

# **PROSTOROVÉ ANALÝZY DAT**

**doc. Dr. Ing. Jiří Horák**

**VŠB-TU Ostrava, HGF, Institut geoinformatiky, 2015  
6.vydání**

## Obsah:

1	Prostorové analýzy - vymezení, rozdělení.....	4
1.1	Definice prostorových analýz.....	4
1.2	Historie prostorových analýz.....	5
1.3	Cíle prostorových analýz.....	6
1.4	Typy používaných metod.....	7
1.4.1	Rozdělení podle využitých postupů (aplikovaných technik).....	7
1.4.2	Rozdělení podle způsobu zpracování dat.....	8
1.4.3	Rozdělení podle typu prostorové reprezentace.....	10
1.4.4	Statistické prostorové analýzy (s.s.).....	10
1.4.4.1	Tři základní typy prostorové distribuce.....	10
1.4.4.2	Podrobnější přehled základních technik pro provádění statistických prostorových analýz	11
1.5	Vztah GIS a prostorové analýzy.....	14
2	Typy dat a vztah k prostorovým analýzám.....	15
2.1	Grafické atributy.....	15
2.2	Popisné atributy.....	16
3	Body (lokalizační data).....	19
3.1	Popisná statistika pro body.....	19
3.1.1	Charakteristiky střední polohy.....	19
3.1.2	Charakteristiky rozptýlení.....	24
3.1.3	Dostupné nástroje v programech.....	27
3.2	Inferenční statistické testy pro body.....	28
3.2.1	Kvadrantové testy náhodnosti.....	28
3.2.1.1	Test při náhodném rozmístění buněk.....	31
3.2.1.2	Index disperze.....	32
3.2.1.3	$\chi^2$ test.....	33
3.2.1.4	Problémy při kvadrantových metodách.....	33
3.2.2	Metoda nejbližších vzdáleností.....	34
3.2.2.1	Teoretický model nejbližších vzdáleností pro CSR.....	35
3.2.2.2	Variety vyhodnocení metody.....	36
3.2.2.3	Hraniční efekt.....	43
3.2.3	K funkce.....	45
3.2.3.1	Definice K-funkce.....	45
3.2.3.2	Hraniční korekce pro K funkci.....	46
3.2.3.3	Vyhodnocení K-funkce.....	47
3.3	Modelování prostorového uspořádání událostí (bodů).....	51
3.3.1	CSR - homogenní Poissonův proces.....	51
3.3.2	Heterogenní Poissonův proces.....	52
3.3.3	Coxův proces.....	52
3.3.4	Poissonův shlukovací proces.....	52
3.3.5	Markovovy bodové procesy.....	53
3.3.6	Srovnání jiných modelů než CSR s prostorovým bodovým procesem.....	54
3.4	Transformace bodové textury do kontinuálního pole.....	56
3.4.1	Kvadrantová metoda.....	56
3.4.2	Jádrové vyhlazení.....	57
3.4.2.1	Adaptivní jádrový odhad.....	63
3.4.2.2	Problémy u jádrového odhadu.....	63
3.5	Analýza vícenásobných typů událostí.....	65
3.5.1	Kvadrantová metoda s $\chi^2$ testem.....	65
3.5.2	Metoda nejbližších vzdáleností.....	66
3.5.2.1	Grafické hodnocení.....	66
3.5.2.2	Numerické hodnocení.....	66
3.5.3	K funkce.....	67

