

UNIVERZITA KOMENSKÉHO V BRATISLAVE  
FAKULTA MATEMATIKY, FYZIKY A INFORMATIKY

Riadkovo-stĺpcové návrhy štatistických experimentov

BAKALÁRSKA PRÁCA

UNIVERZITA KOMENSKÉHO V BRATISLAVE  
FAKULTA MATEMATIKY, FYZIKY A INFORMATIKY

**Riadkovo-stĺpcové návrhy štatistických experimentov**

**BAKALÁRSKA PRÁCA**

Študijný program: Matematika  
Študijný odbor: 1114 Matematika  
Školiace pracovisko: Katedra aplikovanej matematiky a štatistiky  
Vedúci práce: doc. Mgr. Radoslav Harman, PhD.



Univerzita Komenského v Bratislave  
Fakulta matematiky, fyziky a informatiky

---

## ZADANIE ZÁVEREČNEJ PRÁCE

**Meno a priezvisko študenta:** Róbert Druska  
**Študijný program:** matematika (Jednoodborové štúdium, bakalársky I. st., denná forma)  
**Študijný odbor:** matematika  
**Typ záverečnej práce:** bakalárska  
**Jazyk záverečnej práce:** slovenský  
**Sekundárny jazyk:** anglický

**Názov:** Riadkovo-stĺpcové návrhy štatistických experimentov  
*Row-column designs of statistical experiments*

**Anotácia:** Prvým cieľom je analyzovať vlastnosti odhadov parametrov regresného modelu pre takzvaný riadkovo-stĺpcový experimentálny návrh. Druhým cieľom je navrhnúť algoritmus na výpočet optimálneho návrhov tohto typu, v závislosti na požiadavkách experimentátora.

**Vedúci:** doc. Mgr. Radoslav Harman, PhD.  
**Katedra:** FMFI.KAMŠ - Katedra aplikovanej matematiky a štatistiky  
**Vedúci katedry:** prof. RNDr. Marek Fila, DrSc.  
**Dátum zadania:** 15.10.2019

**Dátum schválenia:** 18.10.2019

prof. RNDr. Ján Filo, CSc.  
garant študijného programu

.....  
študent

.....  
vedúci práce

Pod'akovanie TODO

# Abstrakt

TODO

**Klíčové slova:** lineární regresný model

# Abstract

TODO

**Keywords:** linear regression

# Obsah

Úvod	8
<b>1 Lineárny regresný model</b>	<b>9</b>
1.1 Metóda najmenších štvorcov . . . . .	9
1.2 Maticová príprava . . . . .	11
1.3 Metodika skúmania . . . . .	12
<b>2 Riadkovo-stĺpcový aditívny model s dvoma typmi ošetrov</b>	<b>13</b>
2.1 Odhadnuteľnosť $h^Tb$ . . . . .	15
2.2 Optimálnosť návrhu modelu . . . . .	18
<b>Záver</b>	<b>20</b>
<b>Zoznam použitej literatúry</b>	<b>21</b>
<b>Príloha A</b>	<b>22</b>

# Úvod

TODO



# 1 Lineárny regresný model

Majme  $n$  nameraných štatistických jednotiek tvaru  $\{y, x_1, \dots, x_p\}$ , ktoré sme dostali ako výsledok experimentu. Lineárny regresný model predpokladá, že medzi jednotlivými prvkami  $y, x_1, \dots, x_p$  je lineárny vzťah. Motiváciou za lineárnym regresným modelom je spravidla aproximovať tento lineárny vzťah.

Aproximácia lineárneho vzťahu nám v praxi ponúkne mechanizmus, ktorým možno predikovať neznámu hodnotu  $y$  na základe známych hodnôt  $y, x_1, \dots, x_p$ , čo v reálnom živote predstavuje často sa vyskytujúci problém.

Označme teda daný lineárny vzťah medzi zložkami nameranej štatistickej jednotky:

$$y_i = b_0 + b_1x_{i1} + \dots + b_px_{ip} + e_i = b^T x_i + e_i$$

kde  $\{y_i, x_i\}$  je  $i$ -ta nameraná jednotka,  $b$  je vektor lineárneho vzťahu a  $e_i$  je chyba merania.

