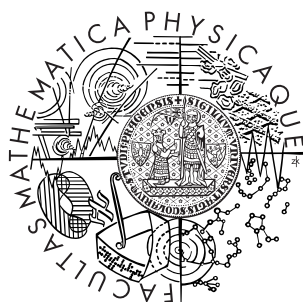


Univerzita Karlova v Praze
Matematicko-fyzikální fakulta

DIPLOMOVÁ PRÁCE



Katarína Mordinová

Errors in variables

Katedra pravděpodobnosti a matematické statistiky

Vedoucí diplomové práce: Mgr. Zdeněk Hlávka, Ph.D.

Studijní program: matematika, finanční a pojistná matematika

2011

Na tomto mieste by som chcela predovšetkým poďakovať vedúcemu diplomovej práce Mgr. Zdeňkovi Hlávkoví, Ph.D. za rady a pripomienky pri vypracovaní tejto práce. Moje poďakovanie patrí aj mojim rodičom za podporu počas celého štúdia.

Prehlasujem, že som svoju diplomovú prácu napísala samostatne a výhradne s použitím citovaných prameňov. Súhlasím so zapožičiavaním práce a jej zverejňovaním.

V Prahe dňa 10. apríla 2011

Obsah

Úvod	5
1 Základné pojmy	6
1.1 Podklady	8
1.1.1 Určenie problému	9
1.1.2 Selekcia potenciálne dôležitých premenných	9
1.1.3 Zber dát	9
1.1.4 Špecifikácia modelu	10
1.1.5 Potvrdenie modelu a posudok	11
1.1.6 Použitie vybraného modelu pri riešení daného problému	12
2 Lineárny regresný model	13
2.1 Odhady parametrov v lineárnom regresnom modeli	14
2.2 Odhad vektora stredných hodnôt	17
2.3 Reziduum	19
2.3.1 Odhad vektora regresných koeficientov β	20
2.4 Normálny lineárny model	30
3 Errors in variables	33
3.1 Lineárne EIV modely	33
3.2 Nelineárne EIV modely	36
3.3 Čiastočne lineárne EIV modely	36
4 Aplikácia EIV modelov v praxi	39
4.1 Výskum rakoviny prsu	39
4.1.1 Rozbor podkladových dát	40
4.1.2 Analýza Errors-in-variables	42
Záver	45

Názov práce: Errors in variables

Autor: Katarína Mordinová

Katedra: Katedra pravděpodobnosti a matematické statistiky

Vedúci bakalárskej práce: Mgr. Zdeněk Hlávka, Ph.D.

e-mail vedúceho: Zdenek.Hlavka@mff.cuni.cz

Abstrakt: Práca pojednáva o regresných modeloch známych ako Errors in variables. V úvodnej kapitole tejto práce definujeme základné pojmy, ktoré v práci používame a základné vzťahy, ktoré s regresnou analýzou súvisia. V druhej kapitole sa venujeme lineárnemu regresnému modelu a jeho vlastnostiam. V tretej kapitole sa zaoberáme jednotlivými typmi modelov errors in variables. V záverečnej kapitole tejto práce si ukážeme možnú aplikáciu modelov v medicíne.

Kľúčové slová: regresná analýza, errors in variables, lineárny regresný model

Title: Errors in variables

Author: Katarína Mordinová

Department: Department of Probability and Mathematical Statistics

Supervisor: Mgr. Zdeněk Hlávka, Ph.D.

Supervisor's e-mail address: Zdenek.Hlavka@mff.cuni.cz

Abstract: The topic of the diploma thesis is Errors in variables. In the opening chapter, we define basic terms used in the thesis and we introduce the regression analysis and basic relations related to this term. In the second chapter, we attend to linear regression model and its characteristics. In the third chapter, we attend to the errors in variables models. In the last chapter of this thesis we present a possible application in medicine.

Keywords: regression analysis, errors in variables, linear regression model

Úvod

V úvodnej kapitole tejto práce sú zadané základné pojmy, ktoré v práci používame. Keďže sa s pojmom regresia stretávame v rôznych oboroch a oblastiach, je zrejmé, že regresná analýza má početné možnosti aplikácie. Štandardne sa s jej aplikáciou stretávame v oboroch financií, ekonómie, obchodu, práva, štatistiky, medicíny a dalo by sa pokračovať ďalšími vednými obormi. V podstate, sa táto jednoduchá metóda dá využiť pri riešení širokého spektra problémov. S niektorými konkrétnymi aplikáciami sa neskôr budeme zaoberať aj v tejto práci.

Po zavedení základných pojmov prejdeme k základným vzťahom medzi skúmanými náhodnými veličinami. Predovšetkým si tieto vzťahy ukážeme na lineárnom regresnom modeli a normálnom regresnom modeli. Tretia kapitola je následne venovaná modelom errors in variables a ich podrobnejším popisom. Finálna, štvrtá kapitola, sa zaoberá aplikáciou modelov errors in variables v praxi. Pre potreby tejto práce sme zvolili aplikáciu errors in variables v medicíne, vo výskume rakoviny pľs.

Kapitola 1

Základné pojmy

Regresia respektívne regres je výraz pochádzajúci z latinského výrazu *regressus*, čo obecné znamená návrat, pohyb späť. S pojmom regresia sa stretávame v rôznych oboroch a oblastiach, keďže regresná analýza má početné možnosti aplikácie. Čiastočný, no určite nie úplný zoznam by zahrnoval financie, ekonómiu, obchod, právo, štatistiku, medicínu a dalo by sa pokračovať ďalšími vednými obormi. V podstate, sa táto jednoduchá metóda dá využiť pri riešení širokého spektra problémov. S niektorými konkrétnymi aplikáciami sa neskôr budeme zaoberať aj v tejto práci.

Samotný pojem regresie zaviedol do štatistiky britský vedec Francis Galton okolo roku 1880. Názov regresnej analýzy, ktorou sa budeme zaoberať, sa odvodzuje práve od jeho článku: *Regression towards mediocrity in hereditary stature*, z roku 1886 [Gal86].

V ktorom sa zaoberal vyšetrovaním závislosti priemernej výšky potomkov na výške rodičov. Avšak so snahou hľadať vzájomný vzťah náhodných veličín sa stretávame už oveľa skôr, napríklad v prácach Galilea Galileiho, R.J.Boscowitcha, Laplaca či Gaussa.

Regresná analýza sa postupom času stala jednou z najviac rozšírených a používaných metód pre analyzovanie dát. Hlavnou príčinou tejto skutočnosti je najmä to, že sa jedná o relatívne jednoduchú a zrozumiteľnú metódu pre skúmanie funkčného vzťahu medzi premennými. Tento vzťah je

vyjadrený vo forme rovnice alebo modelu, ktoré spájajú takzvané závislé premenné a jednu či viac vysvetľujúcich premenných.

Vzhľadom k tomu, že k tejto téme existuje obrovské množstvo literatúry, je občas dosť ťažké sa vyznať v označení jednotlivých premenných vyskytujúcich sa v teórii regresnej analýzy. Najjednoduchší spôsob ako začať je uvažovať na začiatok vzťah medzi jednou závislou a jednou vysvetľujúcou premennou. Použitie výrazu závislá premenná v tomto prípade znamená, že táto premenná je niakym spôsobom naviazaná na hodnotu vysvetľujúcej premennej.

Teda, že napríklad ak vysvetľujúcu premennú vynásobím vhodným koeficientom, dostanem hodnotu závislej premennej. A práve z tejto úvahy vychádzal už spomenutí Francis Galton pri svojom experimente s meraní synov a otcov. Jeho experiment môžeme zjednodušene popísať nasledovne: uvažujme n rodín, v ktorých máme nameranú výšku otca a syna. Pre prehľadnosť začnime s prvou rodinou. Výška otca je pre nás vysvetľujúca premenná, je to náhodná veličina a označíme si ju x_1 . Výška syna je pre nás závislá premenná, ktorá je tiež náhodnou veličinou a označme si ju y_1 . Vzťah medzi x_1 a y_1 potom môžeme popísať rovnicou v tvare:

$$y_1 = \beta_0 + \beta_1 x_1 + \varepsilon_1.$$

V tejto rovnici β_0 a β_1 predstavujú koeficienty, prostredníctvom ktorých vyjadrujeme väzobnú závislosť premenných x_1 a y_1 . Pri meraní, môže dôjsť k chybe a tento fakt je teda nutné reflektovať zahrnutím náhodnej zložky ε_1 do rovnice.

Rovnakou rovnicou možno vyjadriť vzťah medzi výškou otca a syna v každej zo skúmaných n rodín. Tým sa dostávame k ďalšiemu značeniu. Symbolom Y si označíme náhodný vektor, ktorého zložky sú tvorené náhodnými veličinami y_1, \dots, y_n . Teda je to vektor, ktorý nám udáva namerané výšky synov v jednotlivých rodinách. Obdobne potom označme X náhodný vektor, ktorého zložky sú tvorené náhodnými veličinami x_1, \dots, x_n a udáva nám namerané výšky otcov v jednotlivých rodinách. Rovnica závislosti Y a X má potom

tvar:

$$Y = \beta_0 + \beta_1 X + \varepsilon.$$

Tu už ale parametre β_0 a β_1 predstavujú náhodný vektor. Tak isto ε predstavuje už náhodný vektor chýb merania a je tvorený zložkami $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$.

V literatúre zaoberajúcej sa touto problematikou potom môžeme nájsť rôzne pomenovania pre premenné. Závislé premenné y_i nájdeme označenú ako regresandy a odozvy. Vysvetľujúce premenné x_i ako prediktory, regresory alebo nezávislé premenné. Avšak názov nezávislé premenné je najmenej preferovaný a to hlavne z dôvodu, že samotné vysvetľujúce premenné sú len málokedy vzájomne nezávislé. Teda tento názov môže byť mierne zavádzajúci, a v tejto práci sa mu radšej vyhneme.

1.1 Podklady

Pozrime sa teraz obecné na podklady, s ktorými budeme pracovať v tejto práci. Základom pre nás sú niaké dostupné data. Pri skúmaní a vyhodnocovaní týchto dát je vhodné postupovať na základe overenej šablóny, ktorú by sme mohli nazvať krokmi regresnej analýzy. V stručnosti nám totiž poskytuje návod ako postupovať pri výbere dát i modelu, na základe ktorého budeme tieto data vyhodnocovať, a tiež nám sprostredkuje určitý obraz o tom, k akým výsledkom sa na základe zvoleného modelu môžeme no finále dopracovať.