3.6	Časoprostorové analýzy.....	67
3.6.1	Popisná statistika .....	67
3.6.2	Inferenční statistika .....	68
3.6.3	$\chi^2$ test .....	68
3.6.4	Knoxův test.....	68
3.6.5	Mantelův test .....	69
3.6.6	D funkce .....	69
3.6.7	Prostor-čas-atributový analytický stroj.....	70
4	Geostatistické metody pro kontinuální pole .....	72
5	Linie.....	111
5.1	Analýza interakčních dat .....	111
5.1.1	Prostorové vazby mezi zdroji a cíly toků .....	112
5.1.2	Charakteristiky určující velikost toků ze zdrojů.....	112
5.1.3	Charakteristiky určující velikost toků do cílů.....	112
5.2	Teorie grafů a její aplikace pro analýzu interakčních dat.....	112
5.3	Hledání optimální cesty.....	116
5.3.1	Dantzigův algoritmus .....	116
5.3.2	Floydův algoritmus.....	117
5.3.3	Dijkstrův algoritmus.....	119
5.3.4	Základní algoritmus pro nalezení nejlevnější kostry .....	121
5.3.5	Hladový algoritmus .....	121
5.4	Způsoby hodnocení dopravní dostupnosti.....	121
5.4.1	Metrické míry .....	122
5.4.2	Časové míry dostupnosti .....	123
5.4.3	Topologické míry dostupnosti .....	123
5.4.4	Cenové míry dostupnosti.....	124
5.4.5	Vážené míry dostupnosti .....	124
5.4.6	Hodnocení veřejné dopravy.....	124
5.5	Lokalizační a alokační úlohy.....	126
5.6	Gravitační teorie a její zobecnění .....	130
5.6.1	Oboustranně omezený gravitační model .....	131
5.6.2	Jiné gravitační modely a aplikace.....	133
5.7	Popisná statistika .....	134
5.7.1	Popis liniového vzorku .....	134
5.7.2	Náhodnost distribuce pro liniový vzorek.....	135
5.8	Vizualizace interakčních dat.....	136
6	Polygony.....	137
6.1	Popisná statistika pro polygony.....	137
6.2	Problém plošné interpolace .....	137
6.2.1	Absolutní hodnoty .....	138
6.2.2	Relativní hodnoty .....	138
6.3	Problém regionalizace .....	139
6.4	Vyhlazování areálových dat .....	140
6.4.1	Prostorové klouzavé průměry.....	141
6.4.2	Bayesovo vyhlazení.....	144
6.5	Sledování autokorelace.....	147
6.6	Multivariační techniky v prostorových aplikacích .....	151
6.6.1	Shluková analýza.....	151
6.6.2	Analýza hlavních komponent .....	156
6.7	Regresní modely.....	156
6.7.1	Neprostorový regresní model .....	156
6.7.2	Prostorový regresní model.....	158
6.7.3	Generalizované lineární modely .....	163
	Seznam literatury.....	165

# 1 Prostorové analýzy - vymezení, rozdělení

## 1.1 Definice prostorových analýz

Většina informací, se kterými se setkáváme a které využíváme, má prostorový charakter. Jistým způsobem je vázána k určitému místu a reprezentuje ho. Toto místo je třeba chápat v širším slova smyslu - může to být bod, sada bodů, linie, sada (kolekce, svazek) linií, areál. V této souvislosti pak hovoříme o **geoinformacích**. Formalizací informace získáváme data, obdobně formalizací geoinformace získáváme **geodata** (resp. prostorová data). Analogicky bychom mohli specifikovat chápání prostorových analýz jako analýzu prostorových dat, to však není správné, protože ne každá analýza prostorových dat je prostorovou analýzou - pokud např. vytvoříme histogram úplné sady prostorových dat či vypočítáme jejich základní statistické charakteristiky, nevyužili jsme prostorový aspekt těchto dat a nejde tedy o prostorovou analýzu.

Prostorové analýzy představují kolekci technik, které vznikly v různých oborech a jejichž cílem byla analýza dat s důrazem na jejich prostorové vztahy. Významné postavení mezi těmito obory zaujímá statistika, ale řada postupů byla odvozena v geografii, geostatistice, ekonometrii, epidemiologii, v územním plánování a urbanizmu. Tyto postupy jsou používány v ještě širší škále aplikací včetně např. zdravotnictví a kriminalistiky.