Keď lineárne vzťahy pre každú z  $n$  nameraných jednotiek zapíšeme maticovo, dostaneme vzťah

$$y = Xb + e \tag{1}$$

kde  $y = (y_1, y_2, \dots, y_n)^T$ ,  $e = (e_1, e_2, \dots, e_n)^T$  a

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}$$

je matica tvaru  $n \times p$ . V praxi je  $b$  neznámy vektor, ktorý sa snažíme odhadnúť.

Na spočítanie odhadu  $b$  sa používajú rôzne metódy, najčastejšie napr. metóda najmenších štvorcov alebo metóda maximálnej vierohodnosti.

## 1.1 Metóda najmenších štvorcov

Metódou najmenších štvorcov vypočítame odhad  $\hat{b}$  parametra  $b$  nasledovne:

$$\hat{b} = \underset{b \in \mathbb{R}^p}{\operatorname{argmin}} (y - Xb)^T C (y - Xb) = \underset{b \in \mathbb{R}^p}{\operatorname{argmin}} \|y - Xb\|_{C^{-1}}^2$$

kde  $C$  je nejaká kladne definitná matica. Ak  $C = I$ , potom minimalizujeme výraz  $\|y - Xb\|_I^2 = \|y - Xb\|^2 = \sum_{i=1}^n (y_i - X_i \cdot b)^2$ , kde  $X_i$  značí  $i$ -ty riadok matice  $X$ . V našej práci budeme predpokladať nezávislosť chýb a tiež homogenitu ich rozptylu, čo v praxi znamená, že skutočne budeme môcť dosadiť  $C = I$ . Preto maticu  $C$  v ďalšom opise teórie spomínať nebudeme.

Geometricky metódu najmenších štvorcov možno interpretovať ako projekciu vektora  $y$  na stĺpcový priestor matice  $X$ . Hľadáme teda taký vektor  $\hat{b}$ , pre ktorý platí  $X\hat{b} = Py$ , kde  $P$  je matica ortogonálnej projekcie na stĺpcový priestor  $X$ . Z teórie lineárnej algebry vieme, že  $P = X(X^T X)^- X^T$ , kde znamienko  $-$  označuje  $g$ -inverziu. ( $g$ -inverziou matice  $A$  je taká matica  $A^-$ , pre ktorú platí  $AA^-A = A$ ).

Odhad  $\hat{b}$  parametra  $b$  je teda riešením rovnice

$$Xb = X(X^T X)^- X^T y \quad (2)$$

Všetky riešenia tejto rovnice musia mať tvar:

$$(X^T X)^- X^T y$$

kde použitá  $g$ -inverzia je ľubovoľná.

Z uvedeného vyplýva, že v prípade regulárnosti  $X^T X$  je odhad  $\hat{b}$  parametra  $b$  jednoznačný. V našej práci budeme skúmať matice (modely)  $X$ , ktoré nie sú regulárne, takže jednoznačný odhad  $\hat{b}$  nebudeme schopní nájsť (čo v konečnom dôsledku ani nie je naším záujmom). Budeme odhadovať lineárnu funkciu zložiek vektora  $b$ , konkrétne  $h^T b = h_1 b_1 + \dots + h_p b_p$ , ktorá býva jednoznačne odhadnuteľná aj v prípade singularity  $X^T X$ , ak vektor  $h$  spĺňa určité podmienky.

**Definícia 1.1.** *Lineárnu kombináciu  $h^T b$  zložiek vektora  $b$  nazývame odhadnuteľnou, ak pre ľubovoľné riešenia  $b^*$  a  $b^{**}$  rovnice (2) platí  $h^T b^* = h^T b^{**}$ .*

Je niekoľko ekvivalentných podmienok, ktoré stačia na to, aby  $h^T b$  bolo odhadnuteľné. Z nich spomenieme dve v nasledovnej vete, ktorú použijeme neskôr v našej práci.

**Veta 1.2.**  *$h^T b$  je odhadnuteľné, ak platia nasledovné ekvivalentné podmienky:*

1.  $h \in \mathcal{M}(X^T)$
2.  $h \in \mathcal{M}(X^T X)$ ,

kde  $\mathcal{M}$  označuje stĺpcový priestor matice.

