathbf{w}$ (kde $\mathbf{w} = (\mathbf{v} - \mathbf{X}\mathbf{g})$) je pre $\mathbf{w} \neq \mathbf{0}$ vždy kladná. Preto je odhad $\mathbf{t}'\mathbf{b}$ najlepším z iných nestranných lineárnych odhadov. \square

Je dôležité si uvedomiť, že v prípade modelu neúplnej hodnosti, je jednoznačný odhad metódou najmenších štvorcov iba pre strednú hodnotu v bode $\mathbf{t} \in \mathcal{M}(\mathbf{X})$ (teda $\mathbf{t}'\mathbf{b}$ je odhadom $\mathbf{t}'\boldsymbol{\beta}$). Nech je uvažovaný lineárny model úplnej hodnosti s regresnou maticou \mathbf{X}_1 a druhý lineárny model neúplnej hodnosti s regresnou maticou $\mathbf{X}_2 = (\mathbf{X}_1, \mathbf{x}_3)$, kde pre vektor \mathbf{x}_3 platí $\mathbf{x}_3 \in \mathcal{M}(\mathbf{X}_1)$. Oba modely generujú rovnaké regresné priestory $\mathcal{M}(\mathbf{X}_1) = \mathcal{M}(\mathbf{X}_2)$ a preto majú aj ortogonálny priemet na tento priestor totožný (ilustruje obrázok 4). Ako bolo už uvedené, v modeli neúplnej hodnosti nie je možné jednoznačne vyjadriť tento priemet (odhadmi regresných koeficientov) ako je tomu v modeli úplnej hodnosti.

Veta 3.5. Nech $\mathbf{t}'\mathbf{b}$ a $\mathbf{t}'_1\mathbf{b}$ sú odhady metódou najmenších štvorcov pre odhadnuteľné funkcie $\mathbf{t}'\boldsymbol{\beta}$ a $\mathbf{t}'_1\boldsymbol{\beta}$. Potom platí

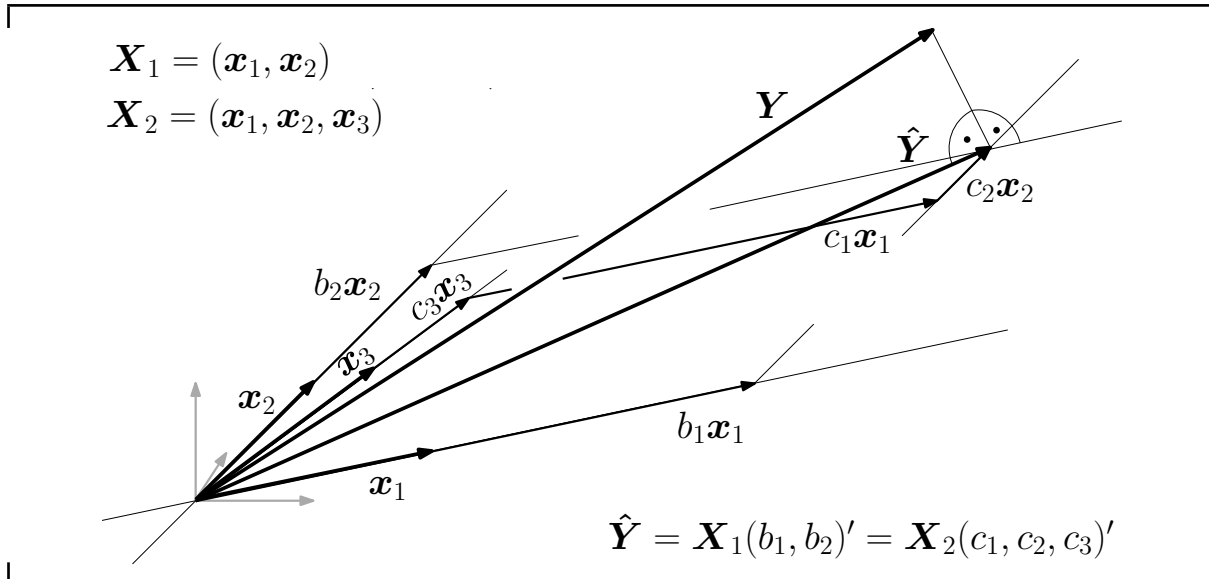
$$D(\mathbf{t}'\mathbf{b}) = \sigma^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-}\mathbf{t}, \quad \text{Cov}(\mathbf{t}'\mathbf{b}, \mathbf{t}'_1\mathbf{b}) = \sigma^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-}\mathbf{t}_1, \quad (3.8)$$

kde matica $(\mathbf{X}'\mathbf{X})^{-}$ je ľubovoľná pseudoinvertovaná matica k $\mathbf{X}'\mathbf{X}$.

Dôkaz. Z odhadnuteľnosti parametrickej funkcie $\mathbf{t}'\boldsymbol{\beta}$ vyplýva $\mathbf{t} \in \mathcal{M}(\mathbf{X}'\mathbf{X})$. Preto existuje \mathbf{g} také že platí $\mathbf{t} = \mathbf{X}'\mathbf{X}\mathbf{g}$. Nech je uvažovaná pseudoinvertovaná matica $(\mathbf{X}'\mathbf{X})^{-}$ k matici $\mathbf{X}'\mathbf{X}$. Potom z definície pseudoinvertovanej matice platí rovnosť $\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{X} = \mathbf{X}'\mathbf{X}$. Transponovaním tohto vzťahu sa vyjadri rovnosť

$$\mathbf{X}'\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-}]'\mathbf{X}'\mathbf{X} = \mathbf{X}'\mathbf{X}. \quad (3.9)$$

3.1. ODHADNUTEĽNÉ PARAMETRICKÉ FUNKCIE



Obrázok 4: Obrázok ilustruje prípad, kedy regresné matice modelu úplnej a neúplnej hodnosti vytvárajú totožné priestory $\mathcal{M}(\mathbf{X}_1) = \mathcal{M}(\mathbf{X}_2)$.

Použitím poslednej rovnosti 3.9 a vzťahu $\mathbf{t} = \mathbf{X}'\mathbf{X}\mathbf{g}$ platí

$$\begin{aligned} D(\mathbf{t}'\mathbf{b}) &= D(\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}) = \mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\text{Var}(\mathbf{Y})\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}]'\mathbf{t} = \\ &= \sigma^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}]'\mathbf{t} = \\ &= \sigma^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}]'\mathbf{X}'\mathbf{X}\mathbf{g} = \\ &= \sigma^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{g} = \sigma^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{t}. \end{aligned} \quad (3.10)$$

Podľa vzťahu 3.7 odhadu $\mathbf{t}'\mathbf{b}$, odhadnuteľnej funkcie $\mathbf{t}'\boldsymbol{\beta}$, metódou najmenších štvorcov platí vzťah $\mathbf{t}'\mathbf{b} = \mathbf{g}'\mathbf{X}'\mathbf{Y}$. Obdobne je možné získať vyjadrenie odhadu pre odhadnuteľnú funkciu $\mathbf{t}'_1\boldsymbol{\beta}$, kde $\mathbf{t}_1 = \mathbf{X}'_1\mathbf{X}\mathbf{g}_1$. Potom pre kovarianciu odhadov $\mathbf{t}'\mathbf{b}$ a $\mathbf{t}'_1\mathbf{b}$ platí

$$\begin{aligned} \text{Cov}(\mathbf{t}'\mathbf{b}, \mathbf{t}'_1\mathbf{b}) &= \text{Cov}(\mathbf{g}'\mathbf{X}'\mathbf{Y}, \mathbf{g}'_1\mathbf{X}'\mathbf{Y}) = \mathbf{g}'\mathbf{X}'\text{Var}(\mathbf{Y})\mathbf{X}\mathbf{g}_1 = \sigma^2\mathbf{g}'\mathbf{X}'\mathbf{X}\mathbf{g}_1 = \\ &= \sigma^2\mathbf{g}'\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{g}_1 = \sigma^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{t}_1. \end{aligned} \quad (3.11)$$

□

V predchádzajúcom dôkaze je možné si povšimnúť, že rozptyl odhad $\mathbf{t}'\mathbf{b}$ by bolo možné vyjadriť bez použitia pseudoinverznej matice vo vzťahu 3.10, takto

$$D(\mathbf{t}'\mathbf{b}) = \sigma^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{t} = \sigma^2\mathbf{g}'\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{g} = \mathbf{g}'(\mathbf{X}'\mathbf{X})\mathbf{g} = \mathbf{g}'\mathbf{t}.$$

Obdobne v odvodzovaní 3.11, bola kovariancia vyjadrená bez použitia pseudoinverznej matice nasledovne

$$\text{Cov}(\mathbf{t}'\mathbf{b}, \mathbf{t}'_1\mathbf{b}) = \sigma^2\mathbf{g}'\mathbf{X}'\mathbf{X}\mathbf{g}_1 = \sigma^2\mathbf{g}'\mathbf{t}_1 = \sigma^2\mathbf{t}'\mathbf{g}_1.$$

Nech je uvažovaných m vektorov $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_m$, ktoré sú z priestoru $\mathcal{M}(\mathbf{X}')$. Potom sú lineárne parametrické funkcie $\mathbf{t}'_1\boldsymbol{\beta}, \mathbf{t}'_2\boldsymbol{\beta}, \dots, \mathbf{t}'_m\boldsymbol{\beta}$ odhadnuteľné a je možné ich zapísať maticovo $\mathbf{T}'\boldsymbol{\beta}$, kde $\mathbf{T} = (\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_m)$. Nech je takto zvolená matica \mathbf{T} , potom pre vektor $\mathbf{T}'\boldsymbol{\beta}$ je vhodný názov - vektor odhadnuteľných parametrických funkcií. Ako príklad

3. LINEÁRNY REGRESNÝ MODEL NEÚPLNEJ HODNOSTI

je možné uviesť samotnú maticu \mathbf{X}' . Každý stĺpec matice \mathbf{X}' je z priestoru $\mathcal{M}(\mathbf{X}')$ a preto vektor parametrických funkcií $\mathbf{X}\boldsymbol{\beta}$ je vždy odhadnuteľný. Preto najlepším odhadom vektoru \mathbf{Y} metódou najmenších štvorcov je

$$\hat{\mathbf{Y}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}. \quad (3.12)$$

Opäť je možné označiť maticou $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'$ a uvažovať, že je nejakou transformáciou (zohľadňujúcou metódu najmenších štvorcov pri vektore \mathbf{Y}). Anděl vo vete 15. v literatúre ([2], str. 69) uviedol, že platí vzťah

$$\mathbf{A}(\mathbf{A}'\mathbf{A})^{-}\mathbf{A}'\mathbf{A} = \mathbf{A}, \quad (3.13)$$

kde teda matica $(\mathbf{A}'\mathbf{A})^{-}\mathbf{A}'$ je pseudoinverzná k matici \mathbf{A} . Ďalej profesor uvádza v tejto vete, že matica

$$\mathbf{A}(\mathbf{A}'\mathbf{A})^{-}\mathbf{A}' \quad (3.14)$$

je symetrická a táto vlastnosť nezávisí na zvolenej pseudoinverzii $(\mathbf{A}'\mathbf{A})^{-}$. Preto je matica $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'$ symetrická. A platí pre túto maticu aj idempotentnosť

$$\mathbf{H}\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}' = \mathbf{H},$$

použitím vzťahu 3.13 uvedeného vyššie. A keďže matica \mathbf{H} je symetrická aj idempotentná, podľa vety 1.8 je matica \mathbf{H} operátorom ortogonálnej projekcie na regresný priestor $\mathcal{M}(\mathbf{X})$. Z idempotentnosti matice \mathbf{H} vyplýva, že $\mathbf{I} - \mathbf{H}$ je tiež idempotentná. A keďže matica \mathbf{H} je symetrická, je zrejmé že, aj matica $\mathbf{I} - \mathbf{H}$ je symetrická. Preto matica $\mathbf{I} - \mathbf{H}$ je projekčnou maticou, podľa vety 1.8, ale na reziduálny priestor $\mathcal{M}(\mathbf{X})^{\perp}$ (podľa vety 1.7). Táto projekčná matica tak, ako v predchádzajúcej kapitole, bude označovaná $\mathbf{M} = \mathbf{I} - \mathbf{H}$.

Z vyššie uvedeného, náhodnú veličinu RSS , je možné vyjadriť nasledovnými vzťahmi

$$RSS = (\mathbf{Y} - \hat{\mathbf{Y}})'(\mathbf{Y} - \hat{\mathbf{Y}}) = \mathbf{Y}'(\mathbf{I} - \mathbf{H})'(\mathbf{I} - \mathbf{H})\mathbf{Y} = \mathbf{Y}'\mathbf{M}\mathbf{Y} = \mathbf{u}'\mathbf{u}.$$

Veta 3.6. $E(RSS) = (n - r)\sigma^2$, takže nestranným odhadom parametra σ^2 je

$$s^2 = \frac{RSS}{(n - r)},$$

kde r označuje hodnotu matice \mathbf{X} .

Dôkaz. Odvodenie je podobné ako v dôkaze vety 2.5. Ak platí $h(\mathbf{X}) = r$, potom platí $h(\mathbf{H}) = r$. Podľa vety 1.5. pre idempotentnú maticu \mathbf{M} platí

$$h(\mathbf{M}) = \text{Tr}(\mathbf{M}) = \text{Tr}(\mathbf{I} - \mathbf{H}) = \text{Tr}(\mathbf{I}) - \text{Tr}(\mathbf{H}) = n - h(\mathbf{H}) = n - r.$$

Zo vzťahu 2.15 vyplýva $E(\mathbf{Y}\mathbf{Y}') = \text{Var}(\mathbf{Y}) + \mathbf{X}\boldsymbol{\beta}(\mathbf{X}\boldsymbol{\beta})'$, aj pre model neúplnej hodnosti. Táto rovnosť bude použitá v nasledujúcom vyjadrovaní strednej hodnoty RSS

$$\begin{aligned} E(RSS) &= E(\mathbf{Y}'\mathbf{M}\mathbf{Y}) = E[\text{Tr}(\mathbf{M}\mathbf{Y}\mathbf{Y}')] = \text{Tr}[\mathbf{M}E(\mathbf{Y}\mathbf{Y}')] = \\ &= \text{Tr}[\mathbf{M}(\text{Var}(\mathbf{Y}) + \mathbf{X}\boldsymbol{\beta}(\mathbf{X}\boldsymbol{\beta})')] = \text{Tr}[\mathbf{M}(\mathbf{I}\sigma^2 + \mathbf{X}\boldsymbol{\beta}(\mathbf{X}\boldsymbol{\beta})')] = \\ &= \sigma^2\text{Tr}(\mathbf{M}) + \text{Tr}[\mathbf{M}\mathbf{X}\boldsymbol{\beta}(\mathbf{X}\boldsymbol{\beta})'] = (n - r)\sigma^2 + (\mathbf{X}\boldsymbol{\beta})'\mathbf{M}\mathbf{X}\boldsymbol{\beta} = \\ &= (n - r)\sigma^2, \end{aligned}$$

kde boli použité vzťahy 1.6 a 2.15. Posledná rovnosť vyplýva z platnosti $\mathbf{M}\mathbf{X} = \mathbf{0}$. \square

3.2. Normálny lineárny regresný model

Nech je uvažovaný $\mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$ a matica \mathbf{T} typu $k \times p$, ktorej stĺpce sú z priestoru $\mathcal{M}(\mathbf{X}')$ a splňa $h(\mathbf{T}) = p$. Teda nebudú zbytočne uvažované nadbytočné zložky vektora parametrických funkcií $\mathbf{T}'\boldsymbol{\beta}$, pretože parametrické funkcie $t'_i\boldsymbol{\beta}$ sú lineárne a ich kombináciami je možné získať ľubovoľný odhad funkcie $\mathbf{t}'\boldsymbol{\beta}$, kde $\mathbf{t} \in \mathcal{M}(\mathbf{T})$. Keďže každý stĺpce matice \mathbf{T} je z priestoru $\mathcal{M}(\mathbf{X}')$, vektor parametrických funkcií $\mathbf{T}'\boldsymbol{\beta}$ je odhadnuteľný. Z odhadnuteľnosti vyplýva, že existuje matica $\mathbf{J}_{k \times p}$ taká, že platí

$$\mathbf{T} = \mathbf{X}'\mathbf{X}\mathbf{J}. \quad (3.15)$$

Aby platilo $h(\mathbf{T}) = p$, vo vzťahu 3.15 musí platiť $h(\mathbf{J}) = p$ a $h(\mathbf{X}\mathbf{J}) = p$. Bodovým odhadom $\mathbf{T}'\boldsymbol{\beta}$ metódou najmenších štvorcov je vektor

$$\mathbf{T}'\mathbf{b} = \mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}, \quad (3.16)$$

pretože homogénne riešenie normálnych rovníc je ortogonálne na priestor $\mathcal{M}(\mathbf{X}')$ a teda aj na $\mathcal{M}(\mathbf{T})$. Použitím vzťahu 3.15 sa ukáže nestrannosť odhadu $\mathbf{T}'\mathbf{b}$ nasledovne

$$\begin{aligned} \mathbb{E}(\mathbf{T}'\mathbf{b}) &= \mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbb{E}(\mathbf{Y}) = \mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \\ &= \mathbf{J}'\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{J}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{T}'\boldsymbol{\beta}. \end{aligned} \quad (3.17)$$

Variančná matica $\text{Var}(\mathbf{T}'\mathbf{b})$ sa vyjadří nasledovne

$$\begin{aligned} \text{Var}(\mathbf{T}'\mathbf{b}) &= \text{Var}(\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}) = \sigma^2\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}]'\mathbf{T} = \\ &= \sigma^2\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}]'\mathbf{X}'\mathbf{X}\mathbf{J} = \\ &= \sigma^2\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{J} = \sigma^2\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{T}, \end{aligned} \quad (3.18)$$

kde v predposlednej rovnosti bol použitý vzťah 3.9. Vzťah 3.16 určuje pre odhad $\mathbf{T}'\mathbf{b}$ lineárnu transformáciu z normálne rozdeleného vektora \mathbf{Y} . Potom p -rozmerný vektor odhadov $\mathbf{T}'\mathbf{b}$ má tiež normálne rozdelenie

$$\mathbf{T}'\mathbf{b} \sim \mathcal{N}_p(\mathbf{T}'\boldsymbol{\beta}, \sigma^2\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{T}). \quad (3.19)$$

Z odvodzovania predchádzajúcich vzťahov vyplýva pre samotný odhad \mathbf{b} rozdelenie

$$\mathbf{b} \sim \mathcal{N}_k((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}]'). \quad (3.20)$$

Variančnú maticu $\text{Var}(\mathbf{T}'\mathbf{b})$ je možné vyjadriť aj bez pseudoinverznej matice $(\mathbf{X}'\mathbf{X})^{-1}$ nasledovne

$$\text{Var}(\mathbf{T}'\mathbf{b}) = \sigma^2\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{T} = \sigma^2\mathbf{J}'\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{J} = \sigma^2\mathbf{J}'\mathbf{X}'\mathbf{X}\mathbf{J}. \quad (3.21)$$

Veta 3.7. *Odhad $\mathbf{T}'\mathbf{b}$ a hodnota RSS sú navzájom nezávislé.*

Dôkaz. Platí $RSS = \mathbf{Y}'\mathbf{M}\mathbf{Y}$ a $\mathbf{T}'\mathbf{b} = \mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$. Myšlienka dôkazu je rovnaká ako pri dôkaze vety 2.10. a v tomto prípade stačí ukázať

$$\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\sigma^2\mathbf{I})\mathbf{M} = \mathbf{O}.$$

Táto rovnosť sa ukáže nasledovne

$$\begin{aligned} \mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\sigma^2\mathbf{I})\mathbf{M} &= \sigma^2\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') = \\ &= \sigma^2\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - \sigma^2\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \\ &= \sigma^2\mathbf{J}'\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - \sigma^2\mathbf{J}'\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \\ &= \sigma^2\mathbf{J}'\mathbf{X}'(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') = \mathbf{O}, \end{aligned}$$

kde v poslednej rovnosti bol využitý vzťah 3.5. □

Veta 3.8. Platí $RSS/\sigma^2 \sim \chi^2(n-r)$, kde r označuje hodnotu matice \mathbf{X} .