Prostorové analýzy můžeme definovat následovně: **Prostorové analýzy jsou souborem technik pro analýzu a modelování lokalizovaných objektů, kde výsledky analýz závisí na prostorovém uspořádání těchto objektů a jejich vlastností.**

Objektem pro tento účel rozumíme geografické objekty a jiné objekty s prostorovou lokalizací (např. hvězdy nebo útvary v obraze), ať již fyzické nebo abstraktní povahy, velmi často i události a jevy.

*Vymezení pojmu prostorové analýzy nebylo dříve tak univerzálně chápáno a často se vztahovalo jen k určité oblasti aplikací či použitých postupů. Jako příklad můžeme uvést několik starších definic:*

*Unwin (1981): „Prostorové analýzy se zabývají uspořádáním prostorových dat na mapách (tedy bodů, linií, ploch, povrchů).“*

*Johnston, Gregory, Smith (1994): „Prostorové analýzy jsou kvantitativní (hlavně statistické) procedury a techniky aplikované v lokalizačních (umístovacích) úlohách.“*

*Goodchild (1988): „Prostorové analýzy jsou techniky umožňující popis uspořádání na mapách a srovnání 2 a více map s cílem identifikace jejich vztahů.“*

*Někteří autoři chápou termín „prostorové analýzy“ jako synonymum pojmu **kvantitativní geografie**, část z nich pak tento pojem uplatňuje jen pro tu část prostorových analýz, která využívá stochastické povahy jevů.*

Bez ohledu na konkrétní vyjádření je zjevné, že prostorové analýzy představují sadu analytických metod, vyžadujících přístup k atributům studovaných objektů i k informaci o jejich lokalizaci. Na rozdíl od jiných forem analýz tedy vyžadují prostorové analýzy atributová data i geografickou lokalizaci objektů.

Prostorové analýzy dat jsou spjaty se studiem uspořádání prostorových dat. Zvláště se zabývají vyhledáváním nových vztahů mezi uspořádáním a atributy objektů nebo geoprvky ve studované oblasti a s modelováním těchto vztahů s cílem dosáhnout jejich lepšího porozumění a předpovídání vývoje v oblasti.

Pozorované uspořádání objektů či jevů je možné označit za **texturu** (*pattern*), což je vhodnější než používání pojmu vzor nebo vzorek, které mají zavádějící homonyma (vzor jako něco co je vhodné následovat, vzorek (*example*) jako výběr jistých zástupců pro sledování celkových charakteristik).

Prostorové analýzy řeší řadu rozdílných prostorových problémů - od korekce obrazu a rozpoznání textury obrazu, přes interpolaci ovzorkovaného surovinového ložiska, průzkum prostorových a

časoprostorových shluků nehod, modelování socioekonomických trendů až po studium migrace zvířat a lidí. Díky svému rozmanitému zaměření postrádají prostorové analýzy jasný systém kodifikace nebo jasný koncepční či teoretický rámec.

Při úvahách o členění prostorových analýz musíme zohlednit i způsob organizace dat, protože některé funkce je možné aplikovat jen pro jistý typ dat.

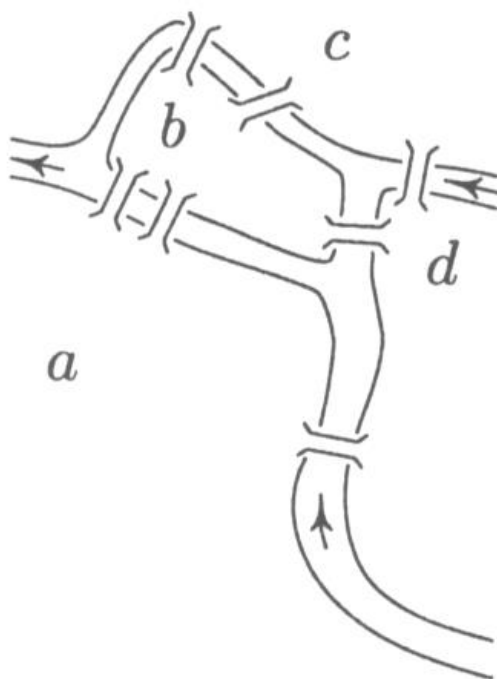
## 1.2 Historie prostorových analýz

Jak již bylo uvedeno, historie prostorových analýz je spojena s celou řadou oborů.