Kroky regresnej analýzy môžeme rozdeliť do nasledujúcich častí:

1. Určenie problému
2. Selekcia potenciálne dôležitých premenných
3. Zber dát
4. Špecifikácia modelu
5. Potvrdenie modelu a posudok
6. Použitie vybraného modelu pri riešení daného problému

1.1.1 Určenie problému

Základným kameňom regresnej analýzy je formulácia daného problému. To zahŕňa definovanie problému a vymedzenie otázok, ktoré majú byť zodpovedané analýzou. Tento krok je prvý a snád najdôležitejší krok v regresnej analýze. Je tak dôležitý, pretože ním sú silne ovplyvnené aj nasledujúce kroky, a chyby v ňom môžu viesť k výberu irelevantnej množiny premenných alebo k nevhodnej voľbe štatistického modelu analýzy.

1.1.2 Selekcia potenciálne dôležitých premenných

Ďalším krokom po vymedzení problému je výber množiny premenných, ktoré sú potrebné pre vysvetlenie alebo predikciu odozvy y . Je nutné sa zamerať na tie data, ktoré nám skutočne vplyvajú na odozvu. To nie je vždy tak jednoduché ako to na prvý pohľad vyzerá.

1.1.3 Zber dát

Zber dát prebieha z prostredia, ktoré je podrobené nášmu štúdiu. Niekedy sú data zozbierané na základe kontrolovaného nastavenia, takže faktory, ktoré nie sú naším primárnym záujmom možno udržiavať na konštantnej úrovni. Menej častý je prístup kedy sú data zozbierané na základe neexperimentálnych podmienok, kedy je len veľmi málo možné zasahovať do ich kontroly. V oboch prípadoch ale zozbierané data pozostávajú z pozorovaní o n zložkách. Každé z týchto n pozorovaní je zaznamenané v tabuľke 1, ktorá má štandardne nasledujúci tvar:

číslo pozorovania	Y	X_1	\dots	X_p	ε
1	y_1	x_{11}	\dots	x_{1p}	ε_1
\cdot	\cdot	\cdot	\cdot	\cdot	\cdot
\cdot	\cdot	\cdot	\cdot	\cdot	\cdot
\cdot	\cdot	\cdot	\cdot	\cdot	\cdot
\cdot	\cdot	\cdot	\cdot	\cdot	\cdot
n	y_n	x_{n1}	\dots	x_{np}	ε_n

Tabuľka 1: Štandardný zápis dát

Stĺpce v tabuľke reprezentujú premennú, zatiaľ čo riadky prislúchajú pozorovaniam, ktoré sú tvorené z množiny $(1 + p)$ hodnôt pre každú zložku. Teda dostávame z toho jednu hodnotu pre odozvu a jednu hodnotu pre každú hodnotu p prediktorov. Pre náš predošlý príklad merania výšky synov by sme teda mali tabuľku zloženú iba zo stĺpcov pre Y a X_1 .

Zápis prediktora x_{ij} prislúcha napozorovanej i -tej hodnote j -tej vysvetľujúcej premennej resp. prediktoru. Prvý index nám hovorí o poradí pozorovania a druhý o poradí vysvetľujúcej premennej. Na základe tohto môžeme pripísať lineárny regresný model pre i -té pozorovanie v tvare:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i,$$

kde $i = 1, \dots, n$ a $j = 1, \dots, p$.

1.1.4 Špecifikácia modelu

Podoba modelu, ktorá súvisí s odozvou a jej vzťahom k množine prediktorov, môže byť špecifikovaná na začiatok na základe už v praxi využívaných modeloch, ktoré vznikli na podklade skúseností a znalostí expertov v danom obore nášho štúdia. Predpokladaný model pravdaže môžeme následne modifikovať alebo úplne zmeniť na základe nazberaných dát. Je nutné myslieť na to, že model je v tomto kroku načrtnutý len v obrysoch, pretože stále môže závisieť na neznámych parametroch.

V tomto kroku ďalej potrebujeme určiť podobu funkcie $f(X_1, X_2, \dots, X_p)$.

Táto funkcia môže byť klasifikovaná v dvoch typoch, a to ako lineárna alebo nelineárna.

Všetky nelineárne funkcie, ktoré možno transformovať do lineárnych funkcií sa nazývajú linearizovateľné funkcie. Avšak nie všetky nelineárne funkcie sú linearizovateľné. V literatúre sa dokonca niekedy stretávame s názorom, že iba nelineárne funkcie, ktoré sú nelinearizovateľné sú skutočne nelineárne funkcie.

1.1.5 Potvrdenie modelu a posudok

Potom, ako sme definovali model a zozbierali data, je naším ďalším krokom odhad regresných parametrov na základe nazbieraných dát. Metóda, ktorá sa najčastejšie využíva pre odhad parametrov sa nazýva metóda najmenších štvorcov (*least-squares method*). Za určitých predpokladov dáva táto metóda odhady s požadovanými vlastnosťami. V niektorých prípadoch, ak je jeden alebo viac predpokladov tejto metódy porušených, sa využívajú aj iné metódy ako napríklad metóda maximálnej vierohodnosti (*maximum likelihood method*), hrebeňová metóda (*ridge method*) alebo metóda hlavných komponent (*principal components method*).

Následne potom aplikujeme vybranú metódu odhadu parametrov na regresné parametre. Odhady regresných parametrov $\beta_0, \beta_1, \dots, \beta_p$ sa označujú $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$. Odhadnutá regresná rovnica má potom tvar:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \dots + \hat{\beta}_p x_{ip}. \quad (1.1)$$

Strieška nad parametrom označuje, že sa jedná o odhad parametru a toto označenie odhadu budeme využívať aj v celej práci. Čo sa týka hodnoty \hat{y}_i , tú nazývame vyrovnaná hodnota (*fitted value*) i -teho pozorovania.

Správnosť štatistickej metódy, akou je i regresná analýza, závisí na určitých predpokladoch. Predpoklady sa zvyčajne vyslovujú o dátach a o modeli. Presnosť rozboru a záverov je odvodená od analýzy zásadne závislej na správnosti týchto predpokladov. Predtým, ako použijeme rovnicu (1.1) za akýmkoľvek

účelom, potrebujeme si najprv overiť či budú špecifické predpoklady platiť. Snažíme sa teda odpovedať na nasledujúce otázky:

1. Aké sú požadované predpoklady?
2. Je každý z týchto predpokladov odôvodnený?
3. Čo budeme robiť ak jeden či viacero predpokladov nebude platiť?

V tomto kroku sa kladie dôraz na to, aby kontrola platnosti predpokladov bola vykonaná predtým ako sú vyvedené akékoľvek závery na základe nášho rozboru. Na regresnú analýzu sa totiž pohliada ako na cyklický proces, ktorého výstupy sa používajú pre diagnostiku, potvrdenie, kritiku a možnú modifikáciu vstupov.

1.1.6 Použitie vybraného modelu pri riešení daného problému

Explicitné určenie regresnej rovnice je najdôležitejší produkt analýzy. Je to zhrnutie vzťahu medzi Y a X_1, X_2, \dots, X_p . Rovnicu možno použiť pre rôzne účely. Možno ju využiť pri zisťovaní dôležitosti individuálneho prediktora pre analyzovanie efektov metódy, čo zahŕňa meniace sa hodnoty prediktorov, alebo pre predpoveď hodnôt odozvy pre danú množinu prediktorov. Aj keď je regresná rovnica finálny produkt, je tu aj množstvo ďalších vedľajších produktov. Vidíme regresnú analýzu ako množinu analytických techník, ktoré sa používajú pre lepšie pochopenie vnútorného vzťahu medzi premennými a daným prostredím. Úlohou regresnej analýzy je naučiť sa čo najviac je možné o prostredí na základe nahromadených dát. Zdôrazňuje sa, že to čo je odkryté počas procesu formulácie rovnice, môže byť niekedy tak cenné a informatívne ako samotná finálna rovnica.

Kapitola 2

Lineárny regresný model

Uvažujme náhodný vektor $Y = (y_1, \dots, y_n)$, ktorého zložky tvoria nezávislé náhodné veličiny y_1, \dots, y_n so spojitým rozdelením. Stredné hodnoty týchto veličín možno vyjadriť v tvare:

$$E y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}, \quad (2.1)$$

kde $i = 1, \dots, n$. Vektor $\beta = (\beta_0, \beta_1, \dots, \beta_k)$ je tvorený parametrami $\beta_0, \beta_1, \dots, \beta_k$, ktoré nazývame regresnými koeficientmi. Jedná sa o neznáme konštanty.

Veličiny x_{ij} sú známe konštanty a možno ich usporiadať do regresnej matice \mathbf{X} typu $(n, k + 1)$:

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1k} \\ 1 & x_{21} & \dots & x_{2k} \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ 1 & x_{n1} & \dots & x_{nk} \end{pmatrix} \quad (2.2)$$

Matica \mathbf{X} je veľmi dôležitá pri odhade vektoru β , kedy konkrétne lineárne kombinácie jej stĺpcov predstavujú strednú hodnotu náhodného vektoru Y , a teda odhad vektoru β .

Ďalej predpokladajme, že platí:

$$\text{Var } y_i = \sigma^2, \quad (2.3)$$

pre $i = 1, \dots, n$, kde $\sigma^2 > 0$ je neznámy parameter, ktorý označujeme ako rozptyl náhodných veličín y_i . Model (2.1) teda môžeme prepísať v tvare:

$$Y \sim (\mathbf{X}\beta, \sigma^2\mathbf{I}), \quad (2.4)$$

kde \mathbf{I} je jednotková matica a výraz $\sigma^2\mathbf{I}$ dáva variančnú maticu vektora Y .

$$\text{Var } Y = \begin{pmatrix} \text{var } y_1 & \text{cov } (y_1, y_2) & \dots & \text{cov } (y_1, y_n) \\ \text{cov } (y_2, y_1) & \text{var } y_2 & \dots & \text{cov } (y_2, y_n) \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \text{cov } (y_n, y_1) & \text{cov } (y_n, y_2) & \dots & \text{var } y_n \end{pmatrix} \quad (2.5)$$

Variančná matica v tomto tvare (2.5) potom udáva rovnaký rozptyl a nekorelovanosť jednotlivých zložiek y_1, \dots, y_n náhodného vektora Y .

Hodnosť matice \mathbf{X} označme ako $h(\mathbf{X}) = p$. V prípade, že $h(\mathbf{X}) = (k + 1)$ hovoríme, že model (2.4) je model s úplnou hodnosťou. Ak $k = 1$ hovoríme o jednoduchej lineárnej regresii, keďže uvažujeme len jednu nezávislú premennú x . Ak $k > 1$ hovoríme o mnohonásobnej lineárnej regresii. Ak $h(\mathbf{X}) < (k + 1)$, potom sa jedná o model s neúplnou hodnosťou. Obecné predpokladáme, že $h(\mathbf{X}) > 0$.