Dôkaz. Postupuje sa rovnako ako v dôkaze vety 2.9. Priestor $\mathcal{M}(\mathbf{X})^\perp$ má dimenziou $n-r$ a preto je možné uvažovať maticu $\mathbf{N}_{n \times n-r}$ splňajúcu $\mathbf{N}\mathbf{N}' = \mathbf{M}$ a $\mathbf{N}'\mathbf{N} = \mathbf{I}_{n-r}$. Potom je zrejmé, že vektor $\mathbf{u} = \mathbf{M}\mathbf{Y}$ má rovnaké rozdelenie ako $\sigma\mathbf{N}\mathbf{Z} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2\mathbf{M})$, kde vektor $\mathbf{Z} \sim \mathbf{N}_{n-r}(\mathbf{0}, \mathbf{I})$. Z toho vyplýva

$$\frac{RSS}{\sigma^2} = \frac{\mathbf{u}'\mathbf{u}}{\sigma^2} \text{ má rovnaké rozdelenie, ako } \mathbf{Z}'\mathbf{N}'\mathbf{N}\mathbf{Z} = \mathbf{Z}'\mathbf{Z} \sim \chi^2(n-r).$$

□

3.2.1. Testy hypotéz a intervaly spoľahlivosti

Najprv bude ukázané testovanie na jednej parametrickej funkcii. Nech je teda daná odhadnuteľná parametrická funkcia $\mathbf{t}'\boldsymbol{\beta}$. Jej bodový odhad metódou najmenších štvorcov bude označený $\mathbf{t}'\mathbf{b}$.

Veta 3.9. Nech $h(\mathbf{X}) = r$ a $\mathbf{t}'\boldsymbol{\beta}$ je odhadnuteľná parametrická funkcia. Potom platí

$$\frac{(\mathbf{t}'\mathbf{b} - \mathbf{t}'\boldsymbol{\beta})}{\sqrt{s^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{t}}} \sim \mathbf{t}(n-r).$$

Dôkaz. Zo vzťahu 3.19 vyplýva $(\mathbf{t}'\mathbf{b} - \mathbf{t}'\boldsymbol{\beta})/\sqrt{\sigma^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{t}} \sim \mathbf{N}(0, 1)$. Z vety 3.8. je známa platnosť $RSS/\sigma^2 \sim \chi^2(n-r)$. A z vety 3.7. nepriamo vyplýva, že hodnota $\mathbf{t}'\mathbf{b}$ a odhad s^2 sú navzájom nezávislé veličiny. Potom platí

$$\frac{(\mathbf{t}'\mathbf{b} - \mathbf{t}'\boldsymbol{\beta})/\sqrt{\sigma^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{t}}}{\sqrt{RSS\sigma^{-2}/(n-r)}} = \frac{(\mathbf{t}'\mathbf{b} - \mathbf{t}'\boldsymbol{\beta})}{\sqrt{s^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{t}}} \sim \mathbf{t}(n-r)$$

kde bol použitý vzťah pre odhad $s^2 = RSS/(n-r)$ neznámeho parametra σ^2 . □

Vetu 3.9 je možné výhodne využiť pri testovaní hypotéz. Nech pre dané číslo $\theta \in \mathbb{R}$, je uvažovaná nulová hypotéza $\mathbf{H}_0 : \mathbf{t}'\boldsymbol{\beta} = \theta$ proti alternatívnej hypotéze $\mathbf{H}_1 : \mathbf{t}'\boldsymbol{\beta} \neq \theta$. Potom \mathbf{H}_0 bude zamietnutá na hladine významnosti α , práve keď platí

$$\frac{(\mathbf{t}'\mathbf{b} - \theta)}{\sqrt{s^2\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{t}}} \geq \mathbf{t}_{1-\frac{\alpha}{2}}(n-r).$$

Obojstranný interval spoľahlivosti pre odhad parametra $\mathbf{t}'\boldsymbol{\beta}$ so spoľahlivosťou $1 - \alpha$ je

$$\left(\mathbf{t}'\mathbf{b} - s\sqrt{\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{t}} \mathbf{t}_{1-\frac{\alpha}{2}}(n-r), \mathbf{t}'\mathbf{b} + s\sqrt{\mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{t}} \mathbf{t}_{1-\frac{\alpha}{2}}(n-r) \right). \quad (3.22)$$

V prípade, že je potrebné testovať viacero odhadnuteľných parametrických funkcií súčasne, je vhodnejšie používať vektorový tvar $\mathbf{T}'\boldsymbol{\beta}$, kde matica $\mathbf{T}_{k \times p}$ bola presne definovaná v odstavci 3.2. Bodovým odhadom získaným metódou najmenších štvorcov vektora parametrických funkcií $\mathbf{T}'\boldsymbol{\beta}$ je vektor $\mathbf{T}'\mathbf{b}$, ktorý má podľa vzťahu 3.19 normálne rozdelenie $\mathbf{N}_p(\mathbf{T}'\boldsymbol{\beta}, \sigma^2\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{T})$.

3.2. NORMÁLNÝ LINEÁRNY REGRESNÝ MODEL

Veta 3.10. *Platí*

$$\frac{1}{ps^2}(\mathbf{T}'\mathbf{b} - \mathbf{T}'\boldsymbol{\beta})'(\mathbf{T}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{T})^{-1}(\mathbf{T}'\mathbf{b} - \mathbf{T}'\boldsymbol{\beta}) \sim F(p, n - r).$$

Dôkaz. Podľa vzťahu 3.21 platí $\text{Var}(\mathbf{T}'\mathbf{b}) = \sigma^2\mathbf{J}'\mathbf{X}'\mathbf{X}\mathbf{J}$, kde na maticu $\mathbf{J}_{k \times p}$ bola položená podmienka $h(\mathbf{X}\mathbf{J}) = p$ (kde $p \leq r$). Variančná matica $\text{Var}(\mathbf{T}'\mathbf{b})$ je potom regulárna (s hodnotou p) pozitívne definitná. Použitím vety 1.4 a kladnosti vlastných čísel matice $\mathbf{J}'\mathbf{X}'\mathbf{X}\mathbf{J}$, je možné urobiť rozklad

$$\mathbf{J}'\mathbf{X}'\mathbf{X}\mathbf{J} = \mathbf{B}\mathbf{B}',$$

kde matica \mathbf{B} je regulárna a jej i -tý stĺpce je $\sqrt{\lambda_i}$ -násobkom vlastného vektora \mathbf{v}_i pôvodnej matice. Potom náhodný vektor $\mathbf{T}'\mathbf{b}$ má rovnaké rozdelenie ako vektor

$$\mathbf{T}'\boldsymbol{\beta} + \sigma\mathbf{B}\mathbf{Z} \sim N_p(\mathbf{T}'\boldsymbol{\beta}, \sigma^2\mathbf{J}'\mathbf{X}'\mathbf{X}\mathbf{J}), \quad (3.23)$$

kde vektor \mathbf{Z} má normálne rozdelenie $N_p(\mathbf{0}, \mathbf{I})$. Pri následnom dosadení náhodného vektora 3.23 do vzťahu

$$\begin{aligned} \frac{1}{\sigma^2}(\mathbf{T}'\mathbf{b} - \mathbf{T}'\boldsymbol{\beta})'(\mathbf{B}\mathbf{B}')^{-1}(\mathbf{T}'\mathbf{b} - \mathbf{T}'\boldsymbol{\beta}) &= \frac{1}{\sigma^2}(\sigma\mathbf{B}\mathbf{Z})'(\mathbf{B}\mathbf{B}')^{-1}(\mathbf{B}\mathbf{Z}\sigma) = \\ &= \mathbf{Z}'\mathbf{B}'(\mathbf{B}\mathbf{B}')^{-1}\mathbf{B}\mathbf{Z} = \mathbf{Z}'\mathbf{I}\mathbf{Z} \sim \chi^2(p), \end{aligned}$$

kde pri odvodzovaní bola využitá rovnosť $\mathbf{B}'(\mathbf{B}\mathbf{B}')^{-1}\mathbf{B} = \mathbf{I}$. Táto rovnosť sa ukáže platnou, podobne ako v dôkaze vety 2.12. vynásobením zľava (regulárnou) maticou \mathbf{B} (resp. sprava maticou \mathbf{B}'). Pomocou výsledkov viet 3.6. a 3.8. je možné odvodiť nasledujúci vzťah

$$\frac{\frac{1}{ps^2}(\mathbf{T}'\mathbf{b} - \mathbf{T}'\boldsymbol{\beta})'(\mathbf{B}\mathbf{B}')^{-1}(\mathbf{T}'\mathbf{b} - \mathbf{T}'\boldsymbol{\beta})}{RSS\sigma^{-2}/(n - r)} = \frac{(\mathbf{T}'\mathbf{b} - \mathbf{T}'\boldsymbol{\beta})'(\mathbf{B}\mathbf{B}')^{-1}(\mathbf{T}'\mathbf{b} - \mathbf{T}'\boldsymbol{\beta})}{ps^2} \sim F(p, n - r),$$

ktorý korešponduje s tvrdením tejto vety. Ľavá strana predchádzajúceho vzťahu odpovedá zadanému Fisherovmu rozdeleniu $\frac{\chi^2(p)/p}{\chi^2(n-r)/(n-r)}$. \square

Podľa vety 3.10. je možné testovať nulovú hypotézu $H_0 : \mathbf{T}'\boldsymbol{\beta} = \boldsymbol{\theta}$ proti alternatívnej hypotéze $H_1 : \mathbf{T}'\boldsymbol{\beta} \neq \boldsymbol{\theta}$, kde $\boldsymbol{\theta}$ je p -prvkový vektor. Nulová hypotéza H_0 bude zamietnutá na hladine α , práve vtedy, keď hodnota výrazu uvedeného vo vete 3.10. prekročí kvantil $F_{1-\alpha}(p, n - k)$.

4. Výber modelu a diagnostika

Pri výbere "rozumného" modelu by mali byť regresory vybrané do modelu tak, aby čo najlepšie popisovali strednú hodnotu náhodného vektora \mathbf{Y} (bez újmy na všeobecnosti, nech je uvažovaný lineárny regresný model neúplnej hodnosti 3.1). Napríklad v prípade, že regresná matica \mathbf{X} ($h(\mathbf{X}) = k$) má príliš veľa stĺpcov a niektoré sa len nevýznamne podieľajú na vyjadrení strednej hodnoty $\mathbf{E}(\mathbf{Y})$. Preto je dôležitý správny výber týchto regresorov. Teda je potrebné uvažovať menší model (takzvaný *podmodel*) s regresnou maticou \mathbf{X}^* s hodnotou k^* , pre ktorú platí

$$\mathcal{M}(\mathbf{X}^*) \subseteq \mathcal{M}(\mathbf{X}) \quad \text{a zároveň} \quad 0 < k^* < k, \quad (4.1)$$

kde $k = h(\mathbf{X})$ a $k^* = h(\mathbf{X}^*)$. Často je takýto podmodel tvorený len istými stĺpcami z \mathbf{X} . V prípade opačnej situácie sa hovorí o *nadmodeli*. Nech je teda uvažovaný model 3.1 a matica \mathbf{X}^* spĺňajúca 4.1. Potom sa povie, že *platí podmodel* modelu 3.1, ak pre nejaký vektor regresných koeficientov $\boldsymbol{\beta}^*$ platí

$$\mathbf{E}(\mathbf{Y}) = \mathbf{X}^* \boldsymbol{\beta}^*.$$

Nech sú označené všetky veličiny vzťahujúce sa na podmodel symbolom $*$ v hornom indexe. Zvára v literatúre ([13], str. 30) uviedol nasledujúce matematické tvrdenie: Ak platí v modeli podmodel a \mathbf{Y} má normálne rozdelenie, potom štatistiky $\hat{\mathbf{Y}}^*$ a $\hat{\mathbf{u}}^*$ sú nezávislé a platí

$$\frac{(RSS^* - RSS)/(k - k^*)}{RSS/(n - k)} \sim F(k - k^*, n - k). \quad (4.2)$$

Dôsledok neúplného výberu regresorov

Nasledujúce pojednanie bude o situácii, kedy budú vytvorené odhady pomocou podmodelu pričom skutočný model bude širší. Nech je teda uvažovaný lineárny regresný model neúplnej hodnosti (podľa 3.1)

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad (4.3)$$

kde $\mathbf{X}_{n \times k}$ je matica daných reálnych čísel hodnosti p a $\boldsymbol{\beta}$ vektor neznámych parametrov. A zároveň platí pre náhodný vektor $\mathbf{E}(\mathbf{e}) = \mathbf{0}$ a $\text{Var}(\mathbf{e}) = \sigma^2 \mathbf{I}$. Avšak skutočný model sa bude riadiť rozšíreným lineárnym modelom

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{e} \quad (4.4)$$

tak, že $\mathbf{Z}_{n \times m}$ je matica daných reálnych čísel a $\boldsymbol{\gamma}$ vektor m -neznámych parametrov. A zároveň pre náhodný vektor tohto rozšíreného modelu platí $\mathbf{E}(\mathbf{e}) = \mathbf{0}$ a $\text{Var}(\mathbf{e}) = \sigma^2 \mathbf{I}$. V ďalšom texte bude regresná matica rozšíreného modelu označená $\mathbf{G}_{n \times (k+m)} = (\mathbf{X}, \mathbf{Z})$ a rozšírený vektor regresných koeficientov $\boldsymbol{\alpha} = (\boldsymbol{\beta}', \boldsymbol{\gamma}')'$. Potom je možné formálne prepísať model na

$$\mathbf{Y} = \mathbf{G}\boldsymbol{\alpha} + \mathbf{e}. \quad (4.5)$$

Aby skutočný rozšírený model rozšíril priestor strednej hodnoty $\mathbf{E}(\mathbf{Y})$, na regresnú maticu \mathbf{G} je kladená podmienka $h(\mathbf{G}) > p$. Samozrejme tiež platí $h(\mathbf{G}) < n$, pretože model 4.5 je lineárnym regresným modelom.

Ďalej budú označené indexom G všetky potrebné štatistiky vzťahujúce sa na rozšírený

model 4.5. \mathbf{H}_G bude označovať projekčnú maticu na rozšírený regresný priestor $\mathcal{M}(\mathbf{G})$ a matica $\mathbf{M}_G = \mathbf{I} - \mathbf{H}_G$ bude označovať projekčnú maticu na (zúžený) reziduálny priestor $\mathcal{M}(\mathbf{G})^\perp$. Potom zrejme platí pre najlepší lineárny nestranný odhad vektora \mathbf{Y} vzťah $\hat{\mathbf{Y}}_G = \mathbf{H}_G \mathbf{Y}$ a pre odhad chyby $\mathbf{u}_G = \mathbf{M}_G \mathbf{Y} = \mathbf{Y} - \hat{\mathbf{Y}}_G$. Pre reziduálny súčet štvorcov rozšíreného modelu $RSS_G = \|\mathbf{u}_G\|^2$. Podľa vzťahu 3.12 platí

$$\hat{\mathbf{Y}}_G = \mathbf{G}(\mathbf{G}'\mathbf{G})^{-1}\mathbf{G}'\mathbf{Y} = \mathbf{X}\mathbf{b}_G + \mathbf{Z}\mathbf{c}_G, \quad (4.6)$$

kde \mathbf{b}_G a \mathbf{c}_G sú odhady metódou najmenších štvorcov pre vektory $\boldsymbol{\beta}$ a $\boldsymbol{\gamma}$. Avšak ak je potrebné vidieť rozdiel v odhadoch regresných koeficientov medzi uvažovaným 4.3 a skutočným rozšíreným modelom, je potrebné aby odhady vektorov regresných koeficientov $\boldsymbol{\beta}$ a $\boldsymbol{\gamma}$ boli nezávislé. Platí $\mathcal{M}(\mathbf{X}, \mathbf{Z}) = \mathcal{M}(\mathbf{X}, \mathbf{MZ})$, pretože \mathbf{M} je projekčná matica do reziduálneho priestoru $\mathcal{M}(\mathbf{X})^\perp$ a teda stĺpce matice \mathbf{Z} sa premietnu do reziduálneho priestoru ako stĺpce matice \mathbf{MZ} . Tým vznikne nezávislosť odhadov regresných koeficientov. Miesto matice \mathbf{G} vo vzťahu 4.6 sa použije matica $(\mathbf{X}, \mathbf{MZ})$, ktorá (ako už bolo ukázané) generuje rovnaký regresný priestor.

$$\begin{aligned} \hat{\mathbf{Y}}_G &= (\mathbf{X}, \mathbf{MZ})((\mathbf{X}, \mathbf{MZ})'(\mathbf{X}, \mathbf{MZ}))^{-1}(\mathbf{X}, \mathbf{MZ})'\mathbf{Y} \\ &= (\mathbf{X}, \mathbf{MZ}) \begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{O} \\ \mathbf{O} & \mathbf{Z}'\mathbf{MZ} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}' \\ \mathbf{Z}'\mathbf{M} \end{pmatrix} \mathbf{Y} \\ &= \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} + \mathbf{MZ}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{M}\mathbf{Y} \\ &= \hat{\mathbf{Y}} + \mathbf{MZ}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{u} \\ &= \mathbf{X}\mathbf{b} + (\mathbf{I} - \mathbf{H})\mathbf{Z}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{u} \\ &= \mathbf{X}\mathbf{b} + (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{Z}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{u} \\ &= \mathbf{X}(\mathbf{b} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{u}) + \mathbf{I}\mathbf{Z}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{u} \\ &= \mathbf{X} \underbrace{(\mathbf{b} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{u})}_{\mathbf{b}_G} + \mathbf{Z} \underbrace{(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{u}}_{\mathbf{c}_G} \end{aligned} \quad (4.7)$$

Použitím vlastnosti matice 3.14, na hore uvedený súčin blokových matíc, sa príde k záveru, že i tento súčin je nezávislý na zvolenej pseudoinverzii. Ďalej je zrejme, že k matici

$$\begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{O} \\ \mathbf{O} & \mathbf{Z}'\mathbf{MZ} \end{pmatrix}$$

je pseudoinverzná matica

$$\begin{pmatrix} (\mathbf{X}'\mathbf{X})^{-1} & \mathbf{O} \\ \mathbf{O} & (\mathbf{Z}'\mathbf{MZ})^{-1} \end{pmatrix}, \quad (4.8)$$

čo sa dá ukázať súčinom blokových matíc z definície pseudoinverznej matice. A práve použitím pseudoinverznej matice 4.8 v odvodzovaní 4.7 sa sprehľadní postup (nezávislý na zvolenej pseudoinverzii). V poslednom riadku odvodzovania 4.7 je možné si všimnúť priamy vzťah medzi odhadmi \mathbf{b} a \mathbf{b}_G . V predchádzajúcom odvodzovaní 4.7 taktiež bola ukázaná platnosť vzťahu $\hat{\mathbf{Y}}_G = \hat{\mathbf{Y}} + \mathbf{MZ}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{u}$. Vektorom \mathbf{d} bude ďalej označený rozdiel odhadov $\hat{\mathbf{Y}}_G$ a $\hat{\mathbf{Y}}$ takto