Snad nejstarší kořeny mají lokalizační úlohy (hledající optimální umístění), kde již v roce 1629 formuloval Fermat úlohu nalezení bodu s minimálním součtem vzdáleností od daných 3 bodů, avšak teprve v roce 1746 ji vyřešil Torricelli. V roce 1837 úlohu zobecnil Steiner pro  $n$ -bodů a v roce 1909 umožnil Weber používat různé váhy bodů v úloze (různá atraktivita bodů) a popsal její ekonomickou interpretaci.

Jeden z prvních dokladů geografické aplikace prostorových analýz je analýza Halleye. Halley již v roce 1686 zobrazil na podkladové mapě území směry větrů a monzuny v blízkosti tropického pásma a snažil se najít jejich fyzikální příčinu.

V oblasti teorie grafů již roku 1736 řešil německý matematik Euler tzv. úlohu 7 mostů (mosty v městě Kaliningrad). Úloha má za cíl nalézt takovou trasu, která obsahuje všechny hrany (tj. mosty) právě jednou a která začíná a končí ve stejném místě. V roce 1859 pak Hamilton řešil tzv. problém obchodního cestujícího, kdy je nutné navštívit všechny uzly v systému (vybranou sadu míst) a minimalizovat procestovanou trasu.



Obr. 1-1 Mosty ve městě Kaliningrad

Řada metod byla odvozena v epidemiologii, kde první uváděnou analýzou je studie doktora Snowa z roku 1854 o způsobu přenosu cholery ve vztahu ke změnám v pozorované mortalitě v Londýně. Společně zobrazil umístění studní, výskyt nemoci a úmrtí a síť ulic a snažil se najít vztah mezi nimi (obr. 1-2). Jednou z ranných geografických analýz je i použití geografické analýzy k prokázání vztahu mezi nedostatkem slunečního záření a výskytem křivice, kterou prováděl Palm v roce 1890.



Obr. 1-2 Mapa Dr. Snowa - úmrtí na cholera (tečky) a výskyt studní (křížky, nakažená studna v centru obrázku), londýnská čtvrt Soho 1854 (Goodchild, Janelle 2004)

Mezi prostorové analýzy je možné zařadit i využití kvadrantové metody, kdy v roce 1907 Student pomocí ní sledoval distribuci částic v kapalině a zjistil, že počet částic v kvadrantu odpovídá Poissonově distribuci.

Velmi významný metodologický přínos také představuje oblast geostatistiky, jejíž základy formuloval Matheron v letech 1962 a 1963 s aplikací poznatků z ložiskové geologie, výpočtu zásob ložisek rud, matematiky a statistiky.

### 1.3 Cíle prostorových analýz

Cíle prostorových analýz se opět značně liší podle oblasti aplikací a je obtížné nalézt univerzální rozdělení. Řada autorů se pouze omezuje na výčet cílů, které reprezentují oblast zájmu pro danou aplikaci. Jako příklad může posloužit **výběr cílů síťových úloh** dle Laurini, Thompson (1994):

- 1) Najděte všechny možné trasy pro nákladní automobil v silniční síti mezi počátkem a cílem cesty.
- 2) Najděte místo, kde je "služba" přerušena díky přerušení nebo špatné funkci sítě.
- 3) Najděte takovou trasu, na které se nachází nejvíce zákazníků.
- 4) Najděte nejbližší elektrickou rozvodnu nebo telefonní ústřednu pro nového zákazníka.
- 5) Umístěte nový servis v dálniční síti (typická lokalizační úloha).
- 6) Rozmístěte děti do nejbližších škol na základě dopravního času při cestování ulicemi (typická alokační úloha)