Model (2.4) možno ekvivalentne zapísať aj pomocou náhodnej zložky $\varepsilon \sim (0, \sigma^2\mathbf{I})$ v tvare:

$$Y = \mathbf{X}\beta + \varepsilon, \quad (2.6)$$

kde $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ je náhodný vektor chýb.

2.1 Odhady parametrov v lineárnom regresnom modeli

Z toho ako sme si zaviedli lineárny regresný model v predošlej podkapitole je jasné, že nás pri skúmaní dát zaujímajú hlavne odhady neznámych parametrov. Teda regresných parametrov $\beta_0, \beta_1, \dots, \beta_k$ odhad rozptylu σ^2 a odhad vektora

stredných hodnôt $\mathbf{X}\beta$ v modeli (2.4).

Aby sme mohli bližšie popísať akým spôsobom odhadujeme požadované parametre, je vhodné si najprv uviesť niekoľko tvrdení, z ktorých budeme vychádzať.

Pri práci s maticami s neúplnou hodnotou sa nám určite zide tzv. *Pravidlo piatich matíc*, ktorého dôkaz môžeme nájsť napríklad v publikácii [And05]. Toto pravidlo môžeme napísať v nasledujúcom tvare:

$$\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X} = \mathbf{X} \quad (2.7)$$

Úpravou potom dostávame tvar:

$$(\mathbf{I} - \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T)\mathbf{X} = \mathbf{0} \quad (2.8)$$

kde označme:

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \quad (2.9)$$

$$\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \quad (2.10)$$

Odkiaľ máme:

$$\mathbf{H}\mathbf{X} = \mathbf{X}$$

$$\mathbf{M}\mathbf{X} = \mathbf{0}.$$

Definícia 1 (Idempotentná matica). *Štvorcovú maticu \mathbf{A} nazývame idempotentnou, ak platí: $\mathbf{A}\mathbf{A} = \mathbf{A}$. Hodnota idempotentnej matice je rovná jej stope, tj. $h(\mathbf{A}) = \text{tr } \mathbf{A}$.*

Definícia 2 (Lineárny priestor). *Uvažujme pevne dané vektory $x_1, \dots, x_k \in R^n$ tj. do n -rozmerného euklidovského priestoru. Podpriestor R^n obsahujúci práve všetky lineárne kombinácie vektorov x_1, \dots, x_k označí ako:*

$$M(x_1, \dots, x_k) = \{x \in R^n : x = \sum_{i=1}^k \alpha_i x_i, \alpha_i \in R^1; i = 1, \dots, k\}$$

nazývame lineárnym priestorom generovaným vektormi x_1, \dots, x_k . Ak sú vektory x_1, \dots, x_k stĺpce matice \mathbf{X} potom píšeme $M(\mathbf{X})$.

Definícia 3 (Ortogonalný doplnok). *Nech je niaky lineárny podpriestor priestoru $\Omega \subset R^n$, tj. s každými dvoma prvkami tam patrí tiež ich lineárna kombinácia. Množina:*

$$\Omega^\perp = \{x \in R^n : (x, y) = 0, \forall y \in \Omega\}$$

je tiež lineárny podpriestor R^n a nazýva sa ortogonálnym doplnkom podpriestoru $\Omega \subset R^n$.

Ďalej platí, že dva lineárne podpriestory Ω_1 a Ω_2 z R^n majú rovnaký ortogonálny doplnok, práve keď sú totožné. Pretože platí:

$$(\Omega^\perp)^\perp = \Omega.$$

Platí, teda že $M(\mathbf{X}\mathbf{X}^T) = M(\mathbf{X})$.

Definícia 4 (Projekcia do podpriestoru). *Nech je Ω podpriestor euklidovského priestoru R^n . Ortogonálnym priemetom vektoru $x \in R^n$ do podpriestoru Ω nazveme taký vektor $r_\Omega(x)$, pre ktorý platí že:*

$$\begin{aligned} r_\Omega(x) &\in \Omega \\ x - r_\Omega(x) &\in \Omega^\perp. \end{aligned}$$

Potom platia nasledujúce tvrdenia, ktorých dôkazy nájdeme v [Zvá89].

Tvrdenie 2.1. :

Priemet $r_\Omega(x)$ je daný jednoznačne.

Tvrdenie 2.2. :

Vektor $x - r_\Omega(x)$ je ortogonálnym priemetom vektora x do podpriestoru Ω^\perp .

Tvrdenie 2.3. :

Platí:

$$\begin{aligned} x \in \Omega &\leftrightarrow r_{\Omega}(x) = x \\ x \in \Omega^{\perp} &\leftrightarrow r_{\Omega}(x) = 0. \end{aligned} \quad (2.11)$$

Tvrdenie 2.4. :

Pre $\forall x \in R^n$ platí:

$$\|x\|^2 = \|r_{\Omega}(x)\|^2 + \|x - r_{\Omega}(x)\|^2. \quad (2.12)$$

Tvrdenie 2.5. :

Nech $x \in R^n$. Ak $z \in \Omega$ také, že pre $\forall y \in \Omega$ platí:

$$\|x - z\|^2 \leq \|x - y\|^2 \quad (2.13)$$

potom tiež $z = r_{\Omega}(x)$.

2.2 Odhad vektora stredných hodnôt

Pozrime sa najprv na odhad vektoru stredných hodnôt $E Y$ v modeli $Y \sim (\mathbf{X}\beta, \sigma^2\mathbf{I})$. Označme si $\mu = \mathbf{X}\beta$. Z definície 1 je zrejmé, že $\mu \in M(\mathbf{X})$. Teda možno vektor odhadnúť μ vektorom \hat{Y} z priestoru $M(\mathbf{X})$, ktorý je k vektoru Y najbližší. Keďže platí tvrdenie (2.13), je \hat{Y} priemetom vektora do $M(\mathbf{X})$. Projekčná maticu si označme \mathbf{H} a platí:

$$\hat{Y} = \mathbf{H}Y.$$

Z (2.9) máme, že:

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$$

Veta 2.1. :

V modeli $Y \sim (\mathbf{X}\beta, \sigma^2\mathbf{I})$ je vektor \hat{Y} najlepším nestranným lineárnym odhadom vektoru stredných hodnôt $\mu = \mathbf{X}\beta$. Je daný jednoznačne a platí:

$$\mathbf{Var} \hat{Y} = \sigma^2\mathbf{H}.$$

Dôkaz. Keďže projekčná matica \mathbf{H} je idempotentná, nestrannosť \hat{Y} plynie z toho, že \mathbf{H} je projekčná matica na $M(\mathbf{X})$. Potom: $\mathbf{Var} \hat{Y} = \mathbf{Var} \mathbf{H}Y = \mathbf{H}\sigma^2\mathbf{I}\mathbf{H} = \sigma^2\mathbf{H}$.

Z vlastností idempotentnej matice máme, že: $\mathbf{H}\mathbf{X} = \mathbf{X}$. Teda:

$$E \hat{Y} = E \mathbf{H}Y = \mathbf{H}\mathbf{X}\beta = \mathbf{X}\beta = \mu.$$

Overme ešte jednoznačnosť odhadu \hat{Y} . Uvažujme niaky lineárny nestranný odhad vektoru μ . Označme ho $\tilde{Y} = \mathbf{A}Y + a$. Pre $\forall\beta$ spĺňa matica \mathbf{A} typu (n, n) rovnosť:

$$E \tilde{Y} = \mathbf{A}\mathbf{X}\beta + a = \mathbf{X}\beta.$$

Keďže je možné aby vektor $\beta = 0$, musí byť teda $a = 0$ a zároveň musí platiť $\mathbf{A}\mathbf{X} = \mathbf{X}$. Teda aj $\mathbf{A}\mathbf{H} = \mathbf{H}$.

K porovnaniu nestranných odhadov vektorového parametru používame ich variančné matice. Platí, že ak sú \hat{Y} a \tilde{Y} dva nestranné odhady vektoru μ , potom je odhad \hat{Y} lepší ak je matica $(\mathbf{Var} \tilde{Y} - \mathbf{Var} \hat{Y})$ pozitívne semidefinitná. Teda, že pre \forall vektor $q \in R^n$ je: $\mathbf{Var} (q^T \hat{Y}) \leq \mathbf{Var} (q^T \tilde{Y})$.

$$\begin{aligned} \mathbf{Var} \tilde{Y} &= \mathbf{Var} \mathbf{A}Y \\ &= \sigma^2\mathbf{A}\mathbf{A}^T \\ &= \sigma^2[\mathbf{H} + (\mathbf{A} - \mathbf{H})][\mathbf{H} + (\mathbf{A} - \mathbf{H})] \\ &= \sigma^2\mathbf{H} + \sigma^2(\mathbf{A} - \mathbf{H})(\mathbf{A} - \mathbf{H})^T, \\ \mathbf{Var} \hat{Y} &= \sigma^2\mathbf{H}. \end{aligned}$$

□

2.3 Reziduum

Vektor $v = Y - \hat{Y}$ sa nazýva vektor reziduí resp. reziduum. Porovnáva nám napozorované hodnoty závislej premennej s odhadom jej stredných hodnôt. Z tvrdenia 2 plynie, že je priemetom do priestoru $M(\mathbf{X})^\perp$, ktorý nazývame reziduálnym priestorom.

Príslušná projekčná matica má tvar:

$$\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T.$$

Reziduálny súčet štvorcov označme:

$$\mathbf{RSS} = \|v\|^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

Je to štvorec vzdialenosti vektorov Y a \hat{Y} , ktorý meria ich nezhodu.

Reziduálny rozptyl potom dostávame ako:

$$\mathbf{S}^2 = \frac{\mathbf{RSS}}{n - p}.$$

Veta 2.2 (Veta o reziduách). :

V lineárnom modeli $Y \sim (\mathbf{X}\beta, \sigma^2\mathbf{I})$ platí:

(i) $v = \mathbf{M}Y = \mathbf{M}\varepsilon$

(ii) $v \sim (\mathbf{0}, \sigma^2\mathbf{M})$

(iii) $\mathbf{RSS} = \varepsilon^T\mathbf{M}\varepsilon$

(iv) $E \mathbf{RSS} = (n - p)\sigma^2$

$$(v) E \mathbf{S}^2 = \sigma^2$$

$$(vi) \mathbf{X}^T v = \mathbf{0}.$$

Dôkaz. Tvrdenie 1 vyplýva z toho, že $v = \mathbf{M}Y = \mathbf{M}(\mathbf{X}\beta + \varepsilon) = \mathbf{M}\varepsilon$, pretože platí: $\mathbf{M}\mathbf{X} = \mathbf{0}$.

Keďže $\varepsilon \sim (0, \sigma^2 \mathbf{I})$ a $v = Y - \hat{Y}$, teda $v = \mathbf{M}Y \sim (\mathbf{0}, \sigma^2 \mathbf{M})$, čím sme dostali tvrdenie 2.