Ak  $h$  patrí do riadkového priestoru matice  $X$ , potom existuje také  $u$ , že  $h = X^T u$ . Potom pre jednoznačný odhad  $h^T \hat{b}$  vektora  $h^T b$  platí:

$$h^T \hat{b} = u^T X \hat{b} = u^T P y = u^T X (X^T X)^- X^T y$$

Vidíme, že v strede výrazu na pravej strane je projekčná matica  $X(X^T X)^- X^T$ , ktorá je vždy jednoznačne určená, nezávisle na voľbe zovšeobecnenej  $g$ -inverzie v danom výraze.

Výsledok predchádzajúcej vety je dôležitý pre našu prácu, pretože nebudeme skúmať odhady  $b$ , ale odhady niektorých lineárnych kombinácií zložiek vektora  $b$ , konkrétne napr. rozdiely medzi parametrami.

Odhad  $\hat{b}$  parametra  $b$ , ako aj odhad  $h^T \hat{b}$  parametra  $h^T b$ , sú lineárne nevychýlené odhady, ktorým prislúcha disperzia (TODO: popremýšľaj, či treba bližšie opísať lineárny nevychýlený odhad). Neskôr v našej práci budeme hľadať také modely  $X$ , pri ktorých je disperzia odhadov  $h^T b$  najmenšia možná, čo nám dá najlepší lineárny nevychýlený odhad. Ak nami navrhované modely  $X$  budú opisovať ten istý experiment, ten model  $X$ , pre ktorý disperzia odhadu  $h^T b$  bude najmenšia, bude svojím spôsobom optimálny.

K nájdeniu optimálneho modelu  $X$  nám poslúži Gaussova-Markovova veta, ktorá určuje minimálnu možnú disperziu odhadu  $h^T b$ .

**Veta 1.3.** (Gaussova-Markovova) *Nech  $h$  je z riadkového priestoru  $X$ . Potom minimálna možná disperzia lineárneho nevychýleného odhadu  $h^T b$  je*

$$m = h^T M^- h$$

kde  $M = X^T X$  je informačná matica parametra  $b$  a  $M^-$  je jej ľubovoľná  $g$ -inverzia.

## 1.2 Maticová príprava

V našej práci budeme často pracovať s poznatkami z lineárnej algebry, preto v krátkosti zhrnieme niektoré najdôležitejšie z nich.

**Definícia 1.4.** *Nech  $A$  je ľubovoľná matica. Potom matica  $A^-$  také, že*

$$A A^- A = A$$

sa nazýva *g-inverzia* (alebo *pseudoinverzia*) matice  $A$ .

*Poznámka:* Ak  $A^{-1}$  neexistuje, tak pre maticu  $A$  existuje nekonečný počet  $g$ -inverzií.

Pseudoinverzie použijeme neskôr v našej práci, keď budeme skúmať, či vektor patrí do stĺpcového priestoru matice. Na to nám poslúži nasledovná veta.

**Veta 1.5.** *Nech rovnica  $Ax = y$  má riešenie a nech  $A^-$  je ľubovoľná  $g$ -inverzia. Potom  $A^-y$  je riešením tejto rovnice.*

Dôkaz: Keďže  $Ax = y$  má riešenie, tak existuje také  $x_0$ , že  $Ax_0 = y$ . Potom

$$A(A^-y) = A(A^-Ax_0) = Ax_0 = y$$

### 1.3 Metodika skúmania

Cieľom nášho skúmania bude nájsť ideálnu maticu vstupu do experimentu (konkrétny problém popíšeme v nasledujúcej kapitole). Pokúsime sa vysloviť závery pre všeobecný prípad matice  $m \times n$ , avšak pre istý prvotný náhľad do problematiky problém popíšeme a preskúmame na menších rozmeroch, konkrétne do rozmeru, do ktorého nám to umožní výpočtová technika. Pracovať budeme s programovacím jazykom *python*.

Od „preskúmania situácie“ do istého rozmeru si sľubujeme dôležité náhľady, ktoré nám umožnia vysloviť závery a hypotézy pre všeobecný prípad  $m \times n$ . Tie sa následne pokúsime dokázať matematickými metódami.

## 2 Riadkovo-stĺpcový aditívny model s dvoma typmi ošetroení

Cieľom tejto práce je skúmať experiment reprezentovaný binárnou maticou. Možná reprezentácia nášho modelu pri matici rozmeru  $m \times n$  je takáto:

Uvažujeme model s dvomi kvalitatívnymi faktormi  $A$  a  $B$ , pričom faktor  $A$  má  $m$  úrovní a faktor  $B$   $n$  úrovní. Pre každú kombináciu faktorov experimentátor zvolí jedno z dvojice ošetroení (angl. treatment). Budeme uvažovať len aditívny model bez interakcií.