Při obecnějším vymezení cílů prostorových analýz můžeme rozlišit následující cíle:

- 1) Popis objektů resp. událostí ve sledovaném prostoru (včetně popisu jejich uspořádání - tj. textury). Zahrnuje odvození statistických charakteristik pozorované textury geoprveků (bodů, linií či areálů) a jejich srovnání; dále testování, zda je pozorovaná distribuce významně odlišná od určité hypotetické textury (což je významné pro následující interpretaci procesů); zkoušení

prostorových vztahů a vazeb mezi entitami, ale i běžný popis vývoje pole např. výpočet hodnoty v neznámých místech (interpolace).

Zajímá nás, proč jsou určité fenomény více seskupeny v některých místech, zda to není jen vliv náhody, jak lze porovnat texturu v různých oblastech, jak lze takový rozdíl kvantifikovat, zda dochází ke změnám v čase.

Někteří autoři kritizují tento cíl, protože většina analýz končí u takového popisu a už se nezabývá vysvětlením procesů, které vedly k pozorovanému uspořádání. Navíc málokdy v přírodě odpovídá vzorek teoretickému modelu.

Zde však nečekáme, že situaci bude přesně vystihovat teoretická distribuce, popis nám ale slouží k nalezení klíčových faktorů, které vedou ke vzniku určitého uspořádání.

- 2) Výběr určitého místa na základě splnění jisté sady podmínek (či obecněji podle jistého rozhodovacího schématu) nebo zkoumání míry splnění daných podmínek v určitém místě nebo území.
- 3) Interpretace procesů, které vedly k pozorovanému stavu uspořádání objektů či událostí ve sledovaném prostoru (systematický průzkum), např. interpretace vzniku pozorovaného uspořádání bodů, vysvětlení vývoje území v čase (jak střední hodnoty tak variability).
- 4) Optimalizace uspořádání objektů/jevů ve sledovaném prostoru např. na lokalizační a alokační úlohy, volba způsobu distribuce toků (rozmístění zaměstnaných, dětí do škol, zboží), ale také např. návrh vhodného systému vzorkování.
- 5) Zlepšení schopnosti předpovídat a kontrolovat objekty či události ve sledovaném prostoru (využití prediktivních modelů).
- 6) Redukce původního množství dat do menší, úspornější a přehlednější sady dat. Provádíme např. generalizaci původních dat pro lepší popis sledovaného jevu nebo jen za účelem snadnější manipulace.

Uvedený přehled cílů prostorových analýz jistě není a ani nemůže být úplný, protože s rozvojem geoinformačních technologií se nacházejí nové formy uplatnění prostorových analýz a s tím i nové cíle.

## 1.4 Typy používaných metod

K rozdělení metod prostorové analýzy je možné přistoupit z různých hledisek, např. z hlediska využitých postupů, způsobu zpracování, úrovně zpracování, počtu současně studovaných jevů či podle typu reprezentace prostorových objektů či jevů.

Uvedeme rozdělení podle:

- použitých postupů (aplikovaných technik),
- způsobu zpracování dat,
- typu prostorové reprezentace.

Významnou skupinu tvoří statistické prostorové analýzy (s.s.).

### 1.4.1 Rozdělení podle využitých postupů (aplikovaných technik)

Má zjevnou úzkou vazbu na členění disciplin, ve kterých techniky vznikaly. Podle tohoto postupu dělíme prostorové analýzy na:

- statistické prostorové analýzy dat (*spatial statistics*) - úzce spjatý s matematickou statistikou,
- mapová analýza - ve smyslu mapové algebry, především překryvné operace
- metody matematického modelování - např. tvorba a analýzy multivariačních či regresních modelů

- interpolační metody
- lokalizační a alokační metody
- síťové analýzy - ve smyslu geografických analýz, kde předmětnou sítí bývají dopravní, hydrologické či inženýrské sítě
- ostatní analýzy okolí a spojitosti - např. techniky zpracování obrazu, používané pro získání geometrických charakteristik obrazu či textury; gravitační analýzy apod.