Pre tvrdenie 3 vychádzame z toho, že $\mathbf{RSS} = \varepsilon^T \mathbf{M}\varepsilon$, kde $\mathbf{M} = \mathbf{M}^T \mathbf{M}$ je idempotentná matica.

$$\text{Potom } \mathbf{RSS} = \|v\|^2 = \|\mathbf{M}\varepsilon\|^2 = (\mathbf{M}\varepsilon)^T (\mathbf{M}\varepsilon) = \varepsilon^T \mathbf{M}^T \mathbf{M}\varepsilon = \varepsilon^T \mathbf{M}\varepsilon.$$

Tvrdenie 4 dostaneme úpravou: $E \mathbf{RSS} = E \|v\|^2 = E v^T v = \text{tr } E v^T v = \text{tr } \text{Var } v = \sigma^2 \text{tr } \mathbf{M} = \sigma^2 \text{tr}(\mathbf{I} - \mathbf{H}) = \sigma^2(n - p)$.

Keďže platí, že: $h(\mathbf{H}) = h(\mathbf{X})$ a stopa tr idempotentnej matice je rovná jej hodnosti.

$$\text{Pre tvrdenie 5 potom máme: } E \mathbf{S}^2 = E \left(\frac{\mathbf{RSS}}{n-p} \right) = \frac{E \mathbf{RSS}}{(n-p)} = \sigma^2 \frac{n-p}{n-p} = \sigma^2.$$

Tvrdenie 6 plynie z toho, že: $\mathbf{M}\mathbf{X} = \mathbf{0}$ a $v = \mathbf{M}\varepsilon$.

□

2.3.1 Odhad vektora regresných koeficientov β

Pri odhade vektora regresných parametrov $\beta = (\beta_0, \beta_1, \dots, \beta_k)$ vychádzame z hodnosti matice \mathbf{X} . Odhad vektoru β totiž vyjadruje strednú hodnotu náhodného vektora Y ako konkrétnu lineárnu kombináciu stĺpcov matice \mathbf{X} .

Ak má matica \mathbf{X} plnú hodnosť, potom možno každý vektor z $M(\mathbf{X})$ vyjadriť pomocou lineárnej kombinácie stĺpcov matice \mathbf{X} nekonečne mnoho spôsobmi.

Jedna z obecně najviac užívaných metód pri odhade regresných parametrov je metóda najmenších štvorcov. Za jej vznikom stáli nezávisle na sebe Carl Friedrich Gauss(1795) a Adrien Marie Legendre(1805).

Prvotné aplikácie tejto metódy boli v astronómii a geodetike. S prvým verejným výtiskom sa stretávame v roku 1805 v Legendreho knihe o určovaní obežných dráh komét.

Pre zjednodušenie si uveďme príklad z praxe, ktorý môžeme nájsť v knihe *Alternative methods of regression* od D. Birkes, Y. Dodge [Dod93].

Pri experimente sa skúmal vzťah medzi dvoma procedúrami merania obsahu kyseliny v chemikálii. Prvá metóda, založená na extrakcii a vážení bola drahšia. Zatiaľ čo druhá, založená na titracii, je oveľa lacnejšia.

Základom pokusu teda bolo zistiť vzájomný vzťah týchto dvoch metód použitím regresie. Pričom ak by vzťah existoval, znamenalo by to v budúcnosti nahradenie drahšej procedúry tou lacnejšou.

Obe procedúry boli použité na 20-tich vzorkách chemikálie, čím sa dospelo k nameraným hodnotám kyseliny v tabuľke 2.

Hodnota y_i predstavuje výsledok i -teho merania drahej metódy a hodnota x_i predstavuje výsledok i -teho merania lacnej metódy merania obsahu kyseliny. Použitím týchto dát sa snažíme získať rovnicu, ktorá by vyjadrovala vzájomný vzťah oboch procedúr.

Pozrime sa najprv na lineárny regresný model, daný rovnicou:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i.$$

Pre tento prípad uvažujme len jednu vysvetľujúcu premennú. Model s týmto nastavením potom označujeme ako jednoduchý lineárny regresný model.

číslo vzorky	Drahá procedura (Y)	Lacná procedura (X)
1	76	123
2	70	109
3	55	62
4	71	104
5	55	57
6	48	37
7	50	44
8	66	100
9	41	16
10	43	28
11	82	138
12	68	105
13	88	159
14	58	75
15	64	88
16	88	164
17	89	169
18	88	167
19	84	149
20	88	167

Tabuľka 2: obsah kyseliny

Rovnica je v tvare:

$$Y = \beta^T X + \varepsilon.$$

Ďalším krokom je určenie odhadov pre regresné koeficienty β_0 a β_1 . Keď následne použijeme lacnejšiu procedúru X na zmeranie množstva kyseliny vo vzorku, potom môžeme pre drahšiu procedúru Y písať odhad i -teho merania v tvare:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i.$$

Pre odhad parametrov β_0 a β_1 použijeme metódu najmenších štvorcov.

Uvažujeme nameranú závislosť ako súbor dvojíc $(x_1, y_1); (x_2, y_2); \dots; (x_n, y_n)$.

Týmito bodmi prekladáme funkciu v tvare $y = f(x, \beta_0, \beta_1, \dots, \beta_p)$, pričom hľadáme štatistické odhady $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ jej parametrov $\beta_0, \beta_1, \dots, \beta_p$.

Základné kritérium je, aby bol súčet štvorcov odchýliek empirických hodnôt y_i , označujeme ho S , minimálny od vyrovnaných hodnôt $y = f(x_i, \beta_0, \beta_1, \dots, \beta_p)$, tj. od teoretickej regresnej funkcie. Teda mala by platiť rovnosť:

$$S(\beta_0, \beta_1, \dots, \beta_p) = \sum_{i=1}^n (y_i - f(x_i, \beta_0, \dots, \beta_p))^2$$

$$\hat{\beta} = \operatorname{agrm} \min_{\beta} \sum_{i=1}^n (y_i - f(x_i, \beta_0, \dots, \beta_p))^2$$

Kedže uvažujeme lineárny regresný model, obmedzíme si tvar funkcie f na lineárnu regresnú funkciu v tvare:

$$f(x_i, \beta_0, \dots, \beta_p) = \beta_0 f_0(x) + \beta_1 f_1(x) + \beta_2 f_2(x) + \dots + \beta_p f_p(x)$$

Pre existenciu minima funkcie S je nutná podmienka:

$$\frac{\partial S}{\partial \beta_j} = 0$$

pre $j = 0, 1, \dots, p$.

Pre parameter β_j dostávame:

$$\sum_{i=1}^n f_j(x_i) f_0(x_i) \beta_0 + \dots + \sum_{i=1}^n f_j(x_i) f_p(x_i) \beta_p = \sum_{i=1}^n y_i f_j(x_i)$$

kde označíme:

$$\sum_{i=1}^n f_j(x_i) f_h(x_i) = a_{jh}$$

$$\sum_{i=1}^n y_i f_j(x_i) = a_j$$

Dostávame sústavu rovníc:

$$\begin{aligned} a_{00}\beta_0 + a_{01}\beta_1 + \dots + a_{0p}\beta_p &= a_0 \\ a_{10}\beta_0 + a_{11}\beta_1 + \dots + a_{1p}\beta_p &= a_1 \\ &\cdot \\ &\cdot \\ &\cdot \\ a_{p0}\beta_0 + a_{p1}\beta_1 + \dots + a_{pp}\beta_p &= a_p \end{aligned}$$

Jej vyriešením dostávame hľadané odhady $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$.

Vráťme sa teraz späť k nášmu meraniu hodnôt kyseliny. Odhady $\hat{\beta}_0$ a $\hat{\beta}_1$ majú potom tvar:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

a

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

Kde $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ a $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ predstavujú priemery.

Pre nami namerané data potom dostávame hodnoty odhadov $\hat{\beta}_0 = 35.46$ a $\hat{\beta}_1 = 0.3216$. Vhodná regresná priamka pre má experiment má teda tvar:

$$\hat{y}_i = 35.46 + 0.3216x_i.$$

Označme si $\hat{\beta}$ odhad vektora regresných parametrov β . Z definície priemetu

platí, že $\mathbf{X}\hat{\beta} \in M(\mathbf{X})$. V zhládom k predošlému, nám $\hat{\beta}$ predstavuje ľubovoľné riešenie sústavy $\mathbf{X}\hat{\beta} = \hat{Y}$. Čo možno prepísať ako:

$$\mathbf{X}\hat{\beta} + v = Y, \quad (2.14)$$

keďže $v = Y - \hat{Y}$.

Z vety o reziduách máme, že $\mathbf{X}^T v = \mathbf{0}$, teda (2.14) možno prepísať v tvare $\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\hat{\beta}) = \mathbf{0}$. To je ekvivalentné s normálnou rovnicou :

$$\mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{X}^T Y. \quad (2.15)$$

Odkiaľ dostávame vyjadrenie pre $\hat{\beta}$ v tvare:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y. \quad (2.16)$$

Riešenie ale nemusí byť jednoznačné, záleží na hodnosti matice \mathbf{X} .

Veta 2.3. : V modeli $Y \sim (\mathbf{X}\beta, \sigma^2 \mathbf{I})$ s úplnou hodnosťou matice \mathbf{X} je vektor $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y$ jediným riešením normálnej rovnice $\mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{X}^T Y$. Vektor $\hat{\beta}$ je najlepším nestranným lineárnym odhadom vektoru regresných koeficientov β .

Dôkaz. Z vlastností ortogonálneho doplnku platí, že $M(\mathbf{X}\mathbf{X}^T) = M(\mathbf{X})$. To nám dáva $h(\mathbf{X}) = h(\mathbf{X}\mathbf{X}^T)$, teda je $\mathbf{X}^T \mathbf{X}$ regulárna matica v modeli s úplnou hodnosťou.

Normálna rovnica má jediné riešenie v tvare $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y$ a vektor $\hat{\beta}$ je jednoznačnou lineárnou funkciou \hat{Y} , pretože platí: $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{Y}$.

Ďalej podľa vety 1 je najlepším nestranným lineárnym odhadom svojej strednej hodnoty.

□

Veta 2.4. Pre momenty $\hat{\beta}$ platí $E \hat{\beta} = \beta$ a $Var \hat{\beta} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$.

Dôkaz. Keďže predpokladáme, že $Y \sim (\mathbf{X}\beta, \sigma^2\mathbf{I})$ a $h(\mathbf{X}) = p$, z predošlej vety 2.3 teda možno písať: $\hat{\beta} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T Y$.