Napríklad faktor  $A$  môže označovať dobrovoľníka (t.j. máme  $m$  dobrovoľníkov) a faktor  $B$  účinnú látku, ktorú nazveme liečivo (t.j. máme  $n$  liečiv). Ošetroenie 0 reprezentuje podanie liečiva orálnym spôsobom a ošetroenie 1 vnútrožilne). Účinok, ktorý po čase nameriame na dobrovoľníkovi, závisí aditívne od efektu samotného dobrovoľníka (jeho predispozícií), efektu liečiva a efektu ošetroenia. Všetky tieto efekty uvažujeme ako neznáme parametre modelu.

Návrh teda možno reprezentovať binárnou maticou (alebo bipartitným grafom, čo opíšeme neskôršie v našej práci).

Cieľom experimentu je odhadnúť rozdiel medzi dvoma efektmi (treatmentami). Cieľom hľadania optimálneho modelu je nájsť taký návrh, pri ktorom rozptyl Gauss-Markovovho odhadu rozdielu medzi efektmi dvojice ošetroení (t.j. odhadu kontrastu ošetroení) bude čo najmenší. (Použijeme pri tom Gaussovu-Markovovu vetu z predošlej kapitoly.)

Pre maticu typu  $m \times n$  dostaneme teda  $mn$  vzťahov tvaru:

$$y_{ij} = a_i + b_j + t + e_{ij}$$

$i = 1, \dots, m$ ;  $j = 1, \dots, n$ ;  $t$  je  $t_0$  alebo  $t_1$  podľa hodnoty na  $ij$ -tej pozícii matice a  $e_{ij}$  je chyba. Budeme odhadovať rozdiel medzi  $t_0$  a  $t_1$ .

Z návrhu modelu v podobe binárnej matice tvaru  $m \times n$  vytvoríme lineárny regresný model s maticou tvaru  $mn \times (m+n+2)$ , pretože máme  $mn$  meraní závislých od  $m+n+2$  premenných.

Príklad:

Maticu návrhu experimentu

$$\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \quad (3)$$

prepíšeme ako (TODO: popíš tými strapatými zátvorkami čo znamená čo):

$$\begin{bmatrix} 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 & 1 & 0 \end{bmatrix} \quad (4)$$

Označme vo všeobecnosti maticu modelu (3)  $B$  a maticu lineárneho zobrazenia (4)  $X$ . Predpokladáme lineárne vzťahy, preto dostávame lineárnu regresiu

$$y = Xb + e$$

kde  $y$  je vektor dát nameraných v experimente,  $b$  je vektor neznámych koeficientov a  $e = (e_{11}, e_{12}, e_{13}, \dots, e_{m1}, \dots, e_{mn})^T$  je vektor chýb merania. Uvedomme si, že v našej práci nebudeme pracovať s konkrétnym vektorom  $y$ ; snažíme sa totiž navrhnúť maticu modelu  $B$  (a z nej vyplývajúcu maticu  $X$ ) ešte pred tým, než experiment vykonáme.

Naším cieľom teraz bude odhadnúť rozdiel medzi efektmi  $t_0$  a  $t_1$ , resp. zistiť, aká môže byť disperzia tohto odhadu pri danej matici návrhu  $B$ . Modely, pre ktoré bude možná disperzia najmenšia, budeme považovať za optimálne.

V danej lineárnej regresii  $y = Xb + e$  teda nebudeme odhadovať celý vektor  $b$ , ale len jeho lineárnu kombináciu  $h^T b$ , kde  $h$  je vektor rozmeru  $m + n + 2$ , ktorý má na  $m + n$  miestach 0 a na zvyšných dvoch miestach +1 a -1, ktoré zodpovedajú efektom  $t_0$  a  $t_1$ . Z algoritmu, ktorým sme z matice návrhu  $B$  vytvorili maticu lineárnej regresie  $X$ , vyplýva, že vektor  $h$  má tvar  $(0, 0, 0, \dots, +1, -1)^T$ .

Na zistenie, či  $h^T b$  je odhadnuteľné, použijeme vetu (1.2), a na následné nájdenie minimálneho rozptylu Gauss-Markovovho odhadu rozdielu medzi efektmi použijeme Gaussovu-Markovovu vetu (1.3).