V případě síťové analýzy je potřebné odlišit geografickou síťovou analýzu od homonyma síťové analýzy z operačního výzkumu, uplatňované zvláště při projektovém řízení. Každá z nich má jiný význam, i když v obou případech se zpravidla řeší problém optimalizace grafů. Analýzy většinou využívají poznatky z teorie grafů, v některých případech to však není možné (nebo to není efektivní) - např. při hledání cesty terénem, který je popsán pomocí kontinuálního modelu.

Při lokalizačních úlohách řešíme problém optimálního umístění (rozmístění, lokalizace) objektů. Naproti tomu alokační úlohy řeší problém zásobování (model řízení zásob). Tyto metody se často řadí mezi síťové, protože využíváme šíření v síti. Příkladem lokalizačního problému je výběr vhodného místa pro stavbu obchodního centra, kdy vycházíme ze známé lokalizace zákazníků a známé dopravní sítě. Optimalizačním kritériem je pak maximalizace zisku plánovaného obchodního centra.

#### 1.4.2 Rozdělení podle způsobu zpracování dat

Z hlediska **způsobu zpracování dat** lze rozlišit metody (odpovídající 3 základním formám zpracování dat):

1. zobrazovací
2. průzkumové
3. modelovací

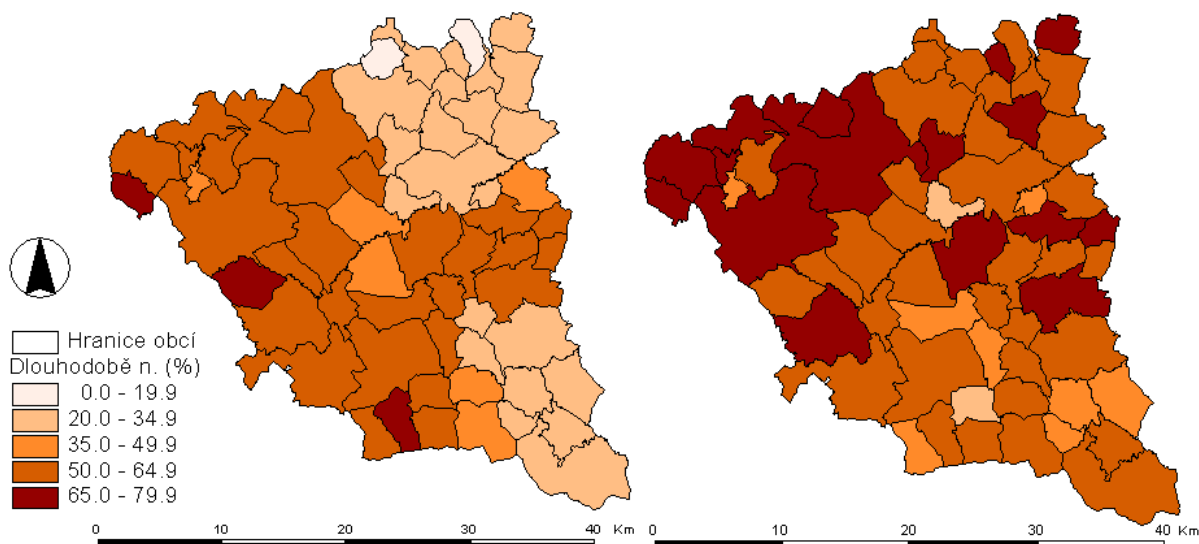
Zpracování dat můžeme chápat jako určitou posloupnost kroků, tvořících určité etapy či úrovně zpracování dat. Prvním krokem v analýze prostorových dat bývá vizualizace primárních či základním způsobem adjustovaných dat. Další etapy, tj. průzkumová či modelovací, mohou nebo nemusí být realizovány, ale pokud se provádí, tak až po této základní zobrazovací etapě. Přitom výsledky těchto etap mohou být (a zpravidla jsou) vizualizovány pomocí zobrazovacích metod.