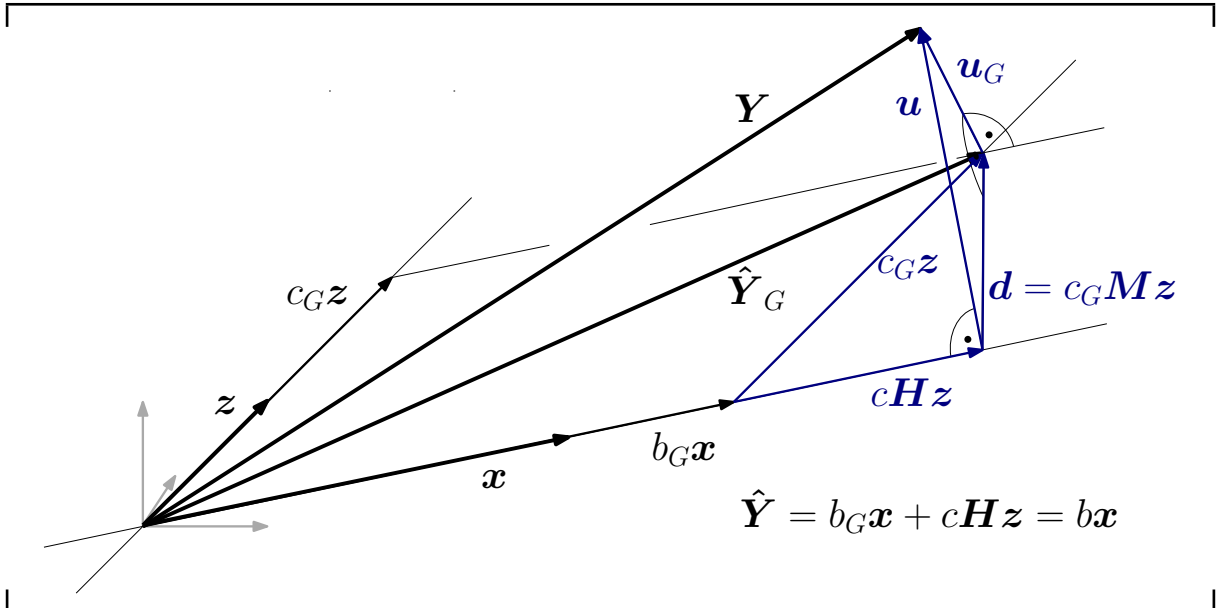
$$\mathbf{d} = \hat{\mathbf{Y}}_G - \hat{\mathbf{Y}} = \mathbf{MZ}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{u}. \quad (4.9)$$

Nech je odčítaný vektor \mathbf{Y} zo vzťahu $\hat{\mathbf{Y}}_G = \hat{\mathbf{Y}} + \mathbf{d}$, potom platí vzťah

$$(\hat{\mathbf{Y}}_G - \mathbf{Y}) = (\hat{\mathbf{Y}} - \mathbf{Y}) + \mathbf{d}$$

a následne rovnosť

$$-\mathbf{u}_G = -\mathbf{u} + \mathbf{d} \quad \text{resp.} \quad \mathbf{u} = \mathbf{u}_G + \mathbf{d}. \quad (4.10)$$



Obrázok 5: Na obrázku je zachytený prípad, kedy regresná matica \mathbf{X} je tvorená len vektorom \mathbf{x} a matica \mathbf{Z} je tvorená len vektorom \mathbf{z} . Označenie korešponduje so zavedením v tejto kapitole. Matica \mathbf{H} je projekčná matica na $\mathcal{M}(\mathbf{x})$ a \mathbf{M} je projekčná matica na $\mathcal{M}(\mathbf{x})^\perp$.

Z platností vzťahov $\mathbf{u}_G \in \mathcal{M}(\mathbf{G})^\perp$ a $\mathbf{d} \in \mathcal{M}(\mathbf{G}) \cap \mathcal{M}(\mathbf{X})^\perp$ vyplýva, že vektory \mathbf{u}_G a \mathbf{d} sú na seba kolmé (zjednodušená situácia je zobrazená na obrázku 5). Preto z kolmosti pre vzťah $\mathbf{u} = \mathbf{u}_G + \mathbf{d}$ z 4.10 platí po znormovaní a umocnení

$$\|\mathbf{u}\|^2 = \|\mathbf{u}_G + \mathbf{d}\|^2 = \|\mathbf{u}_G\|^2 + \|\mathbf{d}\|^2.$$

Z predchádzajúceho priamo vyplýva

$$RSS = RSS_G + \|\mathbf{d}\|^2. \quad (4.11)$$

Podľa vety 3.6. a platnosti modelu 4.5 platí

$$E(RSS_G) = (n - h(\mathbf{G}))\sigma^2. \quad (4.12)$$

Avšak v prípade $E(RSS)$ platí

$$E(RSS) = E\|\mathbf{M}\mathbf{Y}\|^2 = E\|\mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{e})\|^2 = E\|\mathbf{M}\mathbf{Z}\boldsymbol{\gamma} + \mathbf{M}\mathbf{e}\|^2. \quad (4.13)$$

Oba členy v norme, poslednom vyjadrení vzťahu 4.13, sa nachádzajú v navzájom ortogonálnych podpriestoroch a preto je možné ďalej pokračovať vo vyjadrovaní takto

$$E(RSS) = E\|\mathbf{M}\mathbf{Z}\boldsymbol{\gamma}\|^2 + E\|\mathbf{M}\mathbf{e}\|^2 = \|\mathbf{M}\mathbf{Z}\boldsymbol{\gamma}\|^2 + (n - h(\mathbf{X}))\sigma^2. \quad (4.14)$$

Pre odhad strednej hodnoty uvažovaného modelu $E\hat{Y}$ platí

$$E\hat{Y} = HE\hat{Y}_G = H(\mathbf{X}\beta + \mathbf{Z}\gamma) = \mathbf{X}\beta + \mathbf{H}\mathbf{Z}\gamma.$$

V prípade, že rozšírený model skutočne rozširuje pôvodný model (teda neplatí $\mathbf{H}\mathbf{Z} = \mathbf{Z}$) je odhad strednej hodnoty \hat{Y} odhadom vychýleným. Preto platí pre vychýlenie tohto odhadu nasledujúci vzťah

$$\text{bias}(\hat{Y}) = E\hat{Y} - EY = \mathbf{X}\beta + \mathbf{H}\mathbf{Z}\gamma - (\mathbf{X}\beta + \mathbf{Z}\gamma) = (\mathbf{H} - \mathbf{I})\mathbf{Z}\gamma = -\mathbf{M}\mathbf{Z}\gamma. \quad (4.15)$$

Ako uvádza Zvára v literatúre ([13], str. 83), vychýlené odhady nie sú porovnávané pomocou ich rozptylu či variančnej matice, ale pomocou *strednej štvorcovej chyby*. Táto štatistika označovaná MSE (z angl. Mean Square Error) je všeobecne definovaná pre odhad $\hat{\theta}$ parametra θ takto

$$\text{MSE}(\hat{\theta}) = E(\hat{\theta} - \theta)^2 = D(\hat{\theta}) + \text{bias}^2(\hat{\theta}). \quad (4.16)$$

a pre vektor odhadov $\hat{\theta}$ parametrov θ nasledovne

$$\text{MSE}(\hat{\theta}) = E(\hat{\theta} - \theta)(\hat{\theta} - \theta)' = \text{Var}(\hat{\theta}) + \text{bias}(\hat{\theta})\text{bias}(\hat{\theta})'. \quad (4.17)$$

Preto v prípade dvoch nestranných odhadov ($\text{bias}(\hat{\theta}) = 0$ resp. $\text{bias}(\hat{\theta}) = \mathbf{0}$) je ich porovnanie pomocou rozptylov (resp. variančných matíc) a stredných štvorcových chýb totožné. Potom strednú štvorcovú chybu odhadu \hat{Y} pomocou vzťahu 4.15 je možno rozpísať nasledovne

$$\text{MSE}(\hat{Y}) = \text{Var}(\hat{Y}) + \text{bias}(\hat{Y})(\text{bias}(\hat{Y}))' = \sigma^2\mathbf{H} + \mathbf{M}\mathbf{Z}\gamma\gamma'\mathbf{Z}'\mathbf{M}.$$

Pretože odhad \hat{Y}_G je nevychýlený (teda $\text{bias}(\hat{Y}_G) = \mathbf{0}$) platí

$$\begin{aligned} \text{MSE}(\hat{Y}_G) &= \text{Var}(\hat{Y}_G) = (\mathbf{X}, \mathbf{M}\mathbf{Z})\text{Var}((\mathbf{b}'_G, \mathbf{c}'_G)')(\mathbf{X}, \mathbf{M}\mathbf{Z})'\mathbf{Y} \\ &= \sigma^2(\mathbf{X}, \mathbf{M}\mathbf{Z}) \begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{O} \\ \mathbf{O} & \mathbf{Z}'\mathbf{M}\mathbf{Z} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}' \\ \mathbf{Z}'\mathbf{M} \end{pmatrix}, \quad (4.18) \\ &= \sigma^2\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' + \sigma^2\mathbf{M}\mathbf{Z}(\mathbf{Z}'\mathbf{M}\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{M} \\ &= \sigma^2\mathbf{H} + \sigma^2\mathbf{M}\mathbf{Z}(\mathbf{Z}'\mathbf{M}\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{M} \end{aligned}$$

kde bola použitá pseudoinverzia tvaru 4.8. Ako bolo konštatované vyššie, súčin blokových matíc v 4.18 je nezávislý na voľbe pseudoinverzie. Použitím 4.18 ďalej platí

$$\text{MSE}(\hat{Y}_G) - \text{MSE}(\hat{Y}) = \sigma^2(\mathbf{M}\mathbf{Z}(\mathbf{Z}'\mathbf{M}\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{M} - \mathbf{M}\mathbf{Z}\gamma\gamma'\mathbf{Z}'\mathbf{M}/\sigma^2). \quad (4.19)$$

Zvára v literatúre ([13], str. 228) dokázal nasledujúcu vetu: Pre maticu $\mathbf{A}_{m \times n}$ a vektor $\mathbf{c} \in \mathbb{R}^n$ platí nerovnosť $\|\mathbf{A}\mathbf{c}\|^2 \leq 1$ práve vtedy, keď je matica

$$\mathbf{A}(\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}' - \mathbf{A}\mathbf{c}\mathbf{c}'\mathbf{A}' \quad (4.20)$$

pozitívne semidefinitná.

Aplikovaním tejto vety na vzťah 4.19 (položením $\mathbf{A} = \mathbf{M}\mathbf{Z}$ a $\mathbf{c} = \gamma/\sigma$) je možné dôjsť k záveru, že matica 4.19 je pozitívne semidefinitná práve za podmienky $\|\mathbf{M}\mathbf{Z}\gamma/\sigma\|^2 \leq 1$. Túto nerovnosť je možné jednoduchou úpravou previesť na $\|\mathbf{M}\mathbf{Z}\gamma\|^2 \leq \sigma^2$ a použitím vzťahu 4.15 ďalej previesť na vzťah $\|\text{bias}(\hat{Y})\|^2 \leq \sigma^2$ (resp. na vzťah $\|\mathbf{M}\mathbf{Z}\gamma\|^2 \leq \sigma^2$). Potom z toho vyplýva

$$\text{MSE}(\hat{Y}_G) - \text{MSE}(\hat{Y}) \text{ je pozitívne semidefinitná, práve keď, platí } \|\mathbf{M}\mathbf{Z}\gamma\| \leq \sigma. \quad (4.21)$$

To znamená, že pri dostatočne malom vychýlení odhadu strednej hodnoty je menší model prijateľnejší, pretože minimalizuje kvadráty reziduí v budúcom pozorovaní.

4.1. Výber regresorov

Nech je uvažovaná regresná matica \mathbf{X} a vektor závislých veličín \mathbf{Y} . Úlohou je vybrať tie regresory, ktoré dobre popisujú náhodné veličiny v \mathbf{Y} . Prvou dôležitou analýzou dát je zistenie vzájomnej závislosti regresorov. Tu sa využíva korelačný koeficient (resp. korelačná matica), ktorý indikuje mieru závislosti medzi dvoma regresormi. Napríklad môžu byť dva regresory závislé na faktore, ktorý nie je známy. V prípade, že jeden regresor $\mathbf{x}_{\bullet i}$ (faktor) je silne lineárne závislý na druhom $\mathbf{x}_{\bullet j}$, potom si odhady \mathbf{b}_i , \mathbf{b}_j rozdelia pôsobenie jedného faktora. V tomto prípade, kedy sú stĺpce regresnej matice \mathbf{X} lineárne nezávislé, ale majú (v nejakom zmysle) silnú lineárnu závislosť medzi sebou, sa hovorí o *multikolinearite*. Tu je možné využiť metódu hlavných komponentov alebo napríklad faktorovú analýzu.

V prípade modelu úplnej hodnosti je teda situácia výberu regresorov nasledujúca. Variančná matica odhadov pre vektor $\boldsymbol{\beta}$ podľa vzťahu 2.13 bola určený ako $\text{Var}(\mathbf{b}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$. Nech je najprv uvažovaný najjednoduchší prípad, keď stĺpce regresnej matice sú na seba ortogonálne. Potom sú matice $\mathbf{X}'\mathbf{X}$ i $(\mathbf{X}'\mathbf{X})^{-1}$ diagonálne a teda jednotlivé odhady b_i na sebe nezávislé. To znamená, že pri uvažovaní podmodelu zo stĺpcov regresnej matice sa ich príslušné bodové odhady b_i už nebudú meniť. V odstavci 2.1 bolo uvedené, ako t-testom sú testované jednotlivé odhady pre β_i . Avšak ani pri nezávislosti jednotlivých odhadov sa nemožno spoľahnúť na t-test pri vytváraní vhodného podmodelu. Pretože pri testovaní, jednotlivých odhadov na hladine významnosti α , pravdepodobnosť správneho výberu podmodelu z pôvodného modelu s počtom regresorov klesá. V prípade závislosti jednotlivých odhadov b_i , sa táto situácia ešte komplikuje. Pri závislosti odhadov treba pri každom vyradení stĺpca z regresnej matice \mathbf{X} znova vypočítať odhady regresných koeficientov.

F-testom sa testuje nevýznamnosť viacerých koeficientov súčasne. Preto sa zdá byť vhodnou voľbou. Problém môže nastať v prípade, že model obsahuje veľký počet regresorov. V tomto prípade p regresorov je potrebné pomocou F-testu skontrolovať 2^p modelov. V prípade veľkého počtu modelov sa používa takzvaný hladný algoritmus (v angl. Greedy Search). Na tejto myšlienke je založená kroková regresia, ktorá bude následne popísaná.

4.1.1. Kroková regresia

Najprv budú uvedené ideí jednotlivých algoritmov pre *vzostupný výber* (v angl. Forward Selection) a *zostupný výber* (v angl. Backward Elimination), na ktorých je založená kroková regresia. Vzostupný výber je založený na postupnom pridávaní regresorov podľa významnosti. Začína sa teda z prázdny modelom a dá sa zhrnúť nasledovným postupom, ktorý je stručne zhrnutý na obrázku 6.

1. Nech je množine vybraných regresorov (označenej \mathcal{N}_{vyb}) priradená prázdna množina \emptyset . A ďalej nech množine všetkých p regresorov (označenej \mathcal{N}_{nevyb}) je priradená množina $\{1, \dots, p\}$, kde je každému regresoru priradené práve jedno číslo i z $1, \dots, p$ podľa označenia koeficientu β_i . t-testom so stupňom voľnosti $n - 1$ sa vyberie najvýznamnejšie nenulový odhad b_i . Potom do množiny \mathcal{N}_{vyb} sa priradí hodnota i a z množiny \mathcal{N}_{nevyb} sa vyradí hodnota i .

4.1. VÝBER REGRESOROV

2. Vypočíta sa hodnota RSS pre model obsahujúci regresory s indexami z \mathcal{N}_{vyb} , ktorá bude ďalej označená RSS_{mod} . Nech k vyjadruje počet prvkov \mathcal{N}_{vyb} . Obmenou vzťahu 4.2 bude testovaný tentokrát nadmodel za platnosti

$$\frac{RSS_{mod} - RSS^*}{RSS_{mod}/(n - k)} \sim F(1, n - k),$$

kde RSS^* je reziduálny súčet štvorcov pre nadmodel už uvažovaného modelu (vždy rozšírený o jeden regresor s indexom z \mathcal{N}_{nevyb}). Je zrejmé, že pri menšej hodnote RSS^* bude skôr zamietnutá hypotéza. Preto je hľadaná minimálna hodnota RSS^* zo všetkých uvažovaných nadmodelov, ktorá bude ďalej označená RSS_{roz} .

3. Ak

$$\frac{RSS_{mod} - RSS_{roz}}{RSS_{mod}/(n - k)} > F_{1-\alpha}(1, n - k),$$

- hypotéza $\beta_j = 0$ je zamietnutá, index j je priradený do množiny \mathcal{N}_{vyb} a vyradený z množiny \mathcal{N}_{nevyb} ,
- do hodnoty RSS_{mod} sa priradí hodnota RSS_{roz} ,
- prejde sa na krok 2,

inak hypotéza $\beta_j = 0$ nie je zamietnutá a algoritmus sa týmto ukončuje.

Zostupný výber je založený na rovnakej myšlienke, ale vychádza sa z úplného modelu, obsahujúceho všetky regresory. Potom sú postupne vyradované z modelu regresory najmenej významné pomocou podobného F-testu.

Kroková regresia (v angl. Stepwise Regression) funguje na princípe vzostupného výberu s tým, že v každom kroku sa zároveň použije i zostupný výber. Teda v každom kroku sa kontroluje množina \mathcal{N}_{vyb} , či sa v nej nenachádza index už nevýznamného regresného koeficientu. Predošlý postup vzostupného výberu bude rozšírený o nasledujúce dva kroky.

4. RSS_{mod} bude priradená hodnota z kroku 3 rozšíreného modelu RSS_{roz} . Nech opäť k vyjadruje počet prvkov \mathcal{N}_{vyb} . Pomocou vzťahu 4.2 bude testovaný podmodel nasledovne

$$\frac{RSS^* - RSS_{mod}}{RSS_{mod}/(n - k)} \sim F(1, n - k),$$

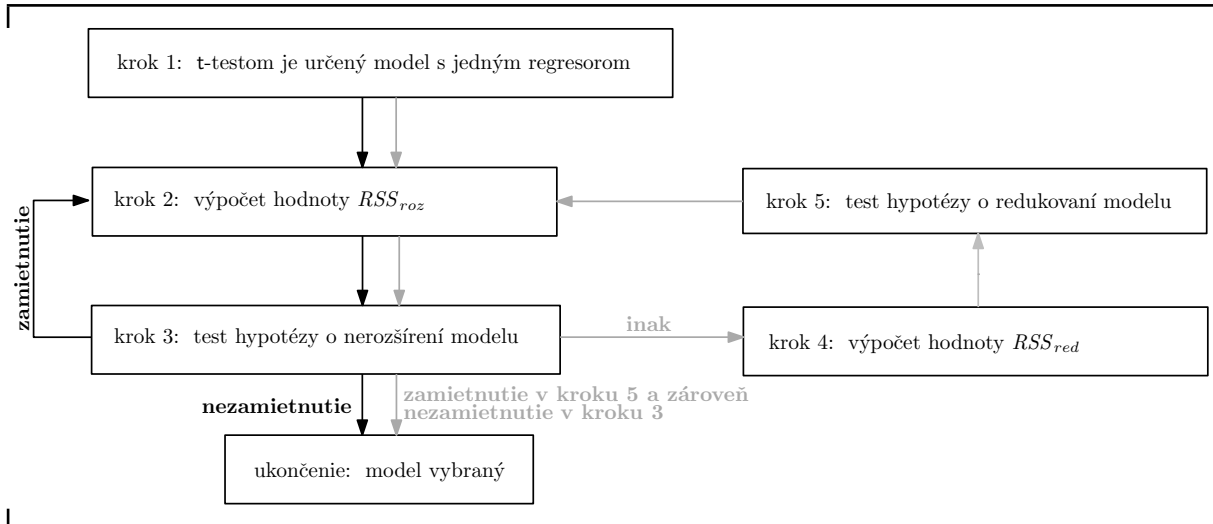
kde RSS^* je reziduálny súčet štvorcov pre podmodel už uvažovaného modelu (vždy redukovaný o jeden regresor s indexom z \mathcal{N}_{vyb}). Opäť je hľadaná minimálna hodnota RSS^* zo všetkých uvažovaných podmodelov, ktorá bude ďalej označená RSS_{red} .