## 2.1 Odhadnuteľnosť $h^T b$

Pokúsime sa preskúmať, kedy matica návrhu umožňuje odhadnuteľnosť hodnoty  $h^T b$  pre vyššie spomenuté  $h$ . Situáciu preskúmame do rozmeru  $4 \times 5$  a na základe našich zistení vyslovíme hypotézu pre všeobecný prípad  $m \times n$ .

Postupovať budeme nasledovne: vyberieme si malý rozmer a pre všetky matice daného rozmeru zistíme, či  $h^T b$  bude alebo nebude odhadnuteľné, a to nasledovným spôsobom. Z danej matice návrhu  $B$  vytvoríme maticu lineárnej regresie  $X$  spôsobom spomenutým vyššie. Hodnota  $h^T b$  bude na základe vety (1.2) odhadnuteľná práve vtedy, keď vektor  $h = (0, 0, 0, \dots, +1, -1)^T$  patrí do stĺpcového priestoru matice  $X^T X$ .

Označme  $M = X^T X$ , túto maticu budeme nazývať informačnou maticou. Ak  $h$  patrí do stĺpcového priestoru  $M$ , potom existuje taký vektor  $v$ , že  $Mv = h$ . Na základe vety (1.5) vieme, že ak toto riešenie existuje, tak sa rovná  $M^- h$ , pričom zvolená  $g$ -inverzia môže byť ľubovoľná. Preto na zistenie, či  $h$  patrí do stĺpcového priestoru  $M = X^T X$  stačí overiť platnosť rovnosti  $M(M^- h) = h$ .

Demonštrujme algoritmus na maticiach rozmeru  $3 \times 3$ . Existuje  $2^9 = 512$  binárnych matíc typu  $3 \times 3$ , a keď sme na nich rozbehli daný algoritmus, dospeli sme k zisteniu, že pri 12 z nich hodnota  $h^T b$  NIE JE odhadnuteľná. To znamená, že optimálne matice návrhu budeme hľadať v zvyšných 500 maticiach.

Ako vyzerajú matice, pri ktorých  $h^T b$  nie je odhadnuteľné? Zoznam všetkých sa nachádza v Prílohe A, tu uvedieme 3 z nich:

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}, \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}, \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Týmto spôsobom sme pri našom výskume získali zoznam matíc, pri ktorých  $h^T b$  nie je odhadnuteľné, až do rozmeru  $4 \times 5$ . Zoznam všetkých sa nachádza v GitHub repozitári spomenutom v úvode. Na základe podobnosti matíc daného typu sme prišli s nasledovnou hypotézou.

**Tvrdenie 2.1.**  $h^T b$  je odhadnuteľné práve vtedy, keď v matici návrhu  $B$  existuje aspoň jeden riadok a aspoň jeden stĺpec taký, ktorý má aspoň jednu 0 a aspoň jednu 1.

Dôkaz: Na základe vety (1.2) vieme, že  $h^T b$  je odhadnuteľné práve vtedy, keď  $h \in \mathcal{M}(X^T)$  (keď  $h$  patrí do riadkového priestoru  $X$ ), preto s danými podmienkami môžeme pracovať ekvivalentne.

$\Rightarrow$  Nech  $h = (0, \dots, +1, -1) \in \mathcal{M}(X^T)$ . Použijeme dôkaz sporom. Nech v matici  $B$  typu  $m \times n$  neexistuje taký riadok, ktorý má aspoň jednu 0 a aspoň jednu 1, t. j. všetky riadky obsahujú buď samé 0 alebo samé 1. Ako vyzerá matica  $X$ ?