Mohli bychom tedy zobrazovací metody ještě rozdělit podle jejich použití před či po vlastním zpracování dat na přípravné (pre-modelling) a výstupní (post-modelling) zobrazovací metody.

**Ad 1) Zobrazovací (vizualizační) metody** se zaměřují na zobrazení prostorových dat bez modifikace grafické složky dat. Nevyžadují statistické zpracování dat, snad pouze při vymezení hranic tříd. Při zobrazovacích metodách se používá často vytváření map, kartogramů nebo kartodiagramů. Výsledné mapové kompozice dokumentují objekty a jevy ve sledovaném území a jsou vizuálně interpretovány. Významné je sledování rozmístění objektů, evidence výskytu shluků nebo anomálně deficitních míst, hledání příčin takového stavu, sledování trendů v prostoru a čase, vliv jednotlivých faktorů na jejich výskyt a uspořádání a jejich korelace.

Příkladem aplikace vizualizačních metod je tvorba statistických map, zvláště kartogramů a kartodiagramů.

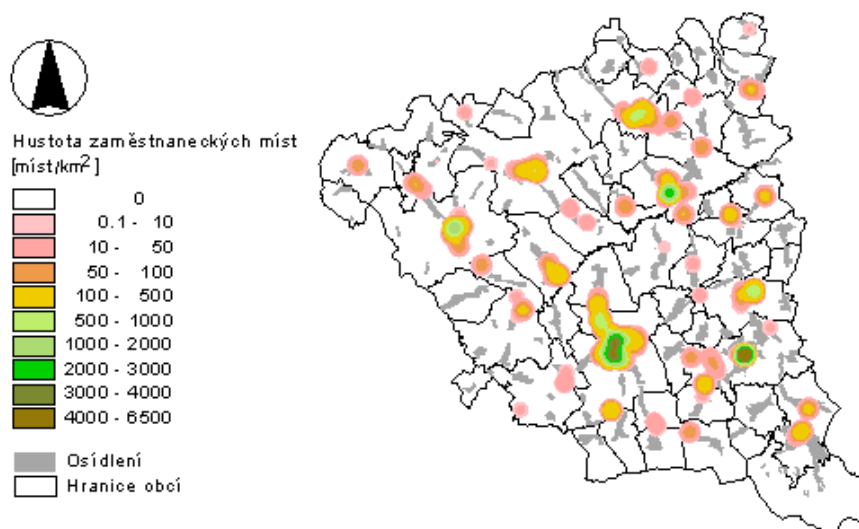




Obr. 1-3 Podíl dlouhodobě nezaměstnaných z celkového počtu nezaměstnaných v obcích okr. Nový Jičín (vlevo k 30.6.1999, vpravo k 30.6.2000). Data UP Nový Jičín.

**Ad 2) Průzkumové metody** nezobrazují původní data, ale používají data prostorově modifikovaná. Využívají tedy takové formy zpracování dat jako je např. vyhlazení, transformace, filtrace apod. Často provádějí sumarizaci hodnot. Zabývají se průzkumem mapové textury, vztahů a detekcí anomálií.

Zobrazovací a průzkumné metody se využívají v průzkumné analýze dat, jejichž cílem je základní průzkum vlastností dat. K hlavním rysům průzkumné analýzy dat patří identifikace vlastností dat potřebných pro účely detekce textury dat, formulování hypotéz na základě evidovaných dat a pro některé aspekty ocenění modelu (např. věrohodnosti modelu). Doporučuje se využívat jednoduché, intuitivní a statisticky robustní metody.



Obr. 1-4 Hustota zaměstnaneckých míst v okr. Nový Jičín u významných zaměