



Počítačová analýza vícerozměrných dat v příkladech

Milan Meloun
Jiří Militký
Martin Hill

ACADEMIA

**Milan Meloun
Jiří Militký
Martin Hill**

Počítačová analýza vícerozměrných dat v příkladech

**v oborech přírodních, technických
i společenských věd**

ACADEMIA 2005

O autorech

Prof. RNDr. **MILAN MELOUN**, DrSc., (*1943) vystudoval přírodovědeckou fakultu Univerzity J. E. Purkyně (dnešní Masarykova) v Brně 1965. Je profesorem analytické chemie a chemometrie na katedře analytické chemie Chemickotechnologické fakulty Univerzity v Pardubicích. Vyučoval statistické metody a analytickou chemii na Bagdádské univerzitě v Iráku a na Královské technice The Royal Institute of Technology ve Stockholmu. Je autorem a spoluautorem více než 150 originálních sdělení, 15 monografií a 11 vysokoškolských učebnic, 8 patentů a zlepšovacích návrhů a na konferencích přednesl více než 200 přednášek. Byl víc jak 10 let členem redakčních rad zahraničních odborných časopisů *Talanta* a *Analytica Chimica Acta* a je předsedou sekce chemometrie při České společnosti chemické.

Většina jeho publikací se týká počítačově orientované analýzy instrumentálních dat při studiu rovnováh v roztocích a určování reakčního produktu. Knižně se uvedl dvoudílnou monografií *Computation of Solution Equilibria*, Folia UJEP Brno 1984 (spoluautor Josef Havel), která pak doplněná o extrakční rovnováhy (spoluautor Erik Högfeltdt) vyšla v roce 1988 v anglickém nakladatelství Ellis Horwood, Chichester. K této problematice se spolu s Josefem Havlem vrátil ve dvou obsáhlých kapitolách monografie *Computational Methods for the Determination of Formation Constants*, Plenum Press, New York 1985, jejímž editorem byl David Leggett.

Ve vědecké práci prof. Melouna představuje počítač spojovací článek mezi statistikou a analytickou chemií. Výsledkem je analytická chemometrie, předmět, který přednáší na Univerzitě v Pardubicích od roku 1978. Dokladem je řada učebních textů, např. *Chemometrie*, VŠCHT Pardubice 1985. V zahraničí přednášel dle textů svých učebnic, a to při dvouletém pobytu na univerzitě v Bagdádu *Data Analysis by Statistical and Computing Technique*, University Baghdad Press, Baghdad 1980. Při hostování na Královské univerzitě ve Stockholmu používal svou sbírku příkladů *Introduction to Chemometrics*, která je postavena na interaktivní analýze pomocí statistického systému STATGRAPHICS. Vlastní přístup k analýze experimentálních dat se promítá i do kapitoly *Chemometrics in the Instrumental Laboratory* v monografii, vydané editorem Jaroslavem Churáčkem: *Advanced Instrumental Methods of Chemical Analysis*, Academia, Praha 1993, nebo kapitoly *Hodnocení analytických výsledků* ve Vlácilově sbírce *Příklady z chemické a instrumentální analýzy*, SNTL, Praha 1983. Obsáhla rešerše a zkušenosti spoluautora Jiřího Militkého přinesly řadu novějších postupů ze statistické analýzy dat, průzkumové analýzy a především interaktivní přístup k analýze dat na osobním počítači. Společně tak vzniklo první vydání učebnice *Chemometrie – Zpracování experimentálních dat na IBM PC*. Text byl přeložen do angličtiny a po doplnění o kapitoly vícerozměrné statistiky Michele Forinou vyšel postupně jako dvoudílná učebnice u nakladatelství Ellis Horwood, Chichester 1991, s názvem *Chemometrics for Analytical Chemistry – Volume I. PC-Aided Statistical Data Analysis, Volume II. Regression Model Building and Testing*.

Na Univerzitě v Pardubicích přednáší v řádném studiu chemometrii, organizuje postgraduální a dvouleté licenční studium chemometrie a i krátkodobé intenzivní kurzy chemometrie pro aplikaci v průmyslu. V těchto formách studia chemometrie jsou užívány především 2 učebnice Milan Meloun a Jiří Militký: *Statistické zpracování experimentálních dat*, Finish, Pardubice 1992, PLUS, Praha 1994 a East Publishing 1998, Academia, Praha 2004 a Milan Meloun, Jiří Militký: *Kompendium statistického zpracování dat*, Academia Praha 2002.

Prof. Ing. **JIRÍ MILITKÝ**, Csc., EURING, (*1949) ukončil fakultu textilní, specializaci textilní chemie, na VŠST v Liberci roku 1973 s vyznamenáním. V letech 1974 až 1976 pracoval ve Státním výzkumném ústavu textilním v Liberci, v oddělení matematického modelování textilních struktur. V letech 1976 až 1989 pracoval ve Výzkumném ústavu zušlechťovacím ve Dvoře Králové nad Labem, kde se věnoval převážně zpracování experimentálních dat s využitím výpočetní techniky. Od roku 1990 je vedoucím katedry textilních materiálů na Technické univerzitě v Liberci. V roce 1982 obhájil kandidátskou disertační práci z oblasti fyziky textilních vláken. V roce 1989 byl jmenován docentem a v roce 1992 se habilitoval. V prosinci 1993 byl jmenován řádným profesorem. Od 1.2.1994 do r. 1999 zastával funkci děkana fakulty textilní TUL. Od roku 1999–2003 zastával funkci prorektora pro vědu, výzkum a zahraniční spolupráci a od r. 2003 je opět děkanem fakulty textilní. V roce 1995 byl jmenován akademikem Ukrajinské akademie inženýrských věd a v r. 1996 obdržel titul EURING. Je členem několika vědeckých společností (The Textile Institute, JČMF) a pracuje ve výboru sekce chemometrie při České společnosti chemické.

Jeho publikační činnost zahrnuje oblasti textilního inženýrství, modelování kinetických procesů v pevné fázi a zpracování experimentálních dat. Je autorem nebo spoluautorem 606 vědeckých příspěvků (publikací, monografií, referátů a článků). Jeho první kniha *Modifikovaná PES vlákna* (spoluautoři Jiří Kryštůfek, Jiří Vaniček a Oldřich Hartych) vyšla v SNTL v roce 1984. Zcela přepracované a rozšířené vydání bylo publikováno nakladatelstvím Elsevier v roce 1991. S Jiřím Kryštůfkem zpracoval knihu *Barvení akrylových vláken a směsí*, která vyšla v Praze v nakl. SNTL v roce 1987. Ve spolupráci s Milanem Melounem publikoval učebnice a monografie z oblasti využití interaktivních statistických metod v chemometrii. Jiří Militký publikoval celkem 10 knih, z nichž tři jsou zaměřeny do oblasti zpracování experimentálních dat s využitím výpočetní techniky. Moderní metody interaktivní statistické analýzy dat zpracoval do rozsáhlého seriálu příruček *Statistické metody v textilní praxi I – IV*, vydaného v letech 1982 až 1985 v Domě techniky Pardubice. Přehled metod regrese a matematického modelování publikoval v seriálu skript *Tvorba matematických modelů I – IV*, vydaných v letech 1983 až 1989 v Domě techniky Ostrava. Vytvořil systém programů pro zpracování experimentálních dat v jazyce HPL. Tyto programy jsou charakteristické tím, že, kromě stránky statistické, vycházejí vždy nejdříve z ověřování předpokladů o modelech, datech a použité metodě a využívají také progresivních numerických postupů (zejména v oblasti lineární a nelineární regrese). Tyto algoritmy se později staly jádrem originálního programového systému ADSTAT.

Přednášel na odborných akcích v USA, Kanadě, Japonsku, Tchaj-wanu, Austrálii, Hongkongu, Vietnamu, Egyptě, Maroku a v řadě zemí Evropy. Je aktivně zapojen do činnosti v řadě odborných společností. Je členem výboru „International Textile Academy“, české pobočky „The Textile Institute“ a předsedou českého monitorovacího výboru FEANI. Je členem výboru České statistické společnosti a České chemické společnosti.

Ing. **Martin Hill**, DrSc. (1962). V roce 1986 absolvoval VŠCHT v Praze, obor chemické a energetické zpracování paliv. V letech 1987 až 1992 byl vědeckým aspirantem v Ústavu geotechniky v Praze. V roce 1992 obhájil hodnost kandidát technických věd za práci v oblasti kinetiky zplynování tuhých paliv. Je autorem nebo spoluautorem přes 150 odborných prací převážně v mezinárodních impaktovaných časopisech. Většina z těchto prací je postavena na využití pokročilých metod statistického zpracování vícerozměrných lékařských a biochemických dat. V roce 1996 absolvoval stáž v Conservatoire National des Arts et Métiers v Paříži, kde se zabýval antiglukokortikoidními účinky derivátů DHEA a pregnenolonu. Od roku byl nebo je hlavním řešitelem sedmi grantových projektů z nichž většina byla zařazena do soutěže o cenu ministra zdravotnictví a spoluřešitelem řady dalších z nichž jeden se umístil na prvním místě v uvedené soutěži. Za sérii metodických prací v oblasti analýzy méně běžných steroidů získal v roce 1998 Cenu Endokrinologické společnosti. V roce 1999 absolvoval dvouleté licenční studium chemometrie na Univerzitě Pardubice, které uzavřel prací na téma vícerozměrné statistické metody v analýze dat s negaussovským rozdělením. V roce 2001 obhájil na základě práce „Steroid analysis and data treatment for physiological and diagnostic conclusions“ vědeckou hodnost „Doktor chemických věd“ v oboru Analytická chemie. Má zkušenosti v oblasti kapalinové a plynové chromatografie, hmotové spektrometrie, imunoanalýzy, chemometrie a statistiky, které uplatňuje jako vědecký pracovník v Oddělení steroidních hormonů Endokrinologického ústavu, kde je zaměstnán od r. 1992. V současnosti se zabývá vývojem nových metod analýzy steroidů ve speciální steroidní diagnostice a aplikací moderních metod statistické analýzy při zpracování biochemických dat.

Obsah

O autorech	05
Obsah	09
Předmluva	13
1 Charakter vícerozměrných dat	15
1.1 Nepřímá pozorování a korelace	15
1.2 Zdrojová matice dat	16
1.3 Druhy dat	17
1.3.1 Nestrukturovaná data	17
1.3.2 Strukturovaná data – jedna skupina závisle proměnných	18
1.3.3 Strukturovaná data – více skupin závisle proměnných	18
1.4 Odhady parametrů polohy, rozptýlení a tvaru	19
1.5 Vybočující body	22
2 Předúprava vícerozměrných dat	31
2.1 Formy standardizace dat	31
2.2 Užití statistických vah	35
2.3 Průzkumová analýza vícerozměrných dat	36
2.3.1 Zobrazení vícerozměrných dat	36
1. Zobecněné rozptylové grafy	37
2. Symbolové grafy	38
2.3.2 Ověření normality	42
3 Metody k odhalení struktury ve znacích a objektech	49
4 Analýza hlavních komponent (PCA)	61
4.1 Zaměření metody PCA	61
4.2 Podstata metody PCA	62
4.3 Cíl metody hlavních komponent PCA	62
4.4 Grafické pomůcky analýzy hlavních komponent	66
4.4.1 Cattelův indexový graf úpatí vlastních čísel (Scree Plot)	66
4.4.2 Graf komponentních vah, zátěží (Plot Components Weights)	68
4.4.3 Rozptylový diagram komponentního skóre (Scatterplot)	68
4.4.4 Dvojný graf (Biplot)	69
4.4.5 Graf reziduí jednotlivých objektů	70
4.4.6 Graf celkového reziduálního rozptylu všech objektů	70
4.5 Diagnostika metody hlavních komponent	70
4.6 Řešení častých problémů v PCA	71
5 Faktorová analýza (FA)	99
5.1 Zaměření metody FA	99
5.2 Podstata metody faktorové analýzy FA	101
5.3 Grafické pomůcky faktorové analýzy FA	105
5.4 Diagnostikování metodou FA	105
1. Cíle faktorové analýzy	106
2. Formulace úlohy faktorové analýzy	106
3. Předpoklady faktorové analýzy	107
4. Nalezené řešení a dosažená těsnost proložení	108
5. Interpretace výsledků	111
6. Ověření výsledků	114
7. Využití výsledků faktorové analýzy	114

8. Diagnostikování problémů faktorové analýzy	117
6 Kanonická korelační analýza CCA	145
6.1 Zaměření metody CCA	145
6.2 Podstata metody CCA	145
6.2.1. Test významnosti kanonických korelací	148
6.2.2. Vysvětlení kanonických proměnných	148
6.2.3. Analýza redundance	148
6.2.4. Grafické pomůcky	149
6.3 Postup diagnostikování CCA	149
1. Cíle kanonické korelační analýzy	149
2. Formulace úlohy kanonické korelační analýzy	149
3. Předpoklady kanonické korelační analýzy	150
4. Nalezené řešení a dosažená těsnost proložení	150
5. Interpretace výsledků	151
6. Ověření výsledků	152
7. Diagnostikování problémů kanonické korelační analýzy	152
7 Diskriminační analýza (DA)	179
7.1 Zaměření metody DA	181
7.2 Zařazovací pravidla DA	182
7.3 Lineární (LDA) a kvadratická (QDA) diskriminační funkce	183
1. Lineární diskriminační funkce LDA	184
2. Kvadratická diskriminační funkce QDA	189
7.4 Užití kanonické korelace v diskriminační analýze	192
7.5 Úprava prahového bodu	193
7.6 Volba znaků, diskriminátorů	193
7.7 Kvalita zařazení objektů do tříd	197
7.8 Logistická diskriminace	198
7.9 Průběh diagnostikování DA	200
1. Cíle diskriminační analýzy	200
2. Formulace úlohy a volba diskriminátorů	201
3. Předpoklady diskriminační analýzy	202
4. Nalezené řešení a dosažená těsnost proložení	203
5. Interpretace výsledků	209
6. Ověření výsledků	211
11 Korespondenční analýza (CA)	397
11.1 Zaměření metody CA	397
11.2 Podstata metody CA	398
11.3 Postup korespondenční analýzy	400
1. Cíle korespondenční analýzy	400
2. Formulace úlohy korespondenční analýzy	401
3. Předpoklady korespondenční analýzy	401
4. Nalezené řešení a dosažená těsnost proložení	401
5. Interpretace výsledků	401
6. Ověření výsledků	402
Literatura	425
Hodnocení	430
Proč právě STATISTICA?	430
Produkty STATISTICA	432
STATISTICA Neuronové sítě Cz	433
STATISTICA Analýza síly testu	433
Průmyslová řešení a nástroje Six Sigma	434

STATISTICA Diagramy pro řízení jakosti Cz	434
STATISTICA Analýza procesů Cz	434
STATISTICA Navrhování experimentů Cz	434
Podnikové systémy	434
STATISTICA Data Miner	434
STATISTICA Text Miner	434
STATISTICA QC Miner	435
STATISTICA Vícerozměrné statistické řízení procesů	435
WebSTATISTICA Server	436
STATISTICA Document Management System	436
WebSTATISTICA Data Warehouse	436

Předmluva

Zpracování vícerozměrných dat v technické praxi využívá poznatků *přirodních věd, matematické statistiky a informatiky* v kombinaci se speciálními počítačově orientovanými postupy. Současné výkonné osobní počítače umožňují interaktivnost při zpracování vícerozměrných dat a interpretaci získaných výsledků. To klade stále větší nároky na znalosti pracovníků, kteří data zpracovávají a analyzují. Nabídka a možnosti počítačově orientovaného statistického zpracování dat nutí ke komplexnější analýze problémů, což vede většinou i k radikální změně pohledu na metodiku jejich zkoumání.

Při zpracování reálných vícerozměrných dat se běžně naráží na řadu problémů a omezení:

- (a) rozsahy zpracovávaných dat nejsou vzhledem k rozměrnosti problémů obvykle dostatečně velké,
- (b) v datech se vyskytují výrazné nelinearity, neaditivita a vzájemné vazby, které je třeba identifikovat a popsat,
- (c) rozdělení dat jen zřídka odpovídá normálnímu běžně předpokládanému ve standardní statistické analýze,
- (d) v datech se vyskytují podezřelá a odlehlá měření a různé heterogenity,
- (e) statistické modely se často tvoří na základě předběžných informací z dat (datově orientované přístupy),
- (f) existuje jistá neurčitost při výběru modelu, popisujícího chování dat.

To vše klade zvýšené nároky na techniky umožňující snižování rozměrnosti, hledání vnitřních skrytých vazeb v datech respektive vhodné zobrazení vícerozměrných dat. Pro tyto účely existuje celé spektrum méně či více dokonalých komplexních programů a programových systémů. Některé jsou budovány jako univerzálně použitelné a některé jsou zaměřené na specifické oblasti (chemometrie, biometrie, ekonometrie, medicínská statistika, obchodní statistika, statistika pro sociology, psychology, atd.). Jejich účinné využití není možné bez znalostí alespoň základů příslušných metod, které jsou základem pro interpretaci výsledků. Vlastní interpretace má obvykle jak statistickou stránku tak i stránku související

s daným oborem.

Knih je výsledkem snahy autorů překlenout rozdíl mezi pokrokem ve vývoji softwarových balíků obsahujících metody statistického zpracování vícerozměrných dat na jedné straně a v praxi stále nedostatečně využívanými možnostmi jejich užitečné aplikace na straně druhé. Přesto, že manipulace s daty je u moderních přístrojů a systémů snadná, vlastní výpočet je zpravidla otázkou sekund, počítačové výstupy z příslušného programu bývají přehledné a jejich interpretace ve většině případů nevyžaduje detailní matematickou znalost metodiky. Mezi odbornou, nematematicky zaměřenou veřejností stále přetrvávají obavy z využívání vícerozměrné statistické analýzy dat. Problémy činí především vlastní formulace úlohy a dostatečná interpretace výsledků. Specifickým problémem bývá také nepochopení výsledků a interpretace výstupů statistické vícerozměrné analýzy dat u části odborné veřejnosti. Přesto, že výsledky jsou pro člověka byť jen s povrchní znalostí uvedených metodik zcela zřejmé, stává se, že odborníci v daném oboru nejsou z důvodů vlastní neznalosti schopni těchto informací využít. Přitom efektivita využití informací z dat při použití vícerozměrných statistických technik je podstatně větší než u jednorozměrných dat.

Cílem této knihy je zpřístupnit vícerozměrné statistické techniky, které jsou dnes již běžnou součástí statistických softwarových balíků široké nematematické přírodovědné a technické veřejnosti a zejména studentům. Zvláštní důraz je kladen na sestavení úlohy a na interpretaci výsledků. Proto jsou těžištěm knihy návody s podrobnými postupy a komentované příklady s množstvím grafických výstupů umožňujících kromě diagnostiky kvality vstupních dat také snadnou interpretaci výsledků. Kromě sestavení úlohy bývá, zvláště u dat z oblasti biochemie a medicíny, problémem jejich nesymetrie rozdělení, nekonstantní rozptýlení a výskyt nehomogenit. Z těchto důvodů byla pozornost věnována i kvalitě vstupních dat a možnostem jejich transformace aby tak byly splněny základní předpoklady správného provedení vícerozměrné statistické analýzy. Příklady jsou voleny ze širokého spektra přírodních a technických věd, často z oblasti biochemie ale také z klinické praxe. Je na nich demonstrován zejména praktický přínos vícerozměrných statistických technik a zpravidla také jejich nezastupitelnost jednoduššími metodikami.

Knih vychází a velmi úzce souvisí s nedávno zveřejněnou publikací Meloun M., Militký J.: *Statistická analýza experimentálních dat*, Academia Praha 2004, která obsahuje výklad jednotlivých vícerozměrných metod do větší hloubky ale s menším zaměřením na počítačově orientovanou analýzu resp. výsledky jednotlivých programů. Předpokládáme, že se kniha stane základem zpřístupnění vícerozměrného statistického zpracování dat širokému okruhu čtenářů, pracovníkům přírodovědných a technických oborů, studentům, lékařům resp. dalším specialistům, kteří zpracovávají vícerozměrná statistická data.

Na závěr je naší milou povinností poděkovat všem spolupracovníkům, studentům a doktorandům, které není možné ani všechny jmenovat a kteří nám pomáhali či přispěli praktickými úlohami, radami či konstruktivní kritikou. Vděk patří také všem pedagogům a studentům řádného i licenčního studia, kteří nám poskytli cenné dotazy, podněty a připomínky k řešeným příkladům.

Milan Meloun, Jiří Militký a Martin Hill

Pardubice, Liberec a Praha, leden 2005