5. Ak

$$\frac{RSS_{red} - RSS_{mod}}{RSS_{mod}/(n - k)} < F_{1-\alpha}(1, n - k),$$

- nie je zamietnutá hypotéza o nulovosti vyradeného regresoru, preto je číslo určujúce tento regresor vyradený z \mathcal{N}_{vyb} a priradený do \mathcal{N}_{nevyb} ,
- do hodnoty RSS_{mod} sa priradí hodnota RSS_{red} ,
- prejde sa na krok 2,

inak sa prejde na krok 2.



Obrázok 6: Na obrázku je porovnaný proces vzostupného výberu pomocou čiernych šípiek a proces krokovej regresie pomocou šedých šípiek.

Proces výberu regresorov je ukončený, ak prejde cyklus, v ktorom v krokoch 3 a 5 nebude možné rozšíriť alebo redukovať množinu vybraných regresorov \mathcal{N}_{vyb} .

Keďže tieto tri metódy fungujú na princípe hladného algoritmu, nemusia nájsť v danom zmysle najlepšie riešenie ani rovnaké (oproti ostatným). Napríklad vyber modelu bude závislý na zvolenej hladine významnosti α . Na túto tému Zvára uviedol v literatúre ([13], str. 144) nasledujúcu vetu: Obzvlášť pri krokovej regresii sa doporučuje vyhľadať niekoľko takmer optimálnych modelov a pokúsiť sa nájsť medzi nimi ten, ktorý má najlepšiu interpretáciu. Ďalej Jörnsten uviedla ([6], 5.prednáška) k tejto téme nasledujúcu vetu: Hladné algoritmy môžu viesť na modeli, ktoré sú ťažko interpretovateľné - regresory, ktoré sú korelované, súperia o zúčastnenie sa na tvorbe modelu a ľudské vedomosti (skúsenosti) môžu vybrať rozumné modely, čo tu štatistika nemôže.

V tomto odstavci, boli metódy výberu regresorov založené na F-teste. Rovnako tieto metódy môžu fungovať na základe iných testovacích kritérií. Známe sú napríklad Akaikeho informačné kritérium (známe pod skratkou AIC) alebo Bayesovské informačné kritérium (známe pod skratkou BIC).

4.1.2. Mallowsova C_p štatistika

Štatistika C_p je založená na odhade celkovej strednej štvorcovej chyby $\sum_{i=1}^n \text{MSE}(\hat{Y}_i)$, ktorej je možné sa dopustiť pri uvažovaní podmodelu miesto skutočného modelu. Nech je uvažovaný zavedený model 4.3, z ktorého sa vychádzalo, avšak v skutočnosti platí rozšírený model 4.5. Model 4.3 bude teraz uvažovaný úplnej hodnoty (hodnoty k). Potom podľa 4.16 platí

$$\sum_{i=1}^n \text{MSE}(\hat{Y}_i) = \sigma^2 \sum_{i=1}^n h_{ii} + (\mathbf{M}\mathbf{Z}\boldsymbol{\gamma})' \mathbf{M}\mathbf{Z}\boldsymbol{\gamma} = \sigma^2 \text{Tr}(\mathbf{H}) + \|\mathbf{M}\mathbf{Z}\boldsymbol{\gamma}\|^2. \quad (4.22)$$

Zo vzťahu 4.14 platí

$$\|\mathbf{M}\mathbf{Z}\boldsymbol{\gamma}\|^2 = \text{E}(RSS) - (n - h(\mathbf{X}))\sigma^2.$$

4.1. VÝBER REGRESOROV

Dosadením tohto vzťahu do 4.22 sa získa vzťah

$$\sum_{i=1}^n \text{MSE}(\hat{Y}_i) = \sigma^2 k + \text{E}(RSS) - (n - k)\sigma^2 = (2k - n)\sigma^2 + \text{E}(RSS).$$

Podelením σ^2 predchádzajúceho vzťahu vznikne

$$\frac{1}{\sigma^2} \sum_{i=1}^n \text{MSE}(\hat{Y}_i) = 2k - n + \frac{\text{E}(RSS)}{\sigma^2}. \quad (4.23)$$

Mallowsova C_p štatistika vyjadruje hodnotu 4.23 po dosadení za σ^2 odhadom tohto parametra s_G^2 (zo skutočného rozšíreného modelu) a za $\text{E}(RSS)$ je uvažovaná hodnota RSS z podmodelu nasledovne

$$C_p = 2k - n + \frac{RSS}{s_G^2}. \quad (4.24)$$

Potom pre strednú hodnotu tejto štatistiky a použitím vzťahu 4.14 platí

$$\text{E}(C_p) = 2k - n + \frac{\text{E}(RSS)}{\sigma^2} = 2k - n + \frac{\|\mathbf{M}\mathbf{Z}\boldsymbol{\gamma}\|^2 + (n - k)\sigma^2}{\sigma^2} = k + \frac{\|\mathbf{M}\mathbf{Z}\boldsymbol{\gamma}\|^2}{\sigma^2}. \quad (4.25)$$

Zo vzťahu 4.25 vyplýva, že pri malom vychýlení je stredná hodnota štatistiky C_p približne rovná k (počtu regresorov v podmodeli).

Ďalej bude uvažované, že skutočný model 4.5 je neznámy (väčšinou reálna situácia). No je uvažovaná množina p regresorov, ktorá v nejakom zmysle dobre zodpovedá modelu. Bez újmy na všeobecnosti a pre zachovanie modelu 4.3, nech medzi p regresormi sú všetky regresory z \mathbf{X} ($p \geq k$). Mallows v tomto prípade vo vzťahu 4.24 uvažuje, že všetkých p uvažovaných regresorov vytvára nevychýlený odhad $\hat{\mathbf{Y}}$ ($\text{bias}(\hat{\mathbf{Y}}) = \mathbf{0}$). Preto je z modelu s p regresormi počítaný odhad s_p^2 rozptylu σ^2 .

Pre takýto prípad je C_p štatistika pre uvažovaný model 4.3

$$C_p = 2k - n + \frac{RSS}{s_p^2}. \quad (4.26)$$

Takéto hodnoty štatistiky C_p sa vypočítajú pre všetky rôzne podmnožiny z množiny p regresorov. Pre n pozorovaní sa hodnota Mallowsovej štatistiky 4.26 mení na základe počtu k regresorov v podmodeli a hodnoty RSS podmodelu.

Navrhuje sa zobrazit' hodnoty C_p štatistiky podľa veľkosti podmodelov a zároveň v grafe zobrazit' úsečku z bodu $[0, 0]$ do bodu $[p, p]$. Podľa vzťahu 4.25 by sa hodnoty C_p pre málo vychýlené odhady mali pohybovať v blízkosti uvedenej úsečky (pre podmodely s m regresormi je to hodnota $\text{E}(C_p) = m$). V prípade väčšieho počtu regresorov p sa pre každý počet regresorov v podmodeli zobrazujú iba tie najbližšie k úsečke. Názorné sú grafy v 5. kapitole.

Mallows v článku [8] uvádza, čo je možné z takéhoto grafu vdedukovať. Ak regresné koeficienty sú nezávislé a tie ktoré sú nenulové, dosahujú vyšších hodnôt (vzhľadom k smerodajnej odchýlke ich odhadu), potom sú hodnoty C_p (podmodelov z nich tvorených) blízke okolo uvedenej úsečky (vyššie - až spoločne vytvárajú dobrý odhad strednej hodnoty). Ďalší prípadom je, ak sú regresory x_1 , x_2 a x_3 silne korelované vzájomne a podobne korelované s \mathbf{Y} . Potom každé dva regresory je možné vylúčiť ale nie všetky tri. Preto tieto všetky uvažované modely budú okolo úsečky (nižšie - pre menšie modely), pretože každý jeden regresor môže dobre odhadnúť strednú hodnotu. Obe situácie sú umelé, ale môžu niekedy pomôcť.

4.2. Diagnostika

Diagnostické metódy pomáhajú zisťovať hodnoty z nameraných dat, ktoré sú v nejakom zmysle významne vzdialené od väčšiny nameraných hodnôt. Tieto hodnoty je možné rozdeliť do dvoch základných skupín. Prvé sú nazývané *odľahlými pozorovaniami* (v angl. outliers), ktoré sú odľahlé z pohľadu závislej premennej. Druhé hodnoty sú nazývané *vplyvnými pozorovaniami* (v angl. leverage points), ktoré sú odľahlé z pohľadu nezávislej premennej. Niekedy môžu byť obe skupiny zaujímavé z fyzikálneho hľadiska. Avšak ak sú odľahlé pozorovania značne vzdialené od strednej hodnoty, vypočítanej bez týchto hodnôt, môžu významne ovplyvniť výsledné riešenie metódou najmenších štvorcov. Napríklad je možné pri takejto analýze využiť metódu Jackknife (vždy by sa vypočítala stredná hodnota bez jedného pozorovania). Práve preto, že majú významný vplyv na riešenie metódou najmenších štvorcov je dôležitá ich diagnostika. V prípade jednorozmerného modelu je situácia s diagnostikou jednoduchšia, často sa dajú obe skupiny určiť z bodového grafu. No v prípade mnohorozmerného modelu je diagnostika týchto hodnôt značne komplikovanejšia a tu sú diagnostické metódy veľmi užitočné. V tejto práci budú uvedené nasledujúce dve diagnostiky:

- diagnostika pomocou projekčnej matice \mathbf{H}
- studentizované reziduá

4.2.1. Diagnostika pomocou projekčnej matice \mathbf{H}

Z kapitoly 2. je známe, že projekčná matica \mathbf{H} je idempotentná a symetrická. Preto platí tiež rovnosť $\mathbf{H}\mathbf{H}' = \mathbf{H}$, z ktorej priamo vyplýva, že každý riadok $\mathbf{h}_{i\bullet}$ matice \mathbf{H} má hodnotu $\|\mathbf{h}_{i\bullet}\|^2$ uvedenú v prvku h_{ii} . Zo vzťahu 2.7 je známa platnosť vzťahu $\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$, ktorý je možné ekvivalentne formulovať nasledovne

$$\hat{Y}_i = \mathbf{h}_{i\bullet}\mathbf{Y}, \quad i = 1, \dots, n.$$

Z predchádzajúceho vzťahu je teda zrejmé, že jednotlivé odhady \hat{Y}_i vektora $\hat{\mathbf{Y}}$ sú závislé na váhach podľa riadku $\mathbf{h}_{i\bullet}$ (na vektore pozorovaní \mathbf{Y}). Rovnako stĺpec $\mathbf{h}_{\bullet j}$ matice \mathbf{H} informuje o vplyve pozorovania Y_j na všetky odhady vektora $\hat{\mathbf{Y}}$. A mierou tohto vplyvu pozorovania Y_j na odhady $\hat{\mathbf{Y}}$ môže byť diagonálny prvok h_{ii} .

Pre hodnotu h_{ii} platia nerovnosti $0 \leq h_{ii} \leq 1$. Prvá nerovnosť vyplýva z hore uvedenej normy $\|\mathbf{h}_{i\bullet}\|^2$. Druhá nerovnosť vyplýva z nerovnosti

$$h_{ii} = \|\mathbf{h}_{i\bullet}\|^2 \geq h_{ii}^2,$$

ktorá môže platiť len pre h_{ii} nezáporné a $h_{ii} \leq 1$. Nerovnosti $0 \leq h_{ii} \leq 1$ budú využité v nasledujúcej úvahe.

Pre odhad \mathbf{u} náhodnej chyby \mathbf{e} platí $\mathbf{u} = \mathbf{M}\mathbf{e} = (\mathbf{I} - \mathbf{H})\mathbf{e}$. Túto rovnosť je možné ekvivalentne formulovať nasledovne

$$u_i = (1 - h_{ii})e_i - \sum_{j \neq i} h_{ij}e_j, \quad i, j = 1, \dots, n.$$

Z predchádzajúceho vzťahu vyplýva, že pre hodnoty h_{ii} blízke 1, náhodná chyba e_i v pozorovaní \hat{Y}_i nie je vysvetlená v odhade u_i tejto chyby.

4.2. DIAGNOSTIKA

Z výsledkov Gaussovej-Markovovej vety a vety 2.6. plynú vzťahy pre variančné matice $\text{Var}\hat{\mathbf{Y}} = \sigma^2\mathbf{H}$ a $\text{Var}(\mathbf{u}) = \sigma^2\mathbf{M} = \sigma^2(\mathbf{I} - \mathbf{H})$. Z úvah, ktoré boli urobené v tomto odstavci, je možné usúdiť nasledujúce tvrdenie. V prípade hodnoty h_{ii} blízkej 1, je rozptyl odhadu \hat{Y}_i väčší ale zároveň rozptyl odhadu u_i menší. A rovnako to platí naopak.

4.2.2. Studentizované reziduá

Najprv bude zavedené nové označenie. $\mathbf{X}_{(i)}$ je označenie pre maticu typu $(n-1) \times k$, ktorá vznikla odobraním i -tého riadku matice $\mathbf{X}_{n \times k}$. $\mathbf{Y}_{(i)}$ značí vektor bez i -tého pozorovania. Nech je uvažovaný normálny lineárny regresný model

$$\mathbf{Y}_{(i)} = \mathbf{X}_{(i)}\boldsymbol{\beta}_{(i)} + \mathbf{e}_{(i)}, \quad (4.27)$$

kde $\boldsymbol{\beta}_{(i)}$ sú regresné koeficienty a pre vektor chýb platí $\mathbf{e}_{(i)} \sim \mathbf{N}_{(n-1)}(\mathbf{0}, \sigma^2\mathbf{I})$. Nech je ďalej uvažovaný normálny lineárny regresný model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{j}_i\gamma + \mathbf{e}, \quad (4.28)$$

kde $(\boldsymbol{\beta}', \gamma)$ sú regresné koeficienty, \mathbf{j}_i je i -tý stĺpec jednotkovej matice \mathbf{I}_n a pre vektor chýb platí $\mathbf{e} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2\mathbf{I})$. Ďalej bude uvažovaný prípad $h(\mathbf{X}_{(i)}) = h(\mathbf{X}) = k$, ktorý znamená, že pre vyradený riadok je hodnota parametrickej funkcie $\mathbf{x}_{i\bullet}\boldsymbol{\beta}_{(i)}$ odhadnuteľná. Potom je možné tvrdiť, že odhady stredných hodnôt $\mathbf{X}_{(i)}\mathbf{b}_{(i)}$ a $\mathbf{X}\mathbf{b} + \mathbf{j}_i c$ modelov 4.27 a 4.28 metódou najmenších štvorcov sú totožné. Totiž pozorovanie i nemá v druhom modeli 4.28 žiadny vplyv, pretože odhad c koeficientu γ nadobudne takú hodnotu, aby minimalizoval reziduum v bode $\mathbf{x}_{i\bullet}$. Odhad c teda vyjadruje reziduum v prvom modeli 4.27 v bode $\mathbf{x}_{i\bullet}$.

Pre prvý model sa vypočíta odhad $\mathbf{b}_{(i)}$ regresných koeficientov $\boldsymbol{\beta}_{(i)}$ nasledovne

$$\mathbf{b}_{(i)} = (\mathbf{X}'_{(i)}\mathbf{X}_{(i)})^{-1}\mathbf{X}'_{(i)}\mathbf{Y}_{(i)}.$$

Potom stredná hodnota takéhoto modelu v bode $\mathbf{x}_{i\bullet}$ je $\hat{Y}_{(i)} = \mathbf{x}_{i\bullet}\mathbf{b}_{(i)}$. A nakoniec bude zavedená hodnota $u_{(i)} = Y_i - \hat{Y}_{(i)}$, ktorá je totožná s odhadom c v druhom modeli.

Za podmienky $h(\mathbf{X}_{(i)}) = h(\mathbf{X})$ je druhý model 4.28 nadmodelom modelu 4.27, ak platí $h(\mathbf{X}_{(i)}) < h(\mathbf{X}, \mathbf{j}_i)$. Toto môže byť ďalej predpokladané, inak by sa jednalo o nezaujímavý prípad, kedy nový stĺpec regresnej matice neposkytuje novú informáciu (platilo by $c = 0$ pre nadmodel). Podľa vzťahu 4.7 a označeného odhadu \mathbf{c}_G rozširujúceho modelu totiž platí

$$\mathbf{c}_G = (\mathbf{Z}'\mathbf{M}\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{u}. \quad (4.29)$$

Pre odhad c zo vzťahu 4.29 teda platí

$$c = (\mathbf{j}'_i\mathbf{M}\mathbf{j}_i)^{-1}\mathbf{j}'_i\mathbf{u} = \frac{u_i}{m_{ii}}. \quad (4.30)$$

V predchádzajúcom vzťahu sa formálne jednalo o pseudoinverziu z regulárnej štvorcovej matice (jednoprvkovej), ktorá je totožná s inverznou maticou. Podľa vety, ktorú uviedol Zvára v literatúre ([13], str. 99), je podmienka $h(\mathbf{X}_{(i)}) = h(\mathbf{X})$ ekvivalentná s tvrdením $m_{ii} > 0$. Preto môže byť uvažovaný vzťah 4.30. Použitím vzťahu $\text{Var}(\mathbf{u}) = \sigma^2\mathbf{M}$ z vety 2.6. je možné vyjadriť vzťah

$$\text{Var}(c) = \text{Var}\left(\frac{u_i}{m_{ii}}\right) = \frac{\text{Var}(u_i)}{m_{ii}^2} = \frac{\sigma^2 m_{ii}}{m_{ii}^2} = \frac{\sigma^2}{m_{ii}}. \quad (4.31)$$

4. VÝBER MODELU A DIAGNOSTIKA

Keďže u_i má normálne rozdelenie a platia vzťahy 4.30 a 4.31, potom náhodná veličina $u_i/(\sigma\sqrt{m_{ii}})$ má rozdelenie $\mathbf{N}(0, 1)$. Nech $s_{(i)}$ značí odhad smerodatnej odchýlky v prvom modeli, potom náhodná veličina $s_{(i)}^2/\sigma^2$ má $\chi^2(n - 1 - k)$ rozdelenie. Podľa vety 2.10. sú odhady $\mathbf{b}_{(i)}$ nezávislé na odhade $s_{(i)}^2$, preto aj odhady $\hat{Y}_{(i)} = \mathbf{x}_{i\bullet}\mathbf{b}_{(i)}$ a $s_{(i)}^2$ sú na sebe nezávislé. A keďže obe tieto štatistiky boli vypočítané bez znalosti Y_i , je možné usúdiť, že sú všetky tri Y_i , $\hat{Y}_{(i)}$ a $s_{(i)}^2$ vzájomne nezávislé. Zo vzťahu 4.30 vyplýva vzťah

$$u_i = cm_{ii} = (Y_i - \hat{Y}_{(i)})m_{ii}.$$

Z predchádzajúceho vzťahu a tvrdení vyplýva nezávislosť medzi veličinami u_i a $s_{(i)}^2$. Potom je možné zaviesť novú štatistiku

$$u_i^* = \frac{u_i/(\sigma\sqrt{m_{ii}})}{\sqrt{s_{(i)}^2/\sigma^2}} = \frac{u_i}{\sqrt{m_{ii}s_{(i)}}} \sim \mathbf{t}(n - k - 1), \quad (4.32)$$

ktorá odpovedá definovanému Studentovmu rozdeleniu podľa prvého vyjadrenia. Táto štatistika sa nazýva *studentizované reziduum*. V prípade, kedy je vopred známy index i , je možné klasicky testovať t -testom odľahlosť pozorovania. Hypotéza o neodľahlosti pozorovania i je zamietnutá na hladine α ak platí $|u_i^*| > \mathbf{t}_{1-\alpha/2}(n - k - 1)$. No v prípade, kedy nie je vopred známy index i a jeho výber je závislý na pozorovaní \mathbf{Y} , je potrebné testovanie mierne pozmeniť.