Teda:

$$E \hat{\beta} = E (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T Y = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}\beta = \beta,$$

a $Var \hat{\beta} = Var (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T Y = \sigma^2(\mathbf{X}^T\mathbf{X})^{-1}$. □

Veta 2.5. *Nech $Y \sim (\mu, \mathbf{W}^-)$, kde \mathbf{W} typu (n, n) je pozitívne definitná matica. Nech matica \mathbf{A} typu (n, n) je pozitívne semidefinitná matica, taká že platí $\mu\mathbf{A} = \mathbf{0}$ a $h(\mathbf{A}) = r > 0$. Potom platí, že ak je matica $\mathbf{A}\mathbf{W}^-$ idempotentná:*

$$Y^T \mathbf{A} Y \sim \chi_r^2.$$

Veta 2.6. *Nech $Y \sim (\mu, \mathbf{W}^-)$, kde \mathbf{W} typu (n, n) je pozitívne definitná matica. Nech $\Omega \subset R_n$ a $\mu \in \Omega_{\mathbf{W}}^\perp$. Ak $dim \Omega = p > 0$ potom platí:*

$$\|\mathbf{P}_\Omega Y\|_{\mathbf{W}}^2 \sim \chi_p^2.$$

Dôkaz. Z definície idempotentnej matice vieme, že projekčná matica \mathbf{P} je idempotentná a matica $\mathbf{W}\mathbf{P}$ je symetrická. Teda možno písať:

$$\|\mathbf{P}_\Omega Y\|^2 = Y^T \mathbf{P}_\Omega^T \mathbf{W} \mathbf{P}_\Omega Y = Y^T \mathbf{W} \mathbf{P}_\Omega Y$$

Matica $\mathbf{W}\mathbf{P}_\Omega$ spĺňa predpoklady vety 2.5, teda štvorec dĺžky priemetu $\mathbf{P}_\Omega Y$ pri skalárnom súčine $(\cdot, \cdot)_{\mathbf{W}}$ má rozdelenie χ_p^2 . □

Veta 2.7. *Nech $Y \sim (\mu, \mathbf{W}^-)$, Ω_1, Ω_2 sú podpriestory R^n , pričom platí že:*

$$\Omega_1 \subset \Omega_{2\mathbf{W}}^\perp.$$

Ak sú $\mathbf{P}_1, \mathbf{P}_2$ projekčné matice na podpriestory Ω_1, Ω_2 pri skalárnom súčine $(\cdot, \cdot)_{\mathbf{W}}$, potom platia tvrdenia:

(i) $\mathbf{P}_1 Y, \mathbf{P}_2 Y$ sú nezávislé náhodné vektory

(ii) $\|\mathbf{P}_1 Y\|_{\mathbf{W}}^2$ a $\|\mathbf{P}_2 Y\|_{\mathbf{W}}^2$ sú nezávislé náhodné veličiny.

Dôkaz. Tvrdenie 1 znamená, že pre ľubovoľné $x, z \in R^n$ platí nasledujúce:

$$\mathbf{0} = (\mathbf{P}_1 x, \mathbf{P}_2 z)_{\mathbf{W}} = x^T \mathbf{P}_1^T \mathbf{W} \mathbf{P}_2 z.$$

Teda je $\mathbf{P}_1^T \mathbf{W} \mathbf{P}_2 = \mathbf{0}$. S využitím symetrie projekčných matíc máme symetriu matíc $\mathbf{W} \mathbf{P}_1$, $\mathbf{W} \mathbf{P}_2$ a:

$$\mathbf{0} = \mathbf{P}_1^T \mathbf{W} \mathbf{P}_2 = \mathbf{P}_1^T \mathbf{W} \mathbf{W}^{-1} \mathbf{W} \mathbf{P}_2 = \mathbf{W} \mathbf{P}_1 \mathbf{W}^{-1} \mathbf{P}_2^T \mathbf{W}.$$

Keďže \mathbf{W} je regulárna, platí aj $\mathbf{P}_1 \mathbf{W}^{-1} \mathbf{P}_2^T = \mathbf{0}$. Teda náhodné vektory $\mathbf{P}_1 Y, \mathbf{P}_2 Y$ sú nezávislé.

Tvrdenie 2 je priamym dôsledkom tvrdenia 1. □

Definícia 5 (Overiteľná lineárna hypotéza). :

Nech pre model $Y \sim (\mathbf{X}\beta, \sigma^2 \mathbf{I})$ je $\mu^0 \in M(\mathbf{X})$ pevne zvolený vektor, $\omega \subset M(\mathbf{X})$ je pevne zvolený podpriestor. Potom požiadavok

$$E Y - \mu^0 \in \omega \tag{2.17}$$

sa nazýva overiteľná lineárna hypotéza.

Veta 2.8. Najlepší nestranný lineárny odhad vektora stredných hodnôt $E Y$ v modeli $Y \sim (\mathbf{X}\beta, \sigma^2 \mathbf{I})$ za platnosti overiteľnej lineárnej hypotézy (2.17) je rovný:

$$\hat{Y}_h = \mu^0 + P_\omega(Y - \mu^0).$$

Dôkaz. Pretože je $\mu^0 \in M(\mathbf{X})$, teda existuje $\beta \in R_k$ také, že $\mu^0 = \mathbf{X}\beta^0$. Označme si $z = Y - \mu^0$ a $\gamma = \beta - \beta^0$. Potom je $z \sim (\mathbf{X}\gamma, \sigma^2 \mathbf{I})$. Ak platí overiteľná lineárna hypotéza (2.17), je $E z \in \omega \subset M(\mathbf{X})$. Teda najlepší nestranný lineárny odhad pre $E z$ je:

$$\hat{z}_h = P_\omega z = P_\omega(Y - \mu^0).$$

Keďže $\hat{Y}_h = \hat{z}_h + \mu^0$, potom je $E Y$ najlepší nestranný lineárny odhad za platnosti hypotézy. □

Veta 2.9. Ak je $E Y - \mu^0 \in \omega$ overiteľná lineárna hypotéza v modeli $Y \sim (\mathbf{X}\beta, \sigma^2\mathbf{I})$ potom platia nasledujúce tvrdenia:

$$(i) \hat{Y}_h = \hat{Y} - d$$

$$(ii) v_h = v + d$$

$$(iii) d^T v = 0$$

$$(iv) \mathbf{RSS}_h = \mathbf{RSS} + \|d\|^2,$$

kde vektor d je priemetom vektora $Y - \mu^0$ do $M(\mathbf{X}) \cap \omega^\perp$, a $\mathbf{d} = \dim M(\mathbf{X}) - \dim \omega = \dim M(\mathbf{X}) \cap \omega^\perp$.

Dôkaz. Uvažujme znovu $z = Y - \mu^0$ a $\gamma = \beta - \beta^0$, $\beta \in R_k$ také, že $\mu^0 = \mathbf{X}\beta^0$ a $z \sim (\mathbf{X}\gamma, \sigma^2\mathbf{I})$. Označme si $\Omega = M(\mathbf{X})$. Keďže P je projekčná matica, z jej vlastností platí, že:

$$\hat{z} - \hat{z}_h = (P_\Omega - P_\omega)z = (P_{\Omega \cap \omega^\perp})z.$$

Keďže $z = Y - \mu^0$, potom dostávame:

$$\hat{Y} - \mu^0 - (\hat{Y}_h - \mu^0) = (P_{\Omega \cap \omega^\perp})(Y - \mu^0) = d.$$

tj. tvrdenie 1. Keď ho odčítame od identity $Y = Y$, dostávame tvrdenie 2. Tvrdenie 3 dostaneme ako dôsledok toho, že $d \in \Omega \cap \omega^\perp$ a $v \in \Omega^\perp$. Potom ale môžeme písať:

$$\mathbf{RSS}_h = \|v_h\|^2 = \|v + d\|^2 = \|v\|^2 + \|d\|^2 = \mathbf{RSS} + \|d\|^2.$$

čo je tvrdenie 4. □

Veta 2.10. Uvažujme konzistentnú sústavu lineárnych rovníc v tvare:

$$\mathbf{A}\beta = a, \tag{2.18}$$

kde \mathbf{A} je daná matica typu (t, k) a $a \in R_t$ je daný vektor.

Potom omedzenie (2.18) je overiteľná lineárna hypotéza v modeli $Y \sim (\mathbf{X}\beta, \sigma^2\mathbf{I})$ práve keď platí:

$$M(\mathbf{A}^T) \subset M(\mathbf{X}^T). \quad (2.19)$$

Potom pre ľubovoľné riešenie $\hat{\beta}$ normálnej rovnice platí:

$$\|d\|^2 = (\mathbf{A}\hat{\beta} - a)^T [\mathbf{A}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{A}^T]^{-1} (\mathbf{A}\hat{\beta} - a)$$

a

$$\mathbf{d} = h(\mathbf{A}).$$

Veta 2.11. Ak $Y \sim (\mathbf{X}\beta, \sigma^2\mathbf{I})$ je model s úplnou hodnotou, potom je každá lineárna hypotéza (2.18) overiteľná. Ak má navyše matica \mathbf{A} lineárne nezávislé riadky, platí:

$$(i) \hat{\beta}_h = \hat{\beta} - (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{A}^T[\mathbf{A}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{A}^T]^{-1}(\mathbf{A}\hat{\beta} - a)$$

$$(ii) \mathbf{RSS}_h = \mathbf{RSS} + (\hat{\beta} - \hat{\beta}_h)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \hat{\beta}_h).$$

Dôkaz. Ak $p = k$, potom platí $M(\mathbf{X}^T) = R_k$ a triviálne platí $M(\mathbf{A}^T) \subset M(\mathbf{X}^T)$. Keďže predpokladáme hodnotu matice $h(\mathbf{A}) = t$, má sústava $\mathbf{A}\beta = a$ vždy riešenie. Vzhľadom k lineárnej nezávislosti riadkov matice \mathbf{A} a pozitívnej definitnosti matice $(\mathbf{X}^T\mathbf{X})^{-1}$ dostávame aj pozitívnu definitnosť a regularitu matice $\mathbf{A}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{A}^T$. Dosadením do $\hat{Y}_h = \hat{Y} - d$ a dostávame:

$$\hat{Y}_h = \mathbf{X}\hat{\beta} - \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{A}^T[\mathbf{A}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{A}^T]^{-1}(\mathbf{A}\hat{\beta} - a) = \mathbf{X}\hat{\beta}_h.$$

tj. tvrdenie 1.