Pozrime sa na riadky matice  $X$  zodpovedajúce prvému riadku matice  $B$ . Sú to práve tie riadky, ktoré majú na prvej pozícii 1, pričom si uvedomme, že všetky ostatné riadky majú na prvej pozícii 0.

$$x_1 = (1, 0, \dots, 1, 0, \dots, 0, 1, 0),$$

$$x_2 = (1, 0, \dots, 0, 1, \dots, 0, 1, 0),$$

$$\vdots$$

$$x_n = (1, 0, \dots, 0, 0, \dots, 1, 1, 0)$$

Keďže  $h \in \mathcal{M}(X^T)$ , existuje taká lineárna kombinácia riadkov  $X$ , ktorá dokopy dáva vektor  $h$ . Aké môžu mať v tejto lineárnej kombinácii zastúpenie riadky  $x_1, \dots, x_n$ ? Označme postupne  $a_1, \dots, a_n$  koeficienty vektorov  $x_1, \dots, x_n$  v lineárnej kombinácii riadkov  $X$ , ktorá dáva vektor  $h$ . Keďže  $x_1, \dots, x_n$  sú jediné riadky matice  $X$  s 1 na prvej pozícii a vektor  $h$  má na prvej pozícii 0, musí platiť  $\sum_{i=1}^n a_i = 0$ .

Z predpokladu, že každý z riadkov matice  $B$  obsahuje buď samé 0 alebo samé 1, vyplýva, že každý z vektorov  $x_1, \dots, x_n$  má rovnaké posledné dve hodnoty, a to buď  $(0, 1)$  alebo  $(1, 0)$ . Z toho spolu s rovnosťou  $\sum_{i=1}^n a_i = 0$  vyplýva, že lineárna kombinácia  $\sum_{i=1}^n a_i x_i$  bude mať na posledných dvoch miestach hodnoty  $(0, 0)$ , a teda nijako "neprispeje" do vektoru  $h$ . Vektor  $h$  teda nemožno napísať ako lineárnu kombináciu riadkov  $X$ , čo je podmienka odhadnuteľnosti.

Analogicky môžeme postupovať pre všetky riadky matice  $B$ . Týmto spôsobom pokryjeme všetky riadky matice  $X$ , vďaka čomu dospejeme k záveru, že z riadkov  $X$  nie je možné lineárnou kombináciou zostrojiť vektor  $h$ . To nám dáva spor s predpokladom



$h \in \mathcal{M}(X^T)$ . Preto nemôže platiť, že všetky riadky matice  $B$  obsahujú buď samé 0 alebo samé 1, a platí opačné tvrdenie, t.j. v  $B$  existuje aspoň jeden riadok taký, ktorý obsahuje 0 aj 1.

Rovnakou úvahou pre stĺpce  $B$  dospejeme k záveru, že  $B$  musí takisto obsahovať aspoň jeden stĺpec obsahujúci 0 aj 1.

$\boxed{\Leftarrow}$  Nech matica  $B$  je taká, že existuje aspoň jeden riadok a aspoň jeden stĺpec také, že majú 0 aj 1. Zostrojíme takú lineárnu kombináciu riadkov  $X$ , ktorá sa bude rovnať  $h$ .

Nech riadok, ktorý má 0 aj 1, je  $i$ -ty v poradí a stĺpec s rovnakou vlastnosťou je  $j$ -ty v poradí. Bez ujmy na všeobecnosti, nech na  $ij$ -tom mieste matice  $B$  sa nachádza 0, teda  $B_{ij} = 0$ .

Potom v  $i$ -tom riadku  $B$  sa určite nachádza hodnota  $B_{ik} = 1$  a v  $j$ -tom stĺpci sa nachádza hodnota  $B_{lj} = 1$ , pričom, samozrejme,  $k \neq j$  a  $l \neq i$ . Označme  $x_{ij}, x_{ik}, x_{lj}, x_{lk}$  riadky matice  $X$  prislúchajúce prvkom  $B_{ij}, B_{ik}, B_{lj}, B_{lk}$  matice  $B$ . Potom:

$$x_{ij} - x_{ik} - x_{lj} + x_{lk} = (0, 0, 0, \dots, 0, N_1, N_2)$$

je vektor, ktorý má na prvých  $m + n$  miestach 0, a na posledných dvoch hodnoty  $N_1$  a  $N_2$ , ktoré zatiaľ necháme bokom. Prečo má daný vektor na prvých  $m + n$  miestach 0? Vo všeobecnosti má riadok  $x_{rs}$  matice  $X$  prislúchajúci prvku  $B_{rs}$  matice  $B$  na prvých  $m + n$  miestach dve 1, jednu prislúchajúcu riadku, druhú stĺpcu matice  $B$ . Konkrétne, riadok  $x_{rs}$  má 1 na  $r$ -tom mieste a  $(m + s)$ -tom mieste.