Nech je pre $\delta \in (0, 1)$ uvažovaných n náhodných javov

$$A_i(\delta) = \{|u_i^*| > \mathbf{t}_{1-\delta/2}(n - k - 1)\}, \quad i = 1, 2, \dots, n.$$

Podľa Bonferroniho nerovnosti 1.1 platí pri zvolenej hladine α

$$\mathbf{P}(\cup_{i=1}^n A_i(\alpha/n)) \leq \sum_{i=1}^n \mathbf{P}(A_i(\alpha/n)) = \alpha.$$

Ak bude vybraný index j na základe štatistík z pozorovaní \mathbf{Y} , bude hypotéza o neodľahlosti pozorovania od strednej hodnoty $\hat{\mathbf{Y}}$ zamietnutá ak $|u_j^*| > \mathbf{t}_{1-\alpha/(2n)}(n - k - 1)$.

5. Analýza experimentu

Experiment bol vyhodnotený s použitím programovacieho jazyka a prostredia R. Ide o voľne dostupný jazyk a software, ktorý je určený primárne na štatistické výpočty a grafické zobrazovanie. Pre prácu v tomto prostredí je možné priamo využiť dostupnú konzolu. Avšak ako nadstavbové vyvojové prostredie bolo použité RStudio, ktoré je taktiež voľne šíriteľné. Toto vyvojové prostredie je možno doporučiť pre jednoduchosť ovládania oproti konzole (vzhľadom aj funkcionalitou pripomína vývojové prostredie MATLABu).

Po dohode s odborníkom v danej oblasti experimentu, bol zadaný problém, ktorý bolo potrebné vyriešiť. Jednalo sa o zistenie najlepšej kombinácie fosforečnanových taviacich solí v zmesi tak, aby bola zaistená maximálna tvrdosť. Prípadne nájsť takú množinu, ktorá bude takému zadaniu s nejakou pravdepodobnosťou odpovedať.

Nezávislé premenné x_1 , x_2 a x_3 budú vyjadrovať pomer určitej zložky na zmesi a teda bude platiť $x_1 + x_2 + x_3 = 1$.

V prípade vyhodnocovania experimentu troch zložiek (napr. zmesi) sa používajú ternárne diagramy. Výhodou je napríklad jednoduchá interpretovateľnosť grafických výsledkov. Avšak pre zadanú úlohu je lepšie uvažovať dve zložky, pretože sa týmto nemení stredná hodnota, ktorá má byť vyhodnocovaná. Takýto záver je možné urobiť na základe vzťahu 4.9

$$\mathbf{d} = \hat{\mathbf{Y}}_G - \hat{\mathbf{Y}} = \mathbf{MZ}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{u}, \quad (5.1)$$

ktorý určoval vektor vzdialenosti dvoch odhadov, ak jeden je širší (v uvažovanom prípade bol skutočný). Nech teda regresná matica \mathbf{X} (projekčnej matici \mathbf{M} vo vzťahu 5.1) je z modelu

$$\mathbf{EY} = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_1^2 + \beta_4x_1x_2 + \beta_5x_2^2. \quad (5.2)$$

No je uvažovaný i širší model s hodnotami x_3 . Takýto model by rozšíril strednú hodnotu nasledovne

$$\mathbf{EY} = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3 + \beta_4x_1x_2 + \beta_5x_2x_3 + \beta_6x_1x_3 + \beta_7x_1^2 + \beta_8x_2^2 + \beta_9x_3^2. \quad (5.3)$$

Ale člen x_3 je možné vyjadriť v tvare $1 - x_1 - x_2$, no nevznikne takto žiadny nový člen a bude platiť

$$\begin{aligned} \mathbf{EY} = & (\beta_0 + \beta_3 + \beta_9) + (\beta_1 - \beta_3 + \beta_6 - 2\beta_9)x_1 + (\beta_2 - \beta_3 + \beta_5 - 2\beta_9)x_2 + \\ & + (\beta_7 - \beta_6 + \beta_9)x_1^2 + (\beta_4 - \beta_5 - \beta_6 + 2\beta_9)x_1x_2 + (\beta_8 - \beta_5 + \beta_9)x_2^2. \end{aligned} \quad (5.4)$$

Ako je vidieť, takto uvažovaný model (s podmienkou $x_1 + x_2 + x_3 = 1$) totiž nerozširuje priestor strednej hodnoty. Preto by bola vo vzťahu 5.1 matica \mathbf{Z} nulová. Ale aj keby bol uvažovaný nadmodel, kde \mathbf{Z} by bola tvorená vektorom \mathbf{x}_3 , platilo by $\mathbf{x}_3 \in \mathcal{M}(\mathbf{1}, \mathbf{x}_1, \mathbf{x}_2)$ a teda $\mathbf{MZ} = \mathbf{0}$ vo vzťahu 5.1.

Nech sú odhady pre strednú hodnotu 5.2 označené v hornom indexe *. Keďže 5.2 a 5.4 vyjadrujú totožné stredné hodnoty, platí $\beta_0^* = \beta_0 + \beta_3 + \beta_9$ a podobne ďalších päť vzťahov. Odhady 5.2 sú jednoznačné. Teda sa jedná o sústavu šiestich rovníc pre desať neznámych. Čo nasvedčuje tomu, že neexistuje jednoznačné vyjadrenie pre vyjadrenie strednej hodnoty 5.3.

Podľa 3D bodového grafu, kde boli zobrazené hodnoty \mathbf{x}_1 a \mathbf{x}_2 na \mathbf{Y} , je zrejmé, že odhad strednej hodnoty rovinou nebude postačujúci. Pretože je takýto graf ťažko interpretovateľný do roviny, nebude tu uvedený, ale rozloženie dát je možné pozorovať na ďalších

stránkach s už odhadnutou strednou hodnotou. Nech je teda uvažovaná stredná hodnota tvaru 5.2 i vyššieho stupňa. Za regresor x_1 bude uvažovaný pomer hydrogenfosforečnanu sodného zo zmesi a za regresor x_2 pomer difosforečnanu sodného. Pozorovania \mathbf{Y} sú pre tvrdosť, 8 týždňov zrenia a 30 dní skladovania.

5.1. Polynóm druhého stupňa

Pomocou dvoch funkcií `summary(lm(y ~ x1 + x2))` je možné veľmi rýchlo získať základné informácie o uvedenom modeli (v tejto funkcii) v jazyku R. Ak je uvažovaný model so strednou hodnotou 5.2, potom sú k dispozícii štatistiky vo výstupe 1. Väčšina bola zavedená v 2. kapitole. Sú tu informácie o reziduách (napr. medián a kvartilové rozpätie), jednotlivé odhady regresných koeficientov metódou najmenších štvorcov a smerodajné odchýlky týchto odhadov. Následne hodnoty pre t-test a príslušné p -hodnoty. Hviezdičky určujú významnosť regresoru v modeli na hladine α podľa uvedenej legendy. Nižšie sú uvedené informácie o odhade smerodajnej odchýlky chyby modelu a koeficient determinácie R^2 (i adjustovaný R_{adj}^2) ako bol uvedený v 2.19 v prípade $\mathbf{1} \in \mathcal{M}(X)$. Posledný je F-test, ktorý testuje hypotézu, že všetky regresné koeficienty sú nulové.

Výstup z R 1: Štatistiky pre model so strednou hodnotou s regresormi do 2. stupňa

```
Call:
lm(formula = y ~ x1 + x2 + I(x1^2) + I(x1 * x2) + I(x2^2))

Residuals:
    Min       1Q   Median       3Q      Max
-3.1820 -0.7955  0.0686  0.7832  3.2342

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    6.098      0.482  12.651  <2e-16 ***
x1             -1.110      2.036  -0.545  0.5866
x2             -2.132      2.036  -1.047  0.2971
I(x1^2)        -4.830      1.933  -2.499  0.0137 *
I(x1 * x2)     37.549      3.156  11.898  <2e-16 ***
I(x2^2)        -1.991      1.933  -1.030  0.3048
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.223 on 126 degrees of freedom
Multiple R-squared:  0.8213, Adjusted R-squared:  0.8142
F-statistic: 115.8 on 5 and 126 DF,  p-value: < 2.2e-16
```

F-test jasne zamieta súčasnú nulovosť všetkých regresorov. Ako významné na hladine α vyšli tri regresory. Korelačný koeficient medzi $\mathbf{x}_{\bullet 2}$ a $\mathbf{x}_{\bullet 4}$ vyšiel 0.9481194. Teda regresory x_1, x_1^2 sú silne lineárne závislé a rovnako regresory x_2, x_2^2 .

5.1. POLYNÓM DRUHÉHO STUPŇA

Pomocou Krokovej regresie bol vybraný podmodel zložený zo 4 regresorov. Zdrojový kód použitej metódy je v prílohe. Metóda použila len výber regresorov v kroku 3 v takomto poradí: absolútny člen, x_1x_2 , x_1^2 a x_2^2 . Následne bude uvedené štatistiky pre tento vybraný model vo výstupe 2.

Výstup z R 2: Štatistiky pre vybraný model Krokovou regresiou

```
Call:
lm(formula = y ~ I(x1^2) + I(x1 * x2) + I(x2^2))

Residuals:
    Min       1Q   Median       3Q      Max
-3.1424 -0.7417  0.0556  0.8026  3.2432

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    5.6902     0.1991  28.581 < 2e-16 ***
I(x1^2)       -5.5444     0.5165 -10.735 < 2e-16 ***
I(x1 * x2)    34.9663     1.4857  23.536 < 2e-16 ***
I(x2^2)       -3.8927     0.5165  -7.537 7.72e-12 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

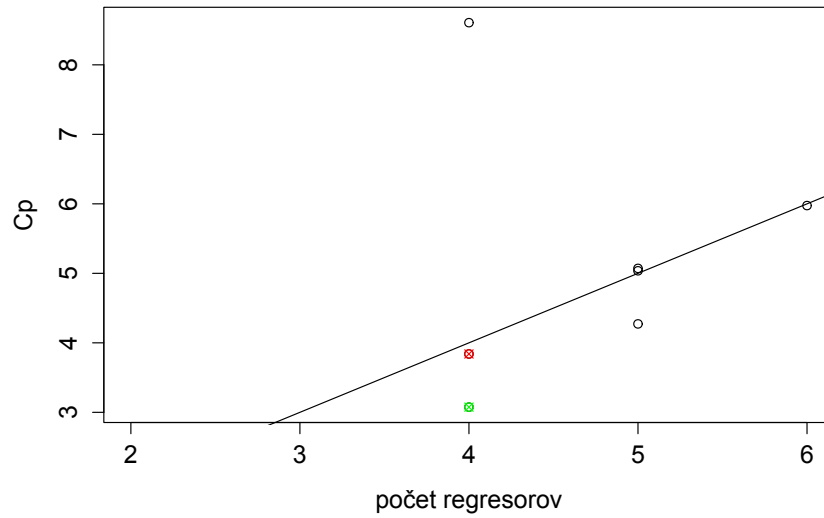
Residual standard error: 1.219 on 128 degrees of freedom
Multiple R-squared:  0.8198, Adjusted R-squared:  0.8155
F-statistic: 194.1 on 3 and 128 DF,  p-value: < 2.2e-16
```

Podľa štatistík, je vybraný model vhodný. Hodnota R^2 klesla nepatrne o 0.0015, ale model bol redukovaný o dva regresory. Všetky regresory sú významné.

V prípade Mallowsvej C_p štatistiky je uvažované, že odhad strednej hodnoty z pôvodného modelu je nestranným odhadom. Bude teda použitý vzťah 4.26

$$C_p = 2k - n + \frac{RSS}{s_p^2},$$

kde $s_p = 1.223$ z výstupu 1. Graf spomínaný v časti o Mallowsvej štatistike je zobrazený v grafe 1, kde je taktiež zobrazený bod pre vybraný model na zeleno. Sú tu vyobrazené hodnoty C_p menšie ako 10. Len informatívne, najväčšia bola hodnota 578.5447. V grafe 1 je vidieť, že k strednej hodnote C_p je bližšia hodnota s počtom regresorov 4 označená na červeno. Ide o model s regresormi: absolútny člen, x_2 , x_1x_2 a x_1^2 . Teda oproti Krokovej regresii ($C_p = 3.074860$) bol vybraný za lepší z pohľadu C_p ($C_p = 3.839803$) regresor x_2 miesto x_2^2 . To nie je nič prekvapivé, keďže oba regresory sú silno lineárne závislé (korelačný koeficient vyšiel 0.9481194). Základné štatistiky pre takto vybraný model sú uvedené vo výstupe 3. Keďže základné štatistiky oboch modelov sú veľmi podobné, bude ďalej uvažovaný model preferovaný na základe Mallowsvej štatistiky. V grafe 2 je uvedený 3D bodový graf s odhadnutou strednou hodnotou. V grafe 3 sú názornejšie zobrazené na-



Graf 1: Zobrazené C_p štatistiky pre pôvodný model. Zelenou farbou je vyznačený bod vybraný Krokovou regresiou.

merané dáta \mathbf{Y} s odhadovanými $\hat{\mathbf{Y}}$. Postupnosť zmesí (v grafe 3) je určená pre dvojicu x_1, x_2 nasledovne

$(1, 0), (0.9, 0.1), (0.8, 0.2), \dots, (0.1, 0.9), (0, 1), (0.9, 0), (0.8, 0.1), \dots, (0.1, 0), (0, 0.1), (0, 0)$.

Výstup z R 3: Štatistiky pre vybraný model Krokovou regresiou

Call:

```
lm(formula = y ~ I(x2) + I(x1^2) + I(x1 * x2))
```

Residuals:

| Min | 1Q | Median | 3Q | Max |
|---------|---------|--------|--------|--------|
| -3.0640 | -0.7895 | 0.0642 | 0.7468 | 3.2037 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 6.1247 | 0.2376 | 25.776 | < 2e-16 *** |
| I(x2) | -3.7697 | 0.5051 | -7.463 | 1.14e-11 *** |
| I(x1^2) | -6.1307 | 0.5586 | -10.975 | < 2e-16 *** |
| I(x1 * x2) | 37.8453 | 1.6264 | 23.269 | < 2e-16 *** |

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

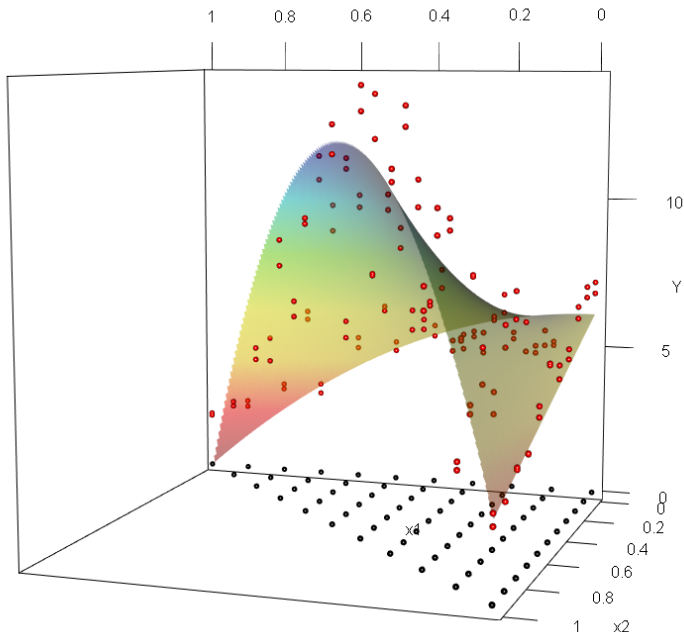
Residual standard error: 1.222 on 128 degrees of freedom

Multiple R-squared: 0.8187, Adjusted R-squared: 0.8144

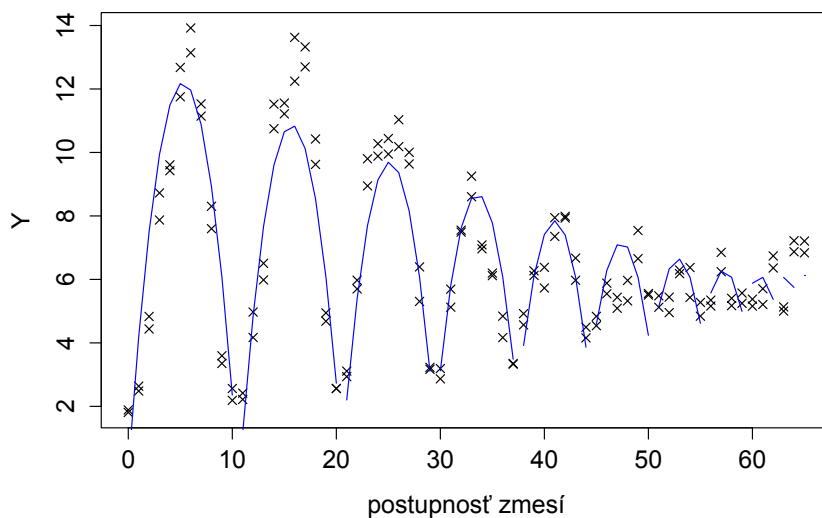
F-statistic: 192.6 on 3 and 128 DF, p-value: < 2.2e-16

5.1. POLYNÓM DRUHÉHO STUPŇA

Z uvedených grafov 2, 3 vyplýva, že odhadnutá stredná hodnota nepokrýva dostatočne maximálne hodnoty Y . Preto je vhodné, pre úlohu najdenia najväčšej tvrdosti, prejsť na model vyššieho stupňa.



Graf 2: 3D bodový graf pre model vybraný s pomocou C_p štatistiky. V grafe je zobrazená stredná hodnota určená metódou najmenších štvorcov.



Graf 3: Zobrazené sú hodnoty Y a po častiach spojitá lineárna funkcia, ktorá nadobúda rovnaké hodnoty s \hat{Y} pre zmesi. Postupnosť zmesí je určená v texte vyššie.

5.2. Polynóm tretieho stupňa

Nech je uvažovaný model so strednou hodnotou

$$EY = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_1 x_2 + \beta_5 x_2^2 + \beta_6 x_1^3 + \beta_7 x_1^2 x_2 + \beta_8 x_1 x_2^2 + \beta_9 x_2^3 \quad (5.5)$$

Základné štatistiky pre model 5.5 sú uvedené vo výstupe 4. Je zrejmé, že odhad smerodajnej odchýlky klesol, pretože sú uvažované nezávislé regresory navyše (4 oproti modelu 5.2).