A ďalej využitím predošlej vety môžeme písať:

$$\begin{aligned} \mathbf{RSS}_h - \mathbf{RSS} &= \|d\|^2 \\ &= (\mathbf{A}\hat{\beta} - a)^T [\mathbf{A}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{A}^T]^{-1} (\mathbf{A}\hat{\beta} - a) \\ &= (\hat{\beta} - \hat{\beta}_h)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \hat{\beta}_h). \end{aligned}$$

□

2.4 Normálny lineárny model

Predpokladajme, že Y má mnohorozmerné normálne rozdelenie, tj. $Y \sim N(\mathbf{X}\beta, \sigma^2\mathbf{I})$. Pre rozdelenie vektor chýb ε , potom dostávame, že $\varepsilon \sim N(0, \sigma^2\mathbf{I})$.

Veta 2.12. *Nech $Y \sim N(\mathbf{X}\beta, \sigma^2\mathbf{I})$, potom platia nasledujúce tvrdenia:*

(i) $\hat{Y} \sim N(\mathbf{X}\beta, \sigma^2\mathbf{H})$

(ii) $v \sim N(\mathbf{0}, \sigma^2\mathbf{M})$

(iii) \hat{Y} a v sú nezávislé náhodné vektory

(iv) **RSS** nezávisí na \hat{Y}

(v) $\frac{\mathbf{RSS}}{\sigma^2} \sim \chi_{n-p}^2$

Dôkaz. Tvrdenie 1 je triviálnym dôsledkom rozdelenia $Y \sim N(\mathbf{X}\beta, \sigma^2\mathbf{I})$, ktoré predpokladáme. Keďže $\hat{Y} = \mathbf{H}Y$, platí $\hat{Y} \sim N(\mathbf{H}\mathbf{X}\beta, \sigma^2\mathbf{H}\mathbf{H})$. Z faktu, že \mathbf{H} je projekčná matica na $M(\mathbf{X})$ dostávame $\mathbf{H}\mathbf{H} = \mathbf{H}$ a $\mathbf{H}\mathbf{X} = \mathbf{X}$.

Pri tvrdení 2 stačí využiť faktu, že matica \mathbf{M} sa premieta na $M(\mathbf{X})^\perp$, a teda platí $\mathbf{M}\mathbf{X} = \mathbf{0}$.

Pre tvrdenie 3 a tvrdenie 4 využijeme vetu 2.7, kedy tvrdenie 3 z nej priamo plynie a tvrdenie 4 je jej bezprostredný dôsledok.

tvrdenie 5 je vlastne upravou vety 2.6, kedy miesto \mathbf{W}^- použijeme $\sigma^2\mathbf{I}$.

□

Veta 2.13. V normálnom lineárnom modeli $Y \sim N(\mathbf{X}\beta, \sigma^2\mathbf{I})$ s úplnou hodnotou platí pre odhad vektora regresných koeficientov:

$$(i) \hat{\beta} \sim N(\beta, \sigma^2\mathbf{C})$$

$$(ii) t_i = \frac{\hat{\beta}_i - \beta_i}{s\sqrt{c_{ii}}} \sim t_{n-p},$$

pre $i = 1, \dots, k$ a $\mathbf{C} = (\mathbf{X}^T\mathbf{X})^{-}$.

Veta 2.14. Nech normálny lineárny model $Y \sim N(\mathbf{X}\beta, \sigma^2\mathbf{I})$ má ľubovoľné riešenie normálnej rovnice $\hat{\beta}$. Potom štatistika:

$$F = \frac{\|\hat{Y} - \mathbf{X}\hat{\beta}\|^2}{\text{RSS}} \frac{n-p}{p} \quad (2.20)$$

má Fischer-Snedecorovo rozdelenie $F_{p, n-p}$.

Na základe poslednej vety možno v normálnom modeli s úplnou hodnotou testovať hypotézy o hodnote celého vektoru β .

Uvažujme náhodný výber (X_1, \dots, X_n) , tj. postupnosť nezávislých náhodných veličín, resp. vektorov, z rozsahu n s rovnakým rozdelením charakterizovaným distribučnou funkciou.

Definícia 6 (Intervalový odhad). :

Interval $(D(X_1, \dots, X_n), H(X_1, \dots, X_n))$ nazveme intervalovým odhadom θ o spoľahlivosti $1 - \alpha$ ak:

$$P_\theta[D(X_1, \dots, X_n) < \theta < H(X_1, \dots, X_n)] = 1 - \alpha,$$

tj. intervalový odhad pokrýva skutočnú hodnotu θ s pravdepodobnosťou $1 - \alpha$.

Krajné body intervalu spoľahlivosti sú funkciami náhodného výberu (X_1, \dots, X_n) , sú to teda náhodné veličiny.

Nech β^0 je vopred daný vektor. Hypotézu $\beta = \beta^0$ zamietame na hladine α v prospech alternatív $\beta \neq \beta^0$, práve keď pre F z vety 2.14 platí:

$$F > F_{k, n-p}(1 - \alpha),$$

kde $F_{k, n-p}(1 - \alpha)$ je $(1 - \alpha)$ -kvantil rozdelenia $F_{k, n-p}$. Ak hypotézu $\beta = \beta^0$ nezamietame, potom množina všetkých takýchto vektorov β^0 vytvorí konfidenčnú množinu pre β s koeficientom spoľahlivosti $1 - \alpha$ v tvare:

$$K_2 = \beta \in R_k : (\beta - \hat{\beta})^T \mathbf{X}^T \mathbf{X} (\beta - \hat{\beta}) \leq ks^2 F_{k, n-p}(1 - \alpha). \quad (2.21)$$

Veta 2.15. *Uvažujme overiteľnú lineárnu hypotézu:*

$$E Y - \mu^0 \in \omega \quad (2.22)$$

v modeli $Y \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I})$ a nech platí $\mathbf{d} = p - \dim \omega > 0$. Potom za platnosti (2.22) má náhodná veličina

$$F = \frac{\mathbf{RSS}_h - \mathbf{RSS}}{\mathbf{RSS}} \frac{n - p}{\mathbf{d}}$$

rozdelenie $F_{\mathbf{d}, n-p}$, náhodná veličina

$$B = \frac{\mathbf{RSS}_h - \mathbf{RSS}}{\mathbf{RSS}_h}$$

má rozdelenie $B_{\mathbf{d}/2, (n-p)/2}$ a náhodná veličina

$$B^* = \frac{\mathbf{RSS}}{\mathbf{RSS}_h}$$

má rozdelenie $B_{(n-p)/2, \mathbf{d}/2}$.

Kapitola 3

Errors in variables

V jednoduchosti možno povedať, že EIV modely sú regresné modely u ktorých sú regresory pozorované s chybami. Navzdory tejto pomerne zjednodušenej definícii, je možné za EIV modelmi vidieť oveľa širšiu škálu. Táto problematika je teda prirodzene podrobne skúmaná viacerými autormi. Avšak medzi najznámejších autorov patria napríklad Wayne A. Fuller [Ful87], Charles E. McCulloch [CE01] alebo Annette J. Dobson [Dob02].

Modely delíme priamo na lineárne EIV modely (*linear EIV models*), nelineárne EIV modely (*nonlinear EIV models*) a čiastočne lineárne modely (*partially linear EIV models*) respektívne v niektorej literatúre ich nájdeme pod anglickým názvom (*partly linear EIV models*). Jednotlivé modely následne možno rozčleniť na základe informácií o správaní sa jednotlivých vstupov. Čím dostávame relatívne zložitú štruktúru tohoto celku regresnej analýzy.

3.1 Lineárne EIV modely

Pri skúmaní jednotlivých častí, začneme najprv modelom, ktorý intuitívne naväzuje na doposiaľ uvedené informácie. Prirodzene sa teda dostávame k lineárnemu EIV modelu, ktorý možno chápať ako priamu nadstavbu predošlých kapitol.

Lineárny EIV model pre jednu závislú premennú je definovaný pre i -te pozorovanie ako:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i,$$

$$w_i = x_i + u_i.$$

Kedže zápis skúmaných dát si možno predstaviť v tvare, ktorý udáva nasledujúca tabuľka:

číslo pozorovania	Y	X	ε	U	W
1	y_1	x_1	ε_1	u_1	w_1
.
.
.
.
n	y_n	x_n	ε_n	u_n	w_n

Tabuľka 2

Prirodzene teda dostávame:

$$Y = \beta^T X + \varepsilon,$$

$$W = X + U,$$

kde Y je vektor závislej premennej, vektor U je náhodný člen respektívne chyba merania a X je vektor regresorov (kovariantov). V tomto modeli sú regresori x_i pozorované s chybou, teda je priamo pozorovaná iba premenná w_i nazývaná tiež *prokázateľná premenná (manifest variable)*. Regresory x_i sa označujú aj ako *latentná premenná (latent variable)*.

Obecne možno lineárne EIV modely rozdeliť aj podľa toho ako sa prejavujú regresori x_i . V prípade, že x_i sú pevné, tj. konštanty, hovoríme o funkcionálnych modeloch (*functional models*). Ak uvažujeme x_i náhodné,

dostávame takzvané štrukturálne modely (*structural models*).

Maticovo možno model prepísať:

$$Y = \mathbf{X}\beta + \varepsilon,$$

$$\mathbf{W} = \mathbf{X} + \mathbf{U},$$

kde \mathbf{X} , \mathbf{U} , \mathbf{W} sú matice v tvare:

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} \quad (3.1)$$

$$\mathbf{U} = \begin{pmatrix} 0 & u_{11} & \dots & u_{1p} \\ 0 & u_{21} & \dots & u_{2p} \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ 0 & u_{n1} & \dots & u_{np} \end{pmatrix} \quad (3.2)$$

$$\mathbf{W} = \begin{pmatrix} 1 & w_{11} & \dots & w_{1p} \\ 1 & w_{21} & \dots & w_{2p} \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ 1 & w_{n1} & \dots & w_{np} \end{pmatrix} \quad (3.3)$$

3.2 Nelineárne EIV modely

Nelineárny EIV model je definovaný ako:

$$E(Y|X_t) = g(X_t),$$

$$W_t = X_t + U_t.$$

V tomto modeli predpokladáme, že vzťah medzi Y a X_t je nelineárny a regresory X_t sú namerané s chybami. Závislú premennú označujeme Y , X_t je premenná pozorovaná s chybou, Z_t je premenná nameraná bez chyby a W_t je *prokazateľná premenná (manifest variable)*.

3.3 Čiastočne lineárne EIV modely

Čiastočne lineárny EIV model je definovaný ako:

$$Y = X_i^T \beta + g(T_i) + \varepsilon_i, \tag{3.4}$$

kde X_i predstavujú vektor kovariátov, T_i sú skalárne kovariáty, funkcia $g(\cdot)$ je neznáma a chyby modelu ε_i sú nezávislé s podmienenou hustotou 0 danou kovariátmi. Počet pozorovaní pre daný model nech je n . Čiastočne lineárny model je predstavený napríklad v článku od autorov Engle, Granger, Rice a Weiss [RE86]. Kde sa jeho tvorcovia zaoberajú súvislosťou aký má vplyv počasie na poptávku po elektrine. Ďalšie zaujímavé štúdie možno nájsť v publikáciách od Heckmanna [E.N86], Chena [Che88], Speckmanna [Spe88], Lianga a Härdla [HL97] a Severiniho a Staniswalissa [A.T94].