Súčet  $x_{ij} + x_{lk}$  má teda štyri jednotky na miestach  $i, j, m + j, m + k$ . Súčet  $x_{ik} + x_{lj}$  má jednotky na tých istých miestach, preto  $x_{ij} - x_{ik} - x_{lj} + x_{lk} = x_{ij} + x_{lk} - (x_{ik} + x_{lj})$  má na prvých  $m + n$  miestach 0.

Takto to vyzerá, keď zanedbáme stĺpce so samými nulami:

$$\begin{aligned} &+(1, 0, 1, 0) \\ &-(1, 0, 0, 1) \\ &-(0, 1, 1, 0) \\ &+(0, 1, 0, 1) \\ &= (0, 0, 0, 0) \end{aligned}$$

Ako vyzerá chvost  $(N_1, N_2)$ ? Keďže  $B_{ij} = 0$  a  $B_{ik} = B_{lj} = 1$ , výraz  $x_{ij} - x_{ik} - x_{lj}$  nám vytvorí chvost  $(1, -2)$  (prípadne  $(-2, 1)$ , na tom ale nezáleží). Potom na základe toho, či na mieste  $B_{lk}$  matice návrhu bola 0 alebo 1, nám výraz  $x_{ij} - x_{ik} - x_{lj} + x_{lk}$  dá na posledných dvoch miestach  $(2, -2)$  alebo  $(1, -1)$ . Vo všetkých prípadoch vektor  $h$  dostaneme ihneď, prípadne po prenasobení konštantou. Našli sme teda lineárnu kombináciu riadkov  $X$ , ktorá nám dala vektor  $h$ , preto  $h \in \mathcal{M}(X^T)$ .  $\square$

## 2.2 Optimálnosť návrhu modelu

Teraz sa pokúsime nájsť modely, ktoré pre nás budú v určitom zmysle optimálne. Najprv ale definujme, čo presne optimalita modelu znamená.

**Definícia 2.2.** Binárnu maticu  $B$  návrhu modelu,  $B \in \mathbb{R}^{m \times n}$ , nazývame optimálnou, ak pre všetky binárne matice  $C \in \mathbb{R}^{m \times n}$  platí:

$$h^T M_B^- h \leq h^T M_C^- h$$

kde  $h = (0, 0, \dots, +1, -1)$  je vektor zodpovedajú lineárnej kombinácii dvoch efektov,  $M_B$  a  $M_C$  sú informačné matice prislúchajúce návrhom  $B, C$  a  $M_B^-, M_C^-$  sú ich ľubovoľné pseudoinverzie.

Nájsť optimálne matice v zmysle definície (2.2) je cieľom tejto práce. Podobne ako pri skúmaní, či hodnotu  $h^T b$  vôbec budeme môcť odhadnúť, preskúmame najprv matice malého rozmeru a na základe výsledkov vyslovíme hypotézu pre všeobecný rozmer.

Postupovať budeme nasledovne: zoberieme si malý rozmer a pre všetky binárne matice tohto rozmeru spočítame hodnotu  $h^T M^- h$ , pokiaľ je to možné (vid' tvrdenie (2.1)). Z množiny týchto hodnôt vyberieme minimum. Optimálne návrhy budú tie, ktorých prislúchajúca hodnota bude práve toto minimum.

Optimálne návrhy sme určili do rozmeru  $6 \times 5$ , tu sú niektoré z nich pre štvorcové rozmery:

$$\begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Ďalšie matice sa nachádzajú v prílohe B. Na základe týchto návrhov sme pre štvorcové matice prišli s nasledovnou hypotézou.

**Hypotéza 2.1.** *(pre štvorcové matice)*

*Binárna matica návrhu  $B$  rozmeru  $2k \times 2k$  je optimálna v zmysle definície (2.2) práve vtedy, keď má v každom riadku a v každom stĺpci práve  $k$  jednotiek (a práve  $k$  núl).*

*Binárna matica návrhu  $B$  rozmeru  $(2k + 1) \times (2k + 1)$  je optimálna práve vtedy, keď všetky jej riadky aj stĺpce majú rovnaký počet jednotiek, a ten počet je buď  $k$  alebo  $k + 1$ .*

**Záver**

TODO

## Zoznam použitej literatúry

- [1] Pázman, A., Lacko, V.: *Prednášky z regresných modelov*, Vydavateľstvo UK, Bratislava, 2012, 2015

## Príloha A