Výstup z R 4: Štatistiky pre celý model

```
Call:
lm(formula = Y ~ x1 + x2 + I(x1^2) + I(x1 * x2) + I(x2^2) + I(x1^3) +
    I(x1^2 * x2) + I(x1 * x2^2) + I(x2^3))

Residuals:
    Min       1Q   Median       3Q      Max
-2.4686 -0.6791  0.0833  0.6273  2.5341

Coefficients:
                Estimate Std. Error t value Pr(>|t|)
(Intercept)      7.4283     0.5767  12.881 < 2e-16 ***
x1             -19.4881     4.6946  -4.151 6.16e-05 ***
x2              -9.5078     4.6946  -2.025 0.045024 *
I(x1^2)         42.4275    10.9122   3.888 0.000165 ***
I(x1 * x2)     78.2282    17.4705   4.478 1.71e-05 ***
I(x2^2)        12.7081    10.9122   1.165 0.246462
I(x1^3)       -30.2817     7.1199  -4.253 4.16e-05 ***
I(x1^2 * x2)  -61.4758    16.5451  -3.716 0.000307 ***
I(x1 * x2^2) -11.5376    16.5451  -0.697 0.486914
I(x2^3)       -10.0715     7.1199  -1.415 0.159742
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.084 on 122 degrees of freedom
Multiple R-squared:  0.8641, Adjusted R-squared:  0.8541
F-statistic: 86.22 on 9 and 122 DF,  p-value: < 2.2e-16
```

Výber regresorov bude uskutočnený najprv na základe Krokovej regresie. Metóda pre $\alpha = 0.05$ postupne vybrala nasledovné regresory: absolutný člen, $x_1 x_2$, x_3^1 , $x_1 x_2^2$ a x_3^2 . Rovnaké regresory vybrala aj pri $\alpha = 0.1$. Výpis štatistík pre takýto model sa nachádza vo výstupe 5.

V grafe 4 sú zobrazené všetky C_p štatistiky, ktoré majú menšiu hodnotu ako 20. Je vidieť, že model vybraný Krokovou regresiou, má najlepšiu C_p štatistiku ($C_p = 18.0526$) z ostatných podmodelov s piatimi regresormi. Ďalej je červenou farbou označený bod, ktorý má síce viac regresorov ale dobre dpovedá strednej hodnote Mallowsovej štatistiky ($C_p = 7.356668$). Tento model zahŕňa regresory: absolutný člen, x_1 , x_2 , x_2^1 , $x_1 x_2$, x_3^1 a $x_1^2 x_2$. Budú teda uvedené štatistiky pre tento model vo výpise 6.

5.2. POLYNÓM TRETIEHO STUPŇA

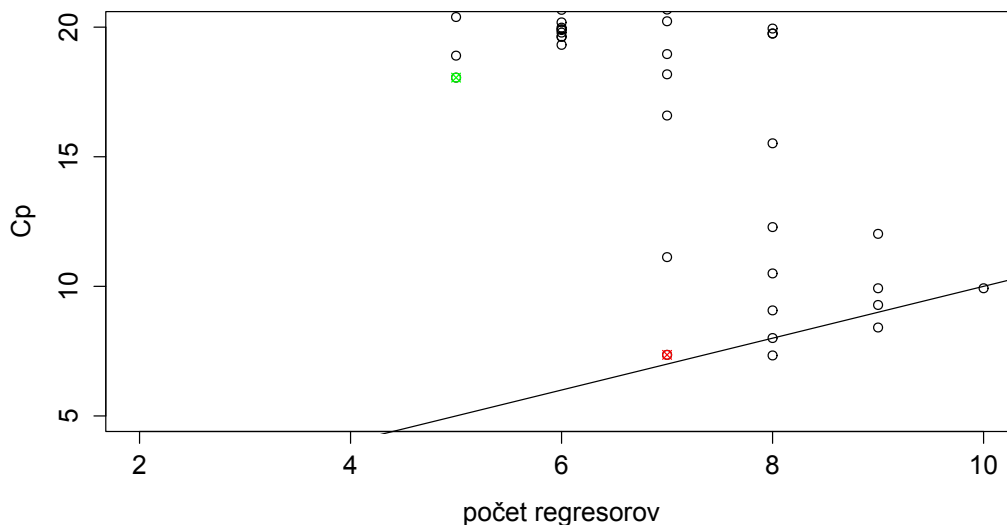
Výstup z R 5: Štatistiky pre celý model

```
Call:
lm(formula = Y ~ I(x1 * x2) + I(x1^3) + I(x1 * x2^2) + I(x2^3))

Residuals:
    Min       1Q   Median       3Q      Max
-2.39084 -0.60531  0.01664  0.70441  2.67335

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    5.5548     0.1741  31.912 < 2e-16 ***
I(x1 * x2)   18.0726     3.3691   5.364 3.72e-07 ***
I(x1^3)      -5.0751     0.5329  -9.524 < 2e-16 ***
I(x1 * x2^2) 28.6603     6.5276   4.391 2.36e-05 ***
I(x2^3)      -5.1508     0.5915  -8.709 1.39e-14 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.138 on 127 degrees of freedom
Multiple R-squared:  0.8439, Adjusted R-squared:  0.839
F-statistic: 171.7 on 4 and 127 DF,  p-value: < 2.2e-16
```



Graf 4: Zobrazené C_p štatistiky pre pôvodný model. Zelenou farbou je vyznačený bod vybraný Krokovou regresiou. Červenou farbou je označená štatistika modelu, ktorý je z pohľadu C_p a teda samotného výberu modelu zaujímavý.

Štatistiky druhého uvažovaného modelu vo výpise 6 sa zlepšili (napr. odhad smerodajnej odchýlky alebo hodnota R^2). Keďže je v modeli silnejšia lineárna závislosť niektorých

regresorov, nie je prekvapujúce, že v podstate iné regresory dokážu rovnako dobre odhadnúť strednú hodnotu. No kvôli odhadu strednej hodnoty by bolo rozumnejšie sa držať navrhovaného modelu podľa C_p . Síce má o dva regresory viac ale menšie vychýlenie podľa Mallowsvej štatistiky. Bude teda zobrazený bodový graf s odhadnutou strednou hodnotou v grafe 5. Bodový graf má body označené červenou farbou. Čiernou farbou sú označené body vyjadrujúce všetky sledované zmesi.

Výstup z R 6: Štatistiky pre vybraný model na základe C_p

```
Call:
lm(formula = Y ~ x1 + x2 + I(x1^2) + I(x1 * x2) + I(x1^3) + I(x1^2 * x2))

Residuals:
    Min       1Q   Median       3Q      Max
-2.9675 -0.7213  0.1014  0.6523  2.5963

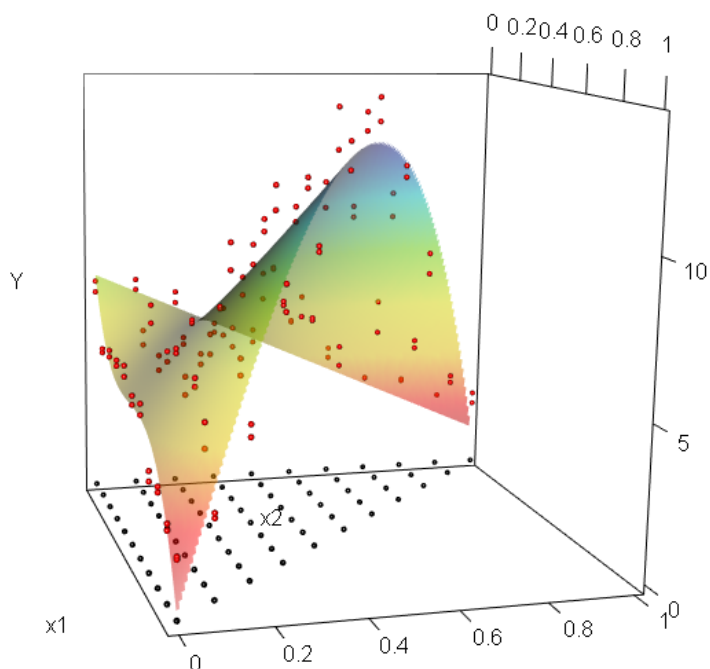
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    7.4256     0.3931  18.891 < 2e-16 ***
x1            -19.5647     3.9397  -4.966 2.19e-06 ***
x2             -6.1460     0.6609  -9.300 5.87e-16 ***
I(x1^2)        42.5981    10.0210   4.251 4.13e-05 ***
I(x1 * x2)     72.4134     5.8947  12.284 < 2e-16 ***
I(x1^3)       -30.3730     6.8037  -4.464 1.77e-05 ***
I(x1^2 * x2) -59.0026     9.8051  -6.018 1.82e-08 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.086 on 125 degrees of freedom
Multiple R-squared:  0.8603, Adjusted R-squared:  0.8536
F-statistic: 128.3 on 6 and 125 DF,  p-value: < 2.2e-16
```

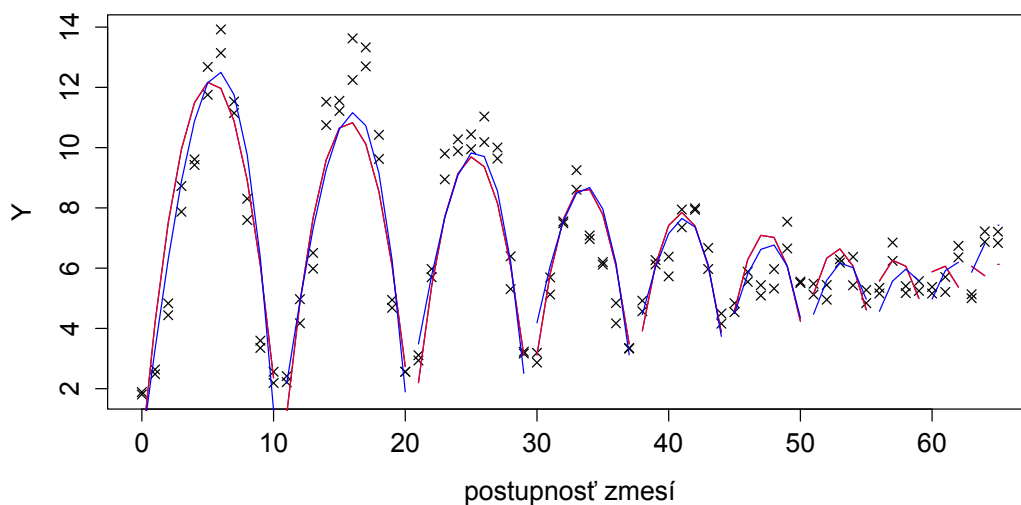
Ani v predchádzajúcom ani v tomto vybranom modeli studentizované reziduá na hladine $\alpha = 0.05$ nezaznamenali odľahlé pozorovania \mathbf{Y} .

V grafoch 5 a 6 je vidieť, že sa podarilo lepšie zachytiť strednú hodnotu pre zmesi s najväčšou tvrdosťou. Vzhľadom k analýze sa bude pokračovať k vyššiemu stupňu polynómu a budú pozorne sledované zmeny. Ak budú zmeny minimálne, bude preferovaný model vybraný v tomto odstavci.

5.2. POLYNÓM TRETIEHO STUPŇA



Graf 5: Zobrazený je 3D bodový graf s odhadnutou strednou hodnotou polynómom tretieho stupňa.



Graf 6: Zobrazené sú hodnoty Y a červenou farbou je označená funkcia totožná s funkciou z grafu 3 (odhad polynómom 2. stupňa). Modrou farbou je označená po častiach spojitá lineárna funkcia, ktorá nadobúda hodnoty \hat{Y} (odhad vybraného polynómu 3. stupňa) v príslušných zmesiach. Postupnosť zmesí bola už určená.

5.3. Polynóm štvrtého stupňa

Nech je uvažovaný model so strednou hodnotou

$$EY = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_1 x_2 + \beta_5 x_2^2 + \beta_6 x_1^3 + \beta_7 x_1^2 x_2 + \beta_8 x_1 x_2^2 + \beta_9 x_2^3 + \beta_{10} x_1^4 + \beta_{11} x_1^3 x_2 + \beta_{12} x_1^2 x_2^2 + \beta_{13} x_1 x_2^3 + \beta_{14} x_2^4. \quad (5.6)$$

Základné štatistiky pre tento model sú vo výstupe 7. Podľa teórie s novými lineárne nezávislými regresormi nevzrastie hodnota RSS (resp. väčšinou klesne). Presnejšie za predpokladu, že nové regresory majú priemety do reziduálneho priestoru pôvodného modelu a tieto priemety nie sú všetky ortogonálne na \mathbf{Y} , potom vždy poklesne hodnota RSS . Preto odhad s smerodajnej odchýlkou klesne a hodnota R^2 vzrastie. Toto je pozorovateľné na uvedených štatistikách (vo výstupoch) pôvodných a rozšírených modelov.

Výstup z R 7: Štatistiky pre úplný model

Call:

```
lm(formula = Y ~ x1 + x2 + I(x1^2) + I(x1 * x2) + I(x2^2) + I(x1^3) +
    I(x1^2 * x2) + I(x1 * x2^2) + I(x2^3) + I(x1^4) + I(x1^3 * x2) +
    I(x1^2 * x2^2) + I(x1 * x2^3) + I(x2^4))
```

Residuals:

| Min | 1Q | Median | 3Q | Max |
|----------|----------|----------|---------|---------|
| -2.08738 | -0.51084 | -0.02468 | 0.48056 | 2.06610 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(> t) | |
|----------------|-----------|------------|---------|----------|-----|
| (Intercept) | 6.4436 | 0.5398 | 11.938 | < 2e-16 | *** |
| x1 | -6.1838 | 7.4585 | -0.829 | 0.408739 | |
| x2 | 9.0341 | 7.4585 | 1.211 | 0.228240 | |
| I(x1^2) | 12.8747 | 31.5518 | 0.408 | 0.683983 | |
| I(x1 * x2) | -151.1174 | 49.0644 | -3.080 | 0.002580 | ** |
| I(x2^2) | -47.7718 | 31.5518 | -1.514 | 0.132703 | |
| I(x1^3) | -19.0753 | 47.9855 | -0.398 | 0.691707 | |
| I(x1^2 * x2) | 451.6944 | 108.8529 | 4.150 | 6.35e-05 | *** |
| I(x1 * x2^2) | 457.8222 | 108.8529 | 4.206 | 5.12e-05 | *** |
| I(x2^3) | 59.6301 | 47.9855 | 1.243 | 0.216474 | |
| I(x1^4) | 7.8460 | 23.7331 | 0.331 | 0.741543 | |
| I(x1^3 * x2) | -330.1447 | 71.7622 | -4.601 | 1.08e-05 | *** |
| I(x1^2 * x2^2) | -445.1972 | 102.1194 | -4.360 | 2.82e-05 | *** |
| I(x1 * x2^3) | -286.3343 | 71.7622 | -3.990 | 0.000116 | *** |
| I(x2^4) | -25.7157 | 23.7331 | -1.084 | 0.280798 | |

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

Residual standard error: 0.8616 on 117 degrees of freedom

Multiple R-squared: 0.9176, Adjusted R-squared: 0.9078

F-statistic: 93.11 on 14 and 117 DF, p-value: < 2.2e-16

5.3. POLYNÓM ŠTVRTEHO STUPŇA

Ďalej bude použitá Kroková regresia. Prvýkrát metóda využila krok 5, kedy vylučovala už vybraný regresor. Metóda postupovala nasledovne:

absolutný člen, x_1x_2 , x_1^3 , x_2^4 , $x_1^3x_2$, $x_1^2x_2$, odobraný x_1x_2 , $x_1x_2^2$, x_1 , odobraný x_1^3 , x_2 , $x_1^2x_2^2$.

Štatistiky pre tento model, ktorý má 8 regresorov, sú uvedené vo výstupe 8. Prvýkrát Kroková regresia vybrala regresory, ktorých odhady koeficientov majú p -hodnoty väčšie ako 0.001. Doteraz veľmi striktný výber regresorov, založený na Krokovej regresii a F-testu ($\alpha = 0.05$), sa zoslabil dôsledkom počtu nezávislých regresorov. Napríklad veľmi významný regresor v menšom modeli, ktorý je korelovaný s novým regresorom, vo väčšom modeli už taký významný nemusí byť (vyššia p -hodnota).

Výstup z R 8: Základné štatistiky pre model vybraný Krokovou regresiou

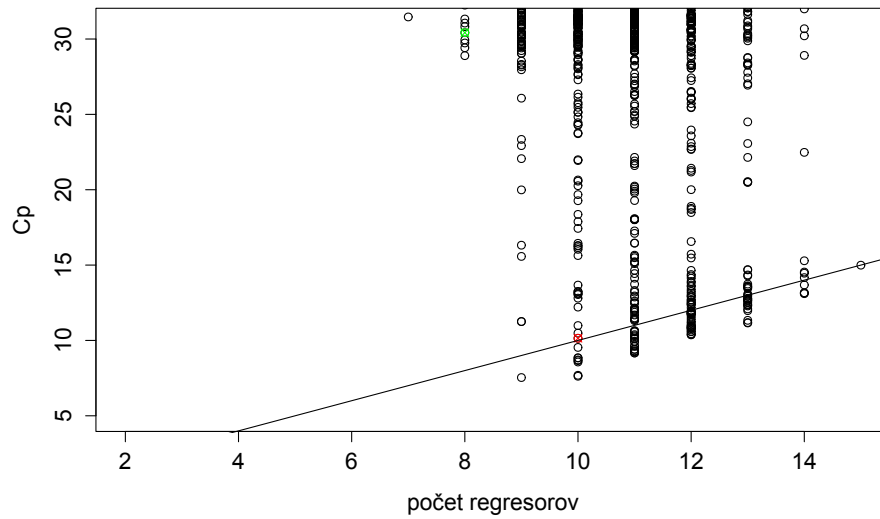
```
Call:
lm(formula = Y ~ x1 + x2 + I(x1^2 * x2) + I(x1 * x2^2) + I(x1^3 *
  x2) + I(x1^2 * x2^2) + I(x2^4))

Residuals:
    Min       1Q   Median       3Q      Max
-1.93603 -0.52415 -0.03862  0.59403  2.47676

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    6.3781     0.2770  23.025 < 2e-16 ***
x1             -4.2619     0.5252  -8.114 4.08e-13 ***
x2             -2.6627     0.8382  -3.177 0.001880 **
I(x1^2 * x2)   188.6860    18.2909  10.316 < 2e-16 ***
I(x1 * x2^2)   24.3791     6.5946   3.697 0.000326 ***
I(x1^3 * x2)  -205.4114    19.7481 -10.402 < 2e-16 ***
I(x1^2 * x2^2) -56.7671    25.6073  -2.217 0.028458 *
I(x2^4)        -2.3669     0.9168  -2.582 0.010996 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9363 on 124 degrees of freedom
Multiple R-squared:  0.8969, Adjusted R-squared:  0.8911
F-statistic: 154.1 on 7 and 124 DF,  p-value: < 2.2e-16
```

V grafe 7 sú označené C_p štatistiky do hodnoty 31. Zelenou farbou je tu označená hodnota $C_p = 30.42169$ vybraného modelu. Ďalej bola označená červenou farbou hodnota $C_p = 10.143824$, pre model s desiatimi regresormi, ktorá je s týmto počtom regresorov najbližšie k strednej hodnote C_p . Tento model má všetky regresory z modelu vybraného Krokovou regresiou a naviac regresory x_1x_2 , $x_1x_2^3$. Základné štatistiky pre tento model sú uvedené vo výstupe 9. Tu je vidieť, že regresor x_2^4 má p -hodnotu približne 0.0805. Teda t -test na hladine menšej ako 0.08 by už nezamietol hypotézu o nulovosti príslušného koeficientu. Keďže je primárne záujem o odhad strednej hodnoty, bude ďalej uvažovaný model vybraný na základe C_p so všetkými desiatimi regresormi.



Graf 7: Zelenou farbou je označená hodnota C_p pre model vybraný Krokovou regresiou a červenou farbou vybraný nový model na základe C_p .