Pri skúmaní modelu nás primárne zaujíma odhad neznámeho parametru β a neznámej funkcie $g(\cdot)$ v modely (3.4), kde kovariáty X_i sú merané s chybami merania. Namiesto priameho pozorovania hodnôt X_i teda pozorujeme hodnoty:

$$W_i = X_i + U_i. \quad (3.5)$$

Pozorované hodnoty W_i pozostávajú teda z dvoch častí, a to z chýb merania, ktoré označíme ako U_i a z kovariátov X_i . Pre uvedené hodnoty U_i platí, že sú nezávislé a stejne rozdelené. Teda nezávislé na (Y_i, X_i, T_i) , s priemerom 0 a kovariančnou maticou Σ_{uu} . Predpokladajme, že kovariančná matica Σ_{uu} je pre nás známa. Ak sú hodnoty X_i pozorovateľné, potom by sme štandardne dostali hodnoty odhadu parametru β . Pre každé pevné β , nech je $\hat{g}(T, \beta)$ odhadom pre $g(T)$. Potom napríklad v implamentácii použitej v práci Severiniho a Staniswalissa [A.T94], $\hat{g}(T, \beta)$ maximalizuje vážený pravdepodobnostný odhad a chyby modelu ε_i majú normálne rozdelenie s váhami danými pomocou jádrových váh zo symetrického jádrového odhadu hustoty $K(\cdot)$ a šírkou pásma h . Keďže máme $\hat{g}(T, \beta)$ je možné odhadnúť β pomocou metódy najmenších štvorcov ako:

$$\min \sum_{i=1}^n [Y_i - X_i^T \beta + g(T_i, \beta)]^2. \quad (3.6)$$

V tomto konkrétnom prípade, odhad pre β je teda určený explicitne. Nech $\hat{g}_{y,h}(\cdot)$ a $\hat{g}_{x,h}(\cdot)$ sú jadra so šírkou pásma h pre Y, X a T . Potom :

$$\hat{\beta}_n = \left[\sum_{i=1}^n (X_i - \hat{g}_{x,h}(T_i))(X_i - \hat{g}_{x,h}(T_i))^T \right]^{-1} \times \sum_{i=1}^n (X_i - \hat{g}_{x,h}(T_i))(Y_i - \hat{g}_{y,h}(T_i)) \quad (3.7)$$

Jednou z dôležitých vlastností odhadu (3.7) je, že nevyžaduje vyhladzovanie, takže šírka pásma v rozsahu $h \sim n^{-1/5}$ vedie k výsledku:

$$n^{1/2}(\hat{\beta}_n - \beta) \rightarrow Normal(0, B^{-1}CB^{-1}). \quad (3.8)$$

Kde B je kovariančná matica pre $X - E(X|T)$ a C je kovariančná matica pre

$\varepsilon[X - E(X|T)]$. Forma zápisu (3.7) môže byť použitá ak neberieme v úvahu chybu merania a nahradíme X pomocou W , potom výsledný odhad je nekonzistentný s β . Tvar, v ako máme odhad vyjadrený, ale ukazuje ešte niečo ďalšie. Je obecné známe, že v lineárnej regresii, nekonzistentnosť spôsobenú chybou merania je možné obísť použitým "korekcie útlumu" (*correction for attenuation*). V kontexte semiparametrických modelov, táto úvaha nás privádza k tomu, aby sme použili odhad v tvare:

$$\hat{\beta}_n = \left[\sum_{i=1}^n (W_i - \hat{g}_{w,h}(T_i))(W_i - \hat{g}_{w,h}(T_i))^T - n\Sigma_{ww} \right]^{-1} \times \sum_{i=1}^n (W_i - \hat{g}_{w,h}(T_i))(Y_i - \hat{g}_{y,h}(T_i)). \quad (3.9)$$

Odhad (3.9) je teda možné odvodiť rovnakým postupom ako je odvodený odhad od Severiniho a Staniswalisa.

Kapitola 4

Aplikácia EIV modelov v praxi

4.1 Výskum rakoviny prsu

Využitie EIV modelov praxi je možné názorne predviesť na nasledujúcej štúdií, pochádzajúcej od Gail Gong, Alice S. Whittemore a Stelly Grosser [GG90], ktorú s aplikáciou na errors in variables modely bližšie rozobrala Alice S. Whittemore [Whi89]. Štúdia sa zaoberala meraním úmrtnosti spojenej s výskytom rakoviny prsu u žien žijúcich v oblasti San Francisca. Zrojom dát bol register nádorových ochorení v oblasti San Francisca. Skúmala sa úmrtnosť v skupine $N = 2495$ žien, ktorým bola diagnostikovaná rakoviná prs. Vek žien sa pohyboval v rozmedzí 55 - 64 rokov.

Pri štúdií sa autori zamerali hlavne na to, ako miera úmrtnosti závisí na závažnosti ochorenia pri diagnóze a na čase, ktorý uplynul od určenia diagnózy. Každá pacientka bola zaradená do jednej z 5-tich skupín, ktoré reflektovali jednotlivé štádiá ochorenia. Pričom väčšie číslo skupiny predstavovalo rozvinutejšiu formu ochorenia, viz tabuľka 3.

Pri určovaní štádia ochorenia sa vychádzalo z informácií uvedených v zdravotných záznamoch jednotlivých pacientiek. Keďže záznamy boli často nekompletné, je prevdepodobné, že niektoré pacientky boli zaradené do nesprávnej skupiny. Pacientky boli sledované po dobu 10 rokov, pričom toto obdobie bolo rozdelené na päť dvojročných cyklov. V rámci jedného cyklu malo každé štádium konštantnú mieru úmrtnosti, tj. ku zmene miery

úmrtnosti daného štádia mohlo dôjsť iba na počiatku nového dvojročného cyklu.

4.1.1 Rozbor podkladových dát

Nech $\lambda = \lambda_{jk}$ je matica typu (J, K) , ktorá nám udáva pravdepodobnosti úmrtnosti. Hodnoty $J = 1, \dots, 5$ predstavujú jednotlivé cykly, a hodnoty $K = 1, \dots, 5$ predstavujú jednotlivé štádia ochorenia, do ktorých sme pacientky rozdelili. Rozloženie smrtí (D) a počtu uplynulých mesiacov od určenie diagnózy (DM) pre sledovaný počet 2 495 pacientiek zachycuje tabuľka 3.

Mesiace	1. N = 637		2. N = 1045		3. N = 441		4. N = 167		5. N = 205	
	D	PM	D	PM	D	PM	D	PM	D	PM
< 24	12	12 773	42	20 456	45	8 359	38	3 072	123	2 788
24 - 48	24	9 048	70	13 096	66	4 485	28	1 508	30	962
48 - 72	21	5 469	33	6 767	14	1 539	8	664	10	267
72 - 96	6	2 639	18	2 651	4	451	1	274	2	131
96 - 120	0	596	2	477	0	82	0	52	1	51

Tabuľka 3: Štádium choroby

Data jednotlivých pacientiek pozostávajú z nasledujúcich údajov: t predstavuje čas do smrti alebo do ukončenia štúdie, premenná δ predstavuje indikátor úmrtia, nadobúda hodnoty 1 pri úmrtí a hodnoty 0 v prípade života pacientky. Posledným údajom je kovariančný vektor indikátorov štádia ochorenia, v ktorom sa pacientka nachádza. Označili sme ho ako $x = (x_1, \dots, x_K)$. Pre pacientku v štádiu k je $x = e_k$, kde e_k je K -dimenzionálny stĺpcový vektor s k -tou komponentou rovnou 1 a ostatnými nulovými komponentami.

Predpokladáme, že štádium, čas do smrti a čas do kontroly sú nezávislé náhodné veličiny. Potom hustota pre (t, δ) daná $x = e_k$ je úmerná k:

$$p(t, \delta, e_k, \lambda) = \exp\left\{\sum_{j=1}^J [\delta m_j(t) \log \lambda_{jk} - M_j(t) \lambda_{jk}]\right\}. \quad (4.1)$$

kde, $m_j(t)$ je indikátor rovný 1 ak t patrí do daného dvojročného cyklu, a rovný 0 inak. $M_j(t)$ je doba trávená v periode j , a konštanta úmrtnosti je nezávislá na λ . Pre hypotetického pacienta, ktorý zomrel v čase $t = 33$ mesiacov je potom $m_2(t) = 1$ a $m_j(t) = 0, j \neq 2$. Okrem toho je $M_1(t) = 24$, $M_2(t) = 9$ a $M_j(t) = 0$ pre $j > 2$. Nasčítaním (4.1) cez N sledovaných subjektov možno napísať log-likelihood kernel dát (t_i, δ_i) pre $i = 1, \dots, N$ ako:

$$S_1 = \sum_{jk} [\Delta_{jk} \log \lambda_{jk} - \mu_{jk} \lambda_{jk}], \quad (4.2)$$

kde Δ_{jk} a μ_{jk} sú počty úmrtí (D) a uplynulé mesiace od určenie diagnózy (MD) v období j pre štádium k pacientov. Teda (4.2) predstavuje log-likelihood kernel pre $J \times K = 25$ poissonových premenných Δ_{jk} , ktoré majú priemere $\mu_{jk} \lambda_{jk}$. Teda regresné modely pre neznáme úmrtnostné miery λ_{jk} možno analyzovať prostredníctvom Poissonových regresných metód.

Najobecnejší model *Model1*, umožňuje aby sa λ_{jk} voľne menili. Graf maximálne vierohodných odhadov hazard rates pre tento model ukazuje tendencie zvyšovania spolu so štádiom, avšak bez jasne viditeľných znakov paralelizmu v rámci daného skúmaného štádia.

The constant hazard model je označený ako *Model2*. Tento model predpokladá, že λ_{jk} sú konštantné v rámci všetkých 5-tich periód pre jednotlivé štádia. Tento model predpokladá nezávislosť hazard rate na čase od určenia diagnózy, čo je značne reštriktívne. Vskutku, likelihood ratio test porovnania *Model1* vs *Model2* vychádza silne rozporne, kedy chi-square = 97.1 s 20-timi stupňami voľnosti a $p < 0.01$.