Výstup z R 9: Základné štatistiky pre model vybraný na základe C_p

Call:

```
lm(formula = Y ~ x1 + x2 + I(x1 * x2) + I(x1^2 * x2) + I(x1 * x2^2)
    + I(x1^3 * x2) + I(x1^2 * x2^2) + I(x1 * x2^3) + I(x2^4))
```

Residuals:

| Min | 1Q | Median | 3Q | Max |
|----------|----------|----------|---------|---------|
| -2.10842 | -0.54463 | -0.01986 | 0.53136 | 2.06035 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(> t) | |
|----------------|-----------|------------|---------|----------|-----|
| (Intercept) | 6.9191 | 0.2777 | 24.914 | < 2e-16 | *** |
| x1 | -4.8015 | 0.4961 | -9.678 | < 2e-16 | *** |
| x2 | -3.6869 | 0.8712 | -4.232 | 4.51e-05 | *** |
| I(x1 * x2) | -109.6728 | 23.6793 | -4.632 | 9.15e-06 | *** |
| I(x1^2 * x2) | 428.2087 | 61.3065 | 6.985 | 1.62e-10 | *** |
| I(x1 * x2^2) | 327.2563 | 62.0074 | 5.278 | 5.77e-07 | *** |
| I(x1^3 * x2) | -335.9864 | 41.2304 | -8.149 | 3.67e-13 | *** |
| I(x1^2 * x2^2) | -376.9328 | 70.8638 | -5.319 | 4.80e-07 | *** |
| I(x1 * x2^3) | -197.5333 | 42.7102 | -4.625 | 9.40e-06 | *** |
| I(x2^4) | -1.7454 | 0.9903 | -1.762 | 0.0805 | . |

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

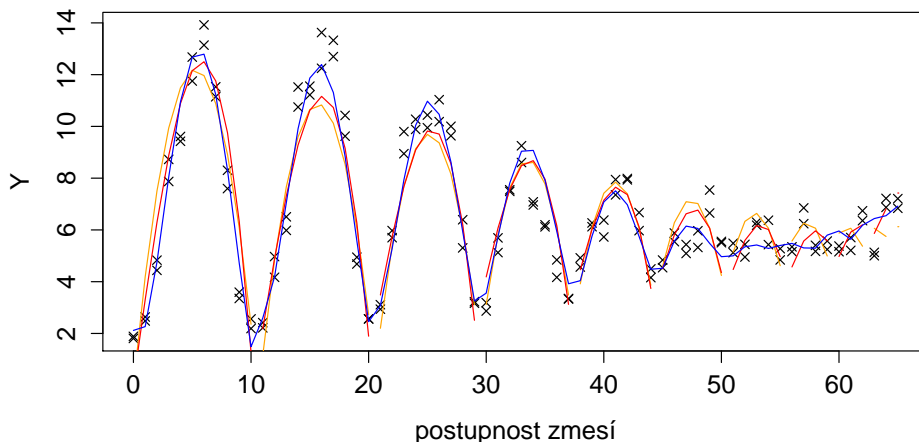
Residual standard error: 0.8621 on 122 degrees of freedom

Multiple R-squared: 0.914, Adjusted R-squared: 0.9077

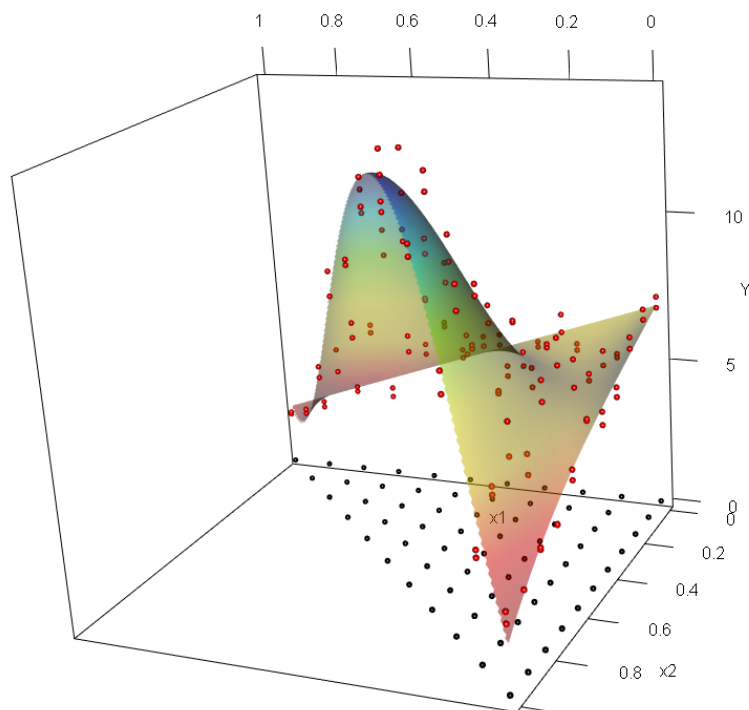
F-statistic: 144.1 on 9 and 122 DF, p-value: < 2.2e-16

5.3. POLYNÓM ŠTVRTÉHO STUPŇA

V grafoch 8 a 9 je vidieť, že sa podarilo zlepšiť odhad strednej hodnoty i pre zmesi s najvyššou tvrdosťou.

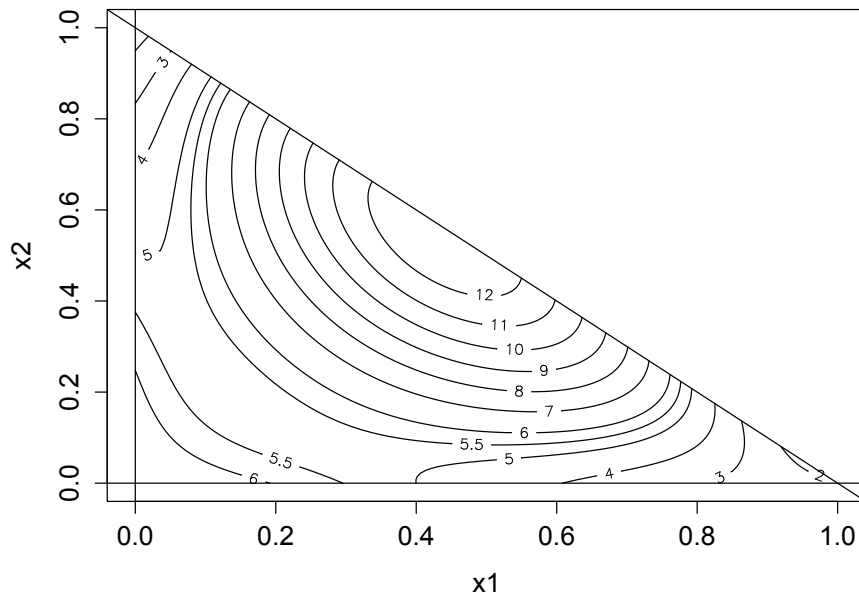


Graf 8: Funkcie označené červenou a oranžovou farbou sú funkcie z grafu 6 (oranžová odhad polynómom 2. stupňa, červená odhad polynómom 3. stupňa). Modrou farbou je označená po častiach spojitá lineárna funkcia, ktorá nadobúda hodnoty \hat{Y} (odhad vybraného polynómu 4. stupňa) v príslušných zmesiach. Postupnosť zmesí bola už určená.



Graf 9: 3D bodový graf s odhadnutou strednou hodnotou polynómom štvrtého stupňa vybraného na základe C_p .

V grafe 10 je zobrazený kontúrový graf pre odhadnutú strednú hodnotu s desiatimi regresormi. Izočiary s hodnotou 12 až po 5.5 nápadne pripomínajú tvar paraboloidu. A keďže sa práve v tejto oblasti (budú uvažované aj hraničné body) nachádza najvyššia tvrdosť, ktorú je potrebné vyšetriť, je možné modelovať strednú hodnotu v tejto oblasti pomocou paraboloidu. O tomto bude pojednávať nasledujúci odstavec.



Graf 10: Kontúrový graf odhadnutej strednej hodnoty.

5.3.1. Modelovanie tvrdosti paraboloidom

Nový model budú tvoriť hodnoty zmesi, ktoré sú označené v kontúrovom grafe 11 tmavomodrou farbou. Všetky sa nachádzajú v danej oblasti vymedzenou izočiariou (strednej hodnoty) s hodnotou 5.5. V tejto oblasti sa nachádza 33 rôznych zmesí a teda 66 pozorovaní z \mathbf{Y} . Pre tieto pozorovania bude uvažovaný model so strednou hodnotou

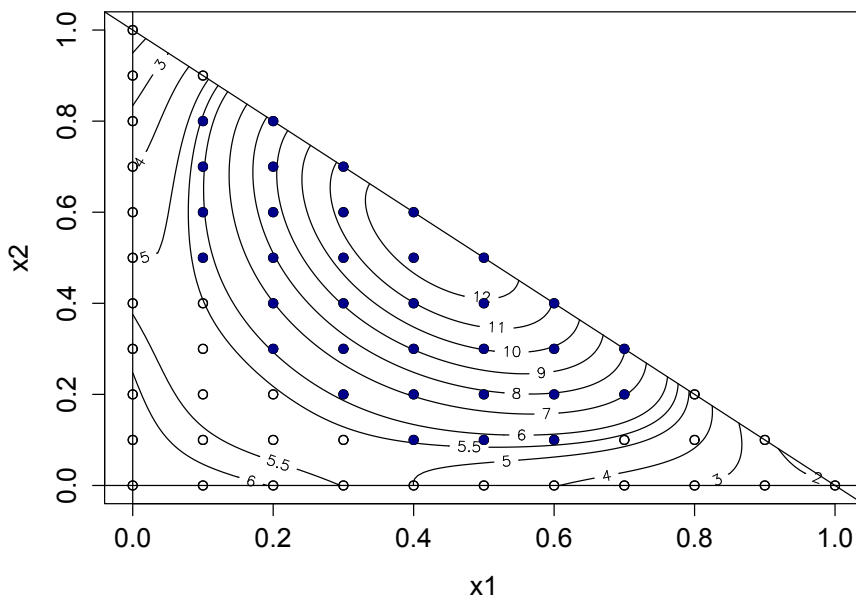
$$EY = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_1 x_2 + \beta_5 x_2^2.$$

Základné štatistiky pre takýto model sa nachádzajú vo výstupe 10.

Kroková regresia pre $\alpha = 0.05$, vybrala 4 regresory (absolútny člen, x_1 , x_1^2 a $x_1 x_2$). C_p štatistika pre tento model vyšla 13.684728 a v grafe 12 je označená červenou farbou. Bude opäť lepšie zamerať sa kvôli strednej hodnote viac na širší model ale menej odchýlený od celkového modelu (so šiestimi regresormi). Totiž v iných prípadoch, pri vyhodnocovaní lineárneho modelu, je dôležitá interpretácia jednotlivých regresorov a ich význam v modeli. Následne sa vyhodnocuje aký majú jednotlivé regresory vplyv na celkový model a z toho sa vyhodnocujú závery. No v prípade, ktorý rieši táto práca, ide primárne o odhad strednej hodnoty a následnom nájdení maximálnej tvrdosti (resp. príslušnej zmesi). Preto je dávaná prednosť metóde výberu pomocou C_p štatistiky, avšak s rozumným počtom regresorov. Táto hodnota je označená modrou farbou a nadobúda hodnotu $C_p = 4.057632$.

5.3. POLYNÓM ŠTVRTÉHO STUPŇA

Je vidieť, že uvažovaný model má päť regresorov (nie je uvažovaný regresor x_1x_2). Pre tento model sú uvedené štatistiky vo výstupe 11.



Graf 11: Kontúrový graf s vybranými hodnotami.

Výstup z R 10: Základné štatistiky pre celý model

Call:

```
lm(formula = Y ~ x1 + x2 + I(x1^2) + I(x1 * x2) + I(x2^2))
```

Residuals:

| Min | 1Q | Median | 3Q | Max |
|----------|----------|----------|---------|---------|
| -2.36715 | -0.67358 | -0.03526 | 0.65396 | 2.17594 |

Coefficients:

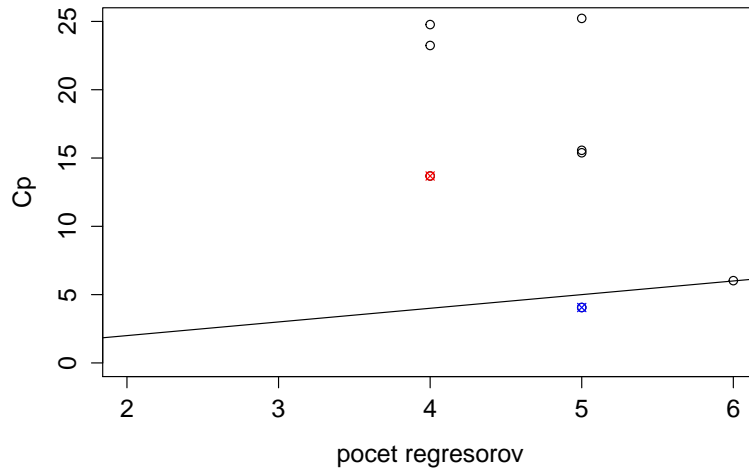
| | Estimate | Std. Error | t value | Pr(> t) | |
|-------------|----------|------------|---------|----------|-----|
| (Intercept) | -9.797 | 3.529 | -2.776 | 0.00732 | ** |
| x1 | 48.102 | 10.450 | 4.603 | 2.21e-05 | *** |
| x2 | 31.250 | 9.200 | 3.397 | 0.00121 | ** |
| I(x1^2) | -45.508 | 8.049 | -5.654 | 4.61e-07 | *** |
| I(x1 * x2) | -2.243 | 12.049 | -0.186 | 0.85296 | |
| I(x2^2) | -21.482 | 6.374 | -3.370 | 0.00132 | ** |

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

Residual standard error: 1.114 on 60 degrees of freedom

Multiple R-squared: 0.8346, Adjusted R-squared: 0.8208

F-statistic: 60.54 on 5 and 60 DF, p-value: < 2.2e-16



Graf 12: Modrou farbou je označený model, ktorý bol vybraný na základe C_p .

Výstup z R 11: Základné štatistiky pre vybraný model na základe C_p .

```
Call:
lm(formula = Y ~ x1 + x2 + I(x1^2) + I(x2^2))

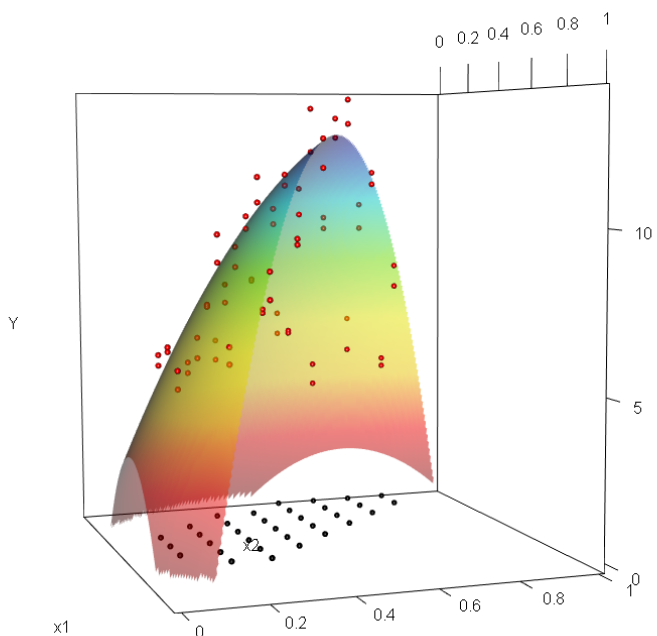
Residuals:
    Min       1Q   Median       3Q      Max
-2.34660 -0.67076 -0.05109  0.65131  2.18474

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)   -9.168     1.015  -9.032 7.54e-13 ***
x1             46.292     3.792  12.209 < 2e-16 ***
x2             29.636     3.048   9.724 5.15e-14 ***
I(x1^2)       -44.281     4.583  -9.662 6.54e-14 ***
I(x2^2)       -20.482     3.403  -6.018 1.09e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

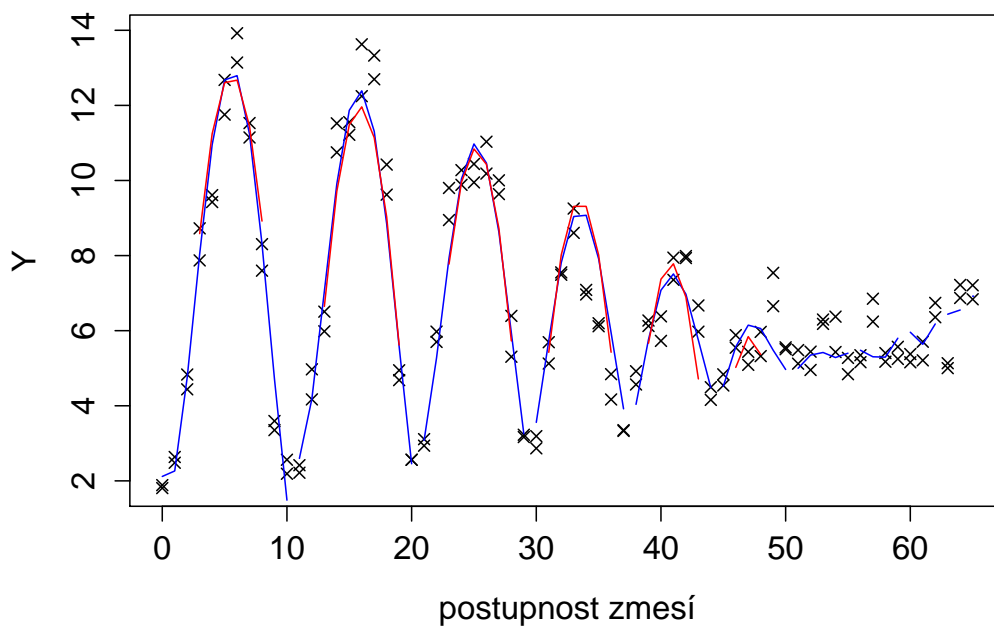
Residual standard error: 1.105 on 61 degrees of freedom
Multiple R-squared:  0.8345, Adjusted R-squared:  0.8236
F-statistic: 76.88 on 4 and 61 DF,  p-value: < 2.2e-16
```

Stredná hodnota pre tento model je zobrazená v grafe 13. V grafe 14 je vidieť, že odhad pomocou paraboloidu skutočne vyšiel dobre. Rozdiely v odhadoch sú v podstate minimálne, zrejme boli v pôvodnom modeli (4. stupňa) korigované vyšším stupňom polynómu. Keďže boli dôsledne prehodnocované zmesi (z odhadu strednej hodnoty polynómom 4. stupňa), ktoré primárne vytvárajú parabolický tvar, výrazne sa zamedzilo pákovému efektu. Je samozrejmé, že lokálny odhad strednej hodnoty (napríklad i na konvexnej množine) nemusí odpovedať odhadu urobenému na základe všetkých pozorovaní.

5.3. POLYNÓM ŠTVRTÉHO STUPŇA



Graf 13: 3D bodový graf s odhadnutou strednou hodnotou.



Graf 14: Porovnanie funkcie modrej farby totožnej s funkciou modrej farby v grafe 8 (odhad polynómom 4. stupňa). Graf funkcie označený červenou farbou (počastiach spojité a lineárna), nadobúda pre sledované zmesi stredných hodnôt odhadovaných paraboloidom.