Nie tak príliš reštriktívny je potom proportional hazard model *Model3*, ktorý predpokladá, že log hazard rates pre prvé štyri štádia sú zhodné. Do výpočtu sa nezahŕňa piate štádium, pretože prežitie pacientov v piatom štádiu, ktorí sú diagnostikovaní s metastázami sa správa inak ako u pacientov zaradených v ostatných štádiách onkologického ochrenia. Likelihood ratio test *Model1* vs *Model3* potom odmieta *Model3*, keďže chi-square = 23.2 s 12-timi stupňami voľnosti a $p < 0.025$.

V tomto bode by sme mohli usudzovať, že obecný hazard model, ktorý sebou síce nesie riziko komplikovaných interpretácií, je vhodným modelom. Avšak autori skúmali aj možnosť, že aj jeden z viac reštriktívnych modelov je vhodný. A jeho nevhodnosť navonok, je na vrub chybám v priradovaní pacientov do jednotlivých štádií ochorenia.

4.1.2 Analýza Errors-in-variables

Nech $\gamma = (\gamma_{kl})$ je matica pravdepodobností typu $(5, 5)$. Potom γ_{kl} je pravdepodobnosť, že pacientka pozorovaná v štádiu l je v skutočnosti v štádiu k , kde $k, l = 1, \dots, 5$. Táto informácia nám obmedzila maticu γ na takú, že v nej nepozorujeme prechodné obdobia, medzi jednotlivými štádiami ochorenia. Má teda tvar:

$$\gamma = \begin{pmatrix} \gamma_{11} & 0 & 0 & 0 & 0 \\ 1 - \gamma_{11} & \gamma_{22} & 0 & 0 & 0 \\ 0 & 1 - \gamma_{22} & \gamma_{33} & 0 & 0 \\ 0 & 0 & 1 - \gamma_{33} & \gamma_{44} & 0 \\ 0 & 0 & 0 & 1 - \gamma_{44} & 1 \end{pmatrix} \quad (4.3)$$

Teda pre pacientku pozorovanú v štádiu 1 predpokladáme, že je v štádiu 1 alebo v štádiu 2 s pravdepodobnosťami γ_{11} a $1 - \gamma_{11}$. To isté potom platí aj pre ostatné štádia ochorenia 2, 3 a 4. Štádium 5, v ktorom pacientky mali metastatický nález, bolo určené jednoznačne.

Parametre (λ, γ) boli odhadnuté maximalizáciou vierohodnosti pozorovaných dát, pomocou pozorovaných štádiových premenných X_1, \dots, X_N :

$$\prod_{i=1}^N \sum_{k=1}^K f_{t, \delta, x|X}(t_i, \delta_i, e_k | X_i; \lambda, \gamma). \quad (4.4)$$

Kde $f_{t,\delta,x|X}$ je podmienená pravdepodobnostná hustota, ktorá je funkciou (t, δ, x) a je daná hodnotou vektora štádií X . Za určitých podmienok je potom úmerná k:

$$\prod_{i=1}^N \sum_{k=1}^K p(t_i, \delta_i, e_k, \lambda) f_{x|X}(e_k | X_i; \gamma). \quad (4.5)$$

kde $p(\cdot)$ je dané (4.1), $f_{x|X}$ je mnohočlenná hustota pre vektor x skutočného štádia, podmienením pozorovaným štádiom a parametrizovaný pomocou γ z (4.3). Konstanta úmernosti pritom nezávisí na (λ, γ) . Úmernosť (4.4) a (4.5) je splnená, za predpokladu, že platia nasledujúce tri predpoklady: (i.) časy do úmrtia a kontroly sú stochasticky nezávislé, (ii.) rozdelenie času do smrti nezávisí na X a (iii.) rozdelenie času do kontroly nezávisí na x . V takom prípade potom možno predpokladať, že (4.4) a (4.5) dávajú oba správne pozorované stavy indikátorov (x, X) .

Prvý predpoklad zabezpečí, že časy do úmrtia majú rozdelenie u skúmaných pacientov v danom pravdivom štádiu, rovnajúcim sa neskúmaným pacientom. Tento postup je bežne používaný v analýze prežitia. Druhý predpoklad ukotvuje nezávislosť prognózy prežitia na presnosti zaradenia množstva skúmaných pacientov do jednotlivých skupín podľa ich skutočného stavu choroby. Tento predpoklad vychádza z analýzy chyby merania. Posledný, tretí predpoklad nakoniec zabezpečí, že skúmaná prognóza je nezávislá na skutočnom stave množstva pacientov pozorovaných v danom štádiu. Predpoklad teda platí približne v prípade, že pravdepodobnosti zaradenia do nesprávneho štádia choroby sú nízke, alebo keď pravdepodobnosť výberu má malý vzťah k rozsahu choroby a stavu štádia.

Odhady $(\hat{\lambda}, \hat{\gamma})$ maximalizace (4.5) nie sú dostupné v priamej forme a možno ich získať použitím EM algoritmu popísaného napríklad v práci od A.P. Dempster, N.M.Laird a D.B.Rubin [A.P77]. *E – krok* (*E – step*) z EM algoritmu napočítava $Q(\lambda^c, \gamma^c)$, predpoklad úplnosti dát logaritmickej vierohodnosti, závisí na pozorovaných datach a aktuálnom odhade (λ^c, γ^c) . Kompletnosť dát v tomto prípade predstavuje hodnoty $(t_i, \delta_i, X_i, x_i)$ pre

$i = 1, \dots, N$. M – krok (M – step) následne maximalizuje Q , tak aby sme dosiahli nový odhad. Z E – kroku dostávame:

$$Q(\lambda, \gamma) = \Xi_1(\lambda) + \Xi_2(\gamma), \quad (4.6)$$

kde:

$$\Xi_1(\lambda) = \sum_{jk} [\varepsilon(\Delta_{jk}) \log \lambda_{jk} - \varepsilon(\mu_{jk}) \lambda_{jk}], \quad (4.7)$$

a ďalej:

$$\Xi_2(\gamma) = \sum_{kl} [\varepsilon(\nu_{kl}) \log \gamma_{kl}]. \quad (4.8)$$

Tu ν_{kl} predstavuje počet k – teho štádia pacientov s pozorovaným stavom diagnózy l a $\varepsilon(\cdot)$ udáva podmienené očakávané hodnoty dané (t_i, δ_i, X_i) pre $i = 1, \dots, N$ parametrizáciou (λ, γ) . Tieto očakávania sú potom splnené jednoduchým nahradením nepozorovaného vektoru x pre každého pacienta, vektorom $\chi = (\chi_1, \dots, \chi_K)$ pravdepodobností, ktoré vyjadrujú skutočný vzťah k jednotlivým stavom choroby. Jednotlivé komponenty χ reprezentujú potom podmienené pravdepodobnosti zavedenia v jednotlivých stavoch, daných pozorovaným stavom pacienta, datami prežitia a odhadmi parametrov [GG90]. Vo výsledku teda možno zhrnúť, že prevedené simulácie v [GG90] odhaľujú nasledujúce dve zistenie: (i) v prípade, že ignorujeme možnosť nesprávneho zaradenia do skupiny, môžeme tým spôsobiť odmietnutie proportional hazard modelu, i keď je tento model správny. Na druhú stranu zistenie (ii) v prípade, že nijakým spôsobom nevalidujeme data, likelihood ratio test je skoro nepoužiteľný pri detekcii nesprávneho zaradenia do skupiny, jedine, že by sme skúmali skutočne veľké množstvo dát.

⁴⁵Záver

V úvodnej kapitole tejto práce sme zdefinovali základné pojmy a pojem regresie a na príklade merania synov a otcov sme si ukázali jednu z možností praktického využitia regresie. V druhej kapitole sme rozobrali lineárny regresný model a jeho vlastnosti, pričom sme sa zaoberali i normálnym lineárnym modelom, jeho charakteristikami a možnosťami odhadu jeho parametrov. Tieto poznatky sme následne v tretej kapitole využili pri zavedení jednotlivých typov modelov errors in variables.

V poslednej kapitole tejto práce sme uviedli aplikáciu modelov error in variables v medicíne, kde sme sa zaoberali výskumom rakoviny prs. Na základe uvedených poznatkov sme dospeli k rozhodnutiu, že použitie analýzy errors in variables nás doviedlo k niekoľkým predpokladom a vlastnostiam, ktoré sú požadované po skúmaných dátach. Obdobne je teda možné nájsť množstvo zaujímavých aplikácií v rôznych oblastiach vedy.

Literatúra

- [And05] J. Anděl. Základy matematické statistiky. MATFYZPRESS, 2005.
- [A.P77] D.B.Rubin A.P.Dempster, N.M.Laird. Maximum likelihood from incomplete data via the em algorithm (with discussion). Journal of the royal statistical society. Series B, Vol.39, pp. 1-39, 1977.
- [A.T94] G.J. Staniswalis A.T.Severini. Quasilikelihood estimation in semi-parametric models. Journal of the American statistical association, Vol. 89, pp. 501-511, 1994.
- [CE01] Shayle R. Searle Charles E.McCulloch. Generalized,linear and mixed models. John Wiley Sons, Inc, 2001.
- [Che88] H. Chen. Convergence rates for parametric components in a partly linear model. The annals of statistics, Vol.16, pp. 136-146, 1988.
- [Dob02] Annette J. Dobson. An introduction to generalized linear models. Chapman Hall/CRC, 2002.
- [Dod93] David Birkes;Yadolah Dodge. Alternative Methods of regression. A Wiley-Interscience Publication, 1993.
- [E.N86] E.N.Heckman. Spline smoothing in partly linear models. Journal of the royal statistical society. Series B, Vol. 48, pp. 244-248, 1986.
- [Ful87] Wayne A. Fuller. Measurement error models. John Wiley Sons, Inc, 1987.
- [Gal86] Francis Galton. Regression towards mediocrity in hereditary stature. Journal of the Anthropological Institute, 1886.

- [GG90] Stella Grosser Gail Gong, Alice S. Whittemore. Censored survival data with misclassified covariates: A case study of breast-cancer mortality. *Journal of the American statistical association*, Vol. 85, No. 409, pp.20-28, 1990.
- [HL97] W. Härdle H. Liang. Asymptotic normality of parametric part in partially linear heteroscedastic regression models. Humboldt university Berlin, DP 33, SFB 373, 1997.
- [RE86] J. Rice A. Weiss R.F. Engle, C.W.J. Granger. Semiparametric estimates of the relation between weather and electricity sales. *Journal of the American statistical association*, Vol. 81, pp.310-320, 1986.
- [Spe88] P. Speckman. Kernel smoothing in partial linear models. *Journal of the royal statistical society. Series B*, Vol.50, pp. 413-436, 1988.
- [Whi89] Alice S. Whittemore. Errors-in-variables regression problems in epidemiology. *American mathematical society*, Vol.112, pp.17-32, 1989.
- [Zvá89] Karel Zvára. Regresní analýza. Academia, 1989.