5.4. Určenie zmesi s najväčšou tvrdosťou

Z predchádzajúceho odstavca 5.3.1 vyplýva, že na uvažovanej oblasti v grafe 11, je možné predpokladať model

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2 + e,$$

kde $e \sim N(0, \sigma^2)$. Na hranici, kde sa nachádza maximum (vid' graf 13), platí $x_2 = 1 - x_1$ a preto platí

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 (1 - x_1) + \beta_3 x_1^2 + \beta_4 (1 - x_1)^2 + e.$$

Úpravou sa prejde na tvar

$$Y = (\beta_0 + \beta_2 + \beta_4) + (\beta_1 - \beta_2 - 2\beta_4)x_1 + (\beta_3 + \beta_4)x_1^2 + e.$$

Pre hľadané maximum platí $\frac{\partial Y}{\partial x_1} = 0$, a preto

$$\frac{\partial Y}{\partial x_1} = (\beta_1 - \beta_2 - 2\beta_4) + 2(\beta_3 + \beta_4)x_1 = 0.$$

Z čoho vyplýva,

$$x_1 = \frac{-\beta_1 + \beta_2 + 2\beta_4}{2(\beta_3 + \beta_4)} = f_1(\beta_1, \beta_2, \beta_3, \beta_4).$$

Keďže z predchádzajúceho odstavca (z výstupu 11) vyplýva pre odhad metódou najmenších štvorcov $\mathbf{b} = (-9.168, 46.292, 29.636, -44.281, -20.482)'$, potom bodovým odhadom na hranici je

$$\hat{x}_1 = \frac{-b_1 + b_2 + 2b_4}{2(b_3 + b_4)} = \frac{-57.62}{-129.526} = 0.4448528.$$

A pre druhú zložku platí odhad

$$\hat{x}_2 = 1 - \hat{x}_1 = 1 - 0.4448528 = 0.5551472.$$

Preto bodovým odhadom pre najväčšiu tvrdosť je binárna zmes zložená

$$(\hat{x}_1, \hat{x}_2) = (0.4448528, 0.5551472),$$

kde x_1 vyjadruje pomer hydrogenfosforečnanu sodného zo zmesi a x_2 pomer difosforečnanu sodného.

Ďalej bude využitá delta metóda ktorá bola uvedená v úvode. V odstavci o normálnom lineárnom modeli úplnej hodnosti bolo vo vete 2.7 uvedené rozdelenie pre odhad regresných koeficientov

$$\mathbf{b} \sim N_k(\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}).$$

Odhad \mathbf{b}_n (z n pozorovaní) má asymptoticky k -rozmerné normálne rozdelenie $N_k(\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1})$ (v tomto odstavci $k = 5$). Aby bolo možné dôjsť k podobnému záveru ako v 1.5 je potrebné vypočítať nasledujúce parciálne derivácie funkcie f_1

$$\begin{aligned} \phi_0 = \frac{\partial f_1}{\partial \beta_0} &= 0, & \phi_1 = \frac{\partial f_1}{\partial \beta_1} &= \frac{-1}{2(\beta_3 + \beta_4)}, \\ \phi_2 = \frac{\partial f_1}{\partial \beta_2} &= \frac{1}{2(\beta_3 + \beta_4)}, & \phi_3 = \frac{\partial f_1}{\partial \beta_3} &= \frac{\beta_1 - \beta_2 - 2\beta_4}{2(\beta_3 + \beta_4)^2}, \\ \phi_4 = \frac{\partial f_1}{\partial \beta_4} &= \frac{4(\beta_3 + \beta_4) - 2(-\beta_1 + \beta_2 + 2\beta_4)}{4(\beta_3 + \beta_4)^2} = \frac{\beta_1 - \beta_2 + 2\beta_3}{2(\beta_3 + \beta_4)^2}, \end{aligned}$$

5.4. URČENIE ZMESI S NAJVÄČŠOU TVRDOSŤOU

kde označenie $\boldsymbol{\phi} = (\phi_0, \phi_1, \phi_2, \phi_3, \phi_4)'$ korešponduje zo zavedením pri uvedení delta metódy. Potom podľa 1.5 platí

$$\sqrt{n}(f_1(\mathbf{b}_n) - f_1(\boldsymbol{\beta})) \text{ konverguje v distribúcii k } \mathbf{N}(0, \sigma^2 \boldsymbol{\phi}'(\mathbf{X}'\mathbf{X})^{-1} \boldsymbol{\phi}). \quad (5.7)$$

Vektor $\boldsymbol{\phi}$ je neznámy a je nahradený $\hat{\boldsymbol{\phi}} = \boldsymbol{\phi}(\mathbf{b}_n)$. Rozptyl σ^2 je taktiež neznámy a nahradený podľa delta metódy $s^2 = \sigma^2(\mathbf{b}_n)$.

Pre odhad rozptylu 5.7 platí $D(\hat{x}_1)$

$$D(\hat{x}_1) = \hat{\boldsymbol{\phi}}' \text{Var}(\mathbf{b}_n) \hat{\boldsymbol{\phi}} = s^2 \hat{\boldsymbol{\phi}}'(\mathbf{X}'\mathbf{X})^{-1} \hat{\boldsymbol{\phi}},$$

kde reálny vektor

$$\hat{\boldsymbol{\phi}} = \left(0, \frac{-1}{2(b_3 + b_4)}, \frac{1}{2(b_3 + b_4)}, \frac{b_1 - b_2 - 2b_4}{2(b_3 + b_4)^2}, \frac{b_1 - b_2 + 2b_3}{2(b_3 + b_4)^2} \right)'$$

Získané boli hodnoty

$$\hat{\boldsymbol{\phi}} = (0, 0.007720451, -0.007720451, 0.006868902, -0.008572)'$$

a

$$(\mathbf{X}'\mathbf{X})^{-1} = \begin{pmatrix} 0.8432522 & -2.2224712 & -1.4219867 & 2.176086 & 1.025125 \\ -2.2224712 & 11.7659634 & -0.5005523 & -13.692503 & 1.776762 \\ -1.4219867 & -0.5005523 & 7.6026069 & 1.028647 & -8.099799 \\ 2.1760862 & -13.6925034 & 1.0286466 & 17.191063 & -1.985065 \\ 1.0251250 & 1.7767623 & -8.0997994 & -1.985065 & 9.480443 \end{pmatrix}.$$

Pre interval spoľahlivosti s bodovým odhadom $\hat{x}_1 = 0.4448528$ platí

$$\left(\hat{x}_1 - u_{1-\alpha/2} s \sqrt{\hat{\boldsymbol{\phi}}'(\mathbf{X}'\mathbf{X})^{-1} \hat{\boldsymbol{\phi}}}, \hat{x}_1 + u_{1-\alpha/2} s \sqrt{\hat{\boldsymbol{\phi}}'(\mathbf{X}'\mathbf{X})^{-1} \hat{\boldsymbol{\phi}}} \right),$$

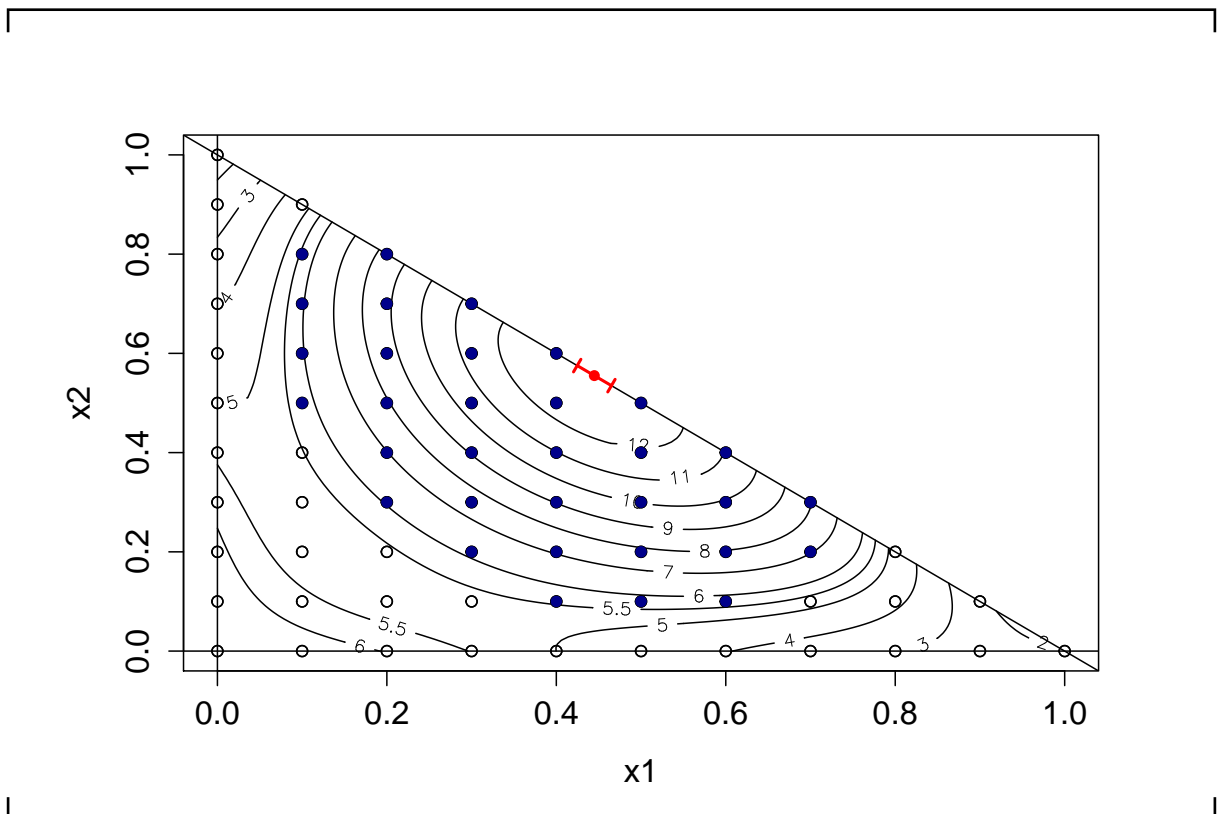
kde hodnota smerodajnej odchýlky $s = 1.105$ je určená z výstupu 11. Po dosadení hodnôt a pre $\alpha = 0.05$ výjde

$$\left(\hat{x}_1 - 1.96 \cdot 1.105 \sqrt{8.701371 \cdot 10^{-5}}, \hat{x}_1 + 1.96 \cdot 1.105 \sqrt{8.701371 \cdot 10^{-5}} \right),$$

a nakoniec

$$(\hat{x}_1 - 0.02020283, \hat{x}_1 + 0.02020283).$$

Jedná sa teda o interval spoľahlivosti pre \hat{x}_1 v binárnej zmesi. V grafe 15 je zaznačený tento interval aj s bodovým odhadom.



Graf 15: Zobrazenie bodového aj intervalového odhadu do kotúrového grafu odhadnutej strednej hodnoty polynómom 4. stupňa (vid' graf 11).

6. Záver

Ciele tejto diplomovej práce sa podarilo naplniť, ako po teoretickej tak i aplikačnej stránke.

V teoretickej časti práce bol v druhej kapitole precízne zavedený lineárny regresný model s úplnou stĺpcovou hodnotou a následne v tretej kapitole model s neúplnou stĺpcovou hodnotou. Ďalšia časť popisuje metódy pre výber regresorov a diagnostiku lineárneho regresného modelu. Celkovo sa podarilo urobiť ucelený výklad teórie s logickým usporiadaním a doplnený vhodnými obrázkami. Z nich niektoré, ako napr. obrázok **3b**, alebo **5**, neboli podľa vedomia autora publikované a môžu napomôcť pri pochopení danej látky.

V aplikačnej časti, ktorá je uvedená v piatej kapitole, sa podarilo pomocou lineárnej regresie úspešne vyhodnotiť experiment. To dokladujú aj výsledky, ktoré sú demonštrované graficky. Bola použitá polynomická regresia a výber regresorov bol založený na Krokovej regresii a Mallowsvej C_p štatistike. Následne bola odhadnutá stredná hodnota. Aby mohla byť odhadnutá zmes taviacich solí s najväčšou tvrdosťou, bola na vybranej oblasti stredná hodnota modelovaná paraboloidom. Využitím tohto odhadu bol nájdený bodový odhad a ďalej pomocou Delta metódy intervalový odhad.

Experiment bol sledovaný na ternárnych zmesiach tvorených hydrogenfosforečnanom sodným, difosforečnanom sodným a polyfosforečnanom sodným. Regresnou analýzou sa ukázalo, že najväčšia tvrdosť je dosiahnutá pri použití binárnej zmesi (8 týždňov zrenia, 30 dní skladovania). Bodovým odhadom je zmes tvorená: 44,485 % hydrogenfosforečnanu sodného a 55,515 % difosforečnanu sodného. Pre hydrogenfosforečnan sodný bol zistený 95% interval spoľahlivosti pre najväčšiu tvrdosť v binárnej zmesi a to interval (42.465, 46.505) v percentách.

Výsledky tejto práce budú využité pri výrobe tavených syrov.

Literatura

- [1] AGRESTI, A. *Categorical Data Analysis*. 2nd edition. Hoboken: Wiley Press, 2002. ISBN 0-471-36093-7
- [2] ANDĚL, J. *Matematická statistika*. 1. vyd. Praha: SNTL/ALFA, 1978.
- [3] ANDĚL, J. *Základy matematické statistiky*. 3. vyd. Praha: MATFYZPRESS, 2011. ISBN 978-80-7378-162-0
- [4] BROOK, R. J., ARNOLD, G. C. *Applied Regression Analysis and Experimental Designs*. 1st edition. New York: Marcel Dekker, 1985. ISBN 0-8247-7252-0
- [5] BUŇKA, F., et al. The effect of ternary emulsifying salt composition and cheese maturity on the textural properties of processed cheese [online]. *International Dairy Journal*, March 2013, vol. 29, iss. 1, pages 1-7 [cit. 2013-12-04]. Dostupné na: <<http://dx.doi.org/10.1016/j.idairyj.2012.09.006>>
- [6] JÖRNSTEN, R. *Linear Statistical Models* (lectures) [online]. Mathematical Sciences, Chalmers University of Technology and the University of Gothenburg, SE-412 96 Göteborg, Sweden. 2011. [cit. 2011-10-11]. Dostupné na: <<http://www.math.chalmers.se/Stat/Grundutb/GU/MSG500/A11/>>
- [7] KAPOOR, R., METZGER, L. E. Process Cheese: Scientific and Technological Aspects — A Review [online]. *Comprehensive Reviews in Food Science and Food Safety*, March 2008, vol. 7, iss. 2, pages 194–214 [cit. 2012-12-10]. Dostupné na: <<http://onlinelibrary.wiley.com/doi/10.1111/j.1541-4337.2008.00040.x/pdf>>
- [8] MALLOWS, C. L. Some Comments on Cp [online]. *Technometrics*, November 1973, vol. 15, No. 4, pages 661-675 [cit. 2013-27-04]. Dostupné na: <<http://www.jstor.org/discover/10.2307/1267380?uid=3737856&uid=2&uid=4&sid=21102278590171>>
- [9] RAO, R. C. *Lineární metody statistické indukce a jejich aplikace*. 1. vyd. Praha: Academia nakladatelství ČSAV, 1978.
- [10] STAUDTE, R. G., SHEATHER, S. J. *Robust estimation and testing*. 1st edition. Wiley Press, INC. 1990. ISBN 0-471-85547-2
- [11] WEISEROVÁ, E., et al. The effect of combinations of sodium phosphates in binary mixtures on selected texture parameters of processed cheese spreads [online]. *International Dairy Journal*, December 2011, vol. 21, iss. 12, pages 979-986 [cit. 2012-12-07]. Dostupné na: <<http://dx.doi.org/10.1016/j.idairyj.2011.06.006>>
- [12] ZVÁRA, K. *Regresní analýza*. 1. vyd. Praha: Academia nakladatelství ČSAV, 1989. ISBN 80-200-0125-5
- [13] ZVÁRA, K. *Regrese*. 1. vyd. Praha: MATFYZPRESS, 2008. ISBN 978-80-7378-041-8

LITERATURA

Príloha

Ako príloha bude uvedený iba zdrojový kód pre Krokovú regresiu, pretože táto metóda bola v odstavci 4.1.1 podrobne popísaná. K analýze experimentu teda nebola použitá v jazyku R na to určená funkcia `step()` (používajúca primárne kritérium AIC), ale tu uvedená verzia založená na F-teste.

V zdrojovom kóde 1 je uvedený krok 1, ktorý na základe t-testu vyberá prvý regresor. Tu sú súbežne počítané hodnoty pre všetky modely s jediným regresorom. Ako kritérium pre najviac významný regresor je uvažovaná maximálna hodnota testovacieho kritéria (v absolútnej hodnote).

Zdrojový kód 1: Kroková regresia (krok 1)

```
#1.krok
Y<-data[,13]                #pozorovania Y
X<-cbind(matrix(1,length(x1),1),x1,x2,x1^2,x1*x2,x2^2) #regresna matica
XY<-t(X)%*%Y
XX<-1/(diag(t(X)%*%X))
b<-XX*XY                    #vektor jednotliych odhadov beta
rezid<-X%*%diag(as.vector(b))-Y%*%t(as.vector(b)) #matica rezidui modelov
RSS<-diag(t(rezid)%*%rezid) #RSS hodnoty pre jednotlivy modely
odh_rozp<-RSS/131          #odhad rozptylu pre jednotlivy modely
test<-b/sqrt(XX*odh_rozp)  #jednotlive statistiky pre t-test
Nvyb <- which.max(abs(test)) #vybraty regresor
```

V zdrojovom kóde 2 je uvedený krok 2 až 5 ako boli popísané v tejto práci. Jediným nepodstatným rozdielom je nepoužitie množiny \mathcal{N}_{nevyb} , ktoré bolo nahradené v zmysle komplementu množiny \mathcal{N}_{vyb} . Označenie je zachované podľa popísanej metódy. Všetky výpočty štatistík sú názorne zapísané vzťahmi tak, ako boli popísané v teoretickej časti.

Zdrojový kód 2: Kroková regresi (krok 2-5)

```
#2.krok
krok5<-TRUE
repeat{
  X1<-X[,Nvyb]           #regresna matica urcena Nvyb
  H1<-X1%%solve(t(X1)%X1)%t(X1)
  M1<-diag(1,length(Y))-H1
  RSSmod<-t(Y)%M1%Y      #RSS vybranych regresorov
  RSSroz<-RSSmod
  for(j in 1:dim(X)[2]) { #vypocet RSS pre vsetky nadmodely(o 1 regresor)
    if (!(j%in%Nvyb)) {
      X2<-X[,c(Nvyb,j)]
      H2<-X2%%solve(t(X2)%X2)%t(X2)
      M2<-diag(1,length(Y))-H2
      RSS2<-t(Y)%M2%Y
      if (RSSroz>RSS2){
        RSSroz<-RSS2
        index<-j }
    }
  }
}

#3.krok
F<-(RSSmod-RSSroz)/(RSSmod/(length(Y)-length(Nvyb)))
#nasleduje F-test o nerozsireni modelu
if(F>qf(0.95,1,length(Y)-length(Nvyb))){
  Nvyb<-c(Nvyb,index)
  RSSmod<-RSSroz
}
if((F<=qf(0.95,1,length(Y)-length(Nvyb)))&&(!krok5)){break}

#4.krok
RSSred<-Inf
for(j in 1:dim(X)[2]) { #vypocet RSS pre vsetky podmodely(o 1 regresor)
  if (j%in%Nvyb){
    X3<-X[,setdiff(Nvyb,j)]
    H3<-X3%%solve(t(X3)%X3)%t(X3)
    M3<-diag(1,length(Y))-H3
    RSS3<-t(Y)%M3%Y
    if (RSSred>RSS3){
      RSSred<-RSS3
      index<-j }
    }
  }

#5.krok
F<-(RSSred-RSSmod)/(RSSmod/(length(Y)-length(Nvyb)))
if(F<qf(0.95,1,length(Y)-length(Nvyb))){
  Nvyb<-setdiff(Nvyb,index)
  krok5<-TRUE }
if(F>=qf(0.95,1,length(Y)-length(Nvyb))){krok5<-FALSE}
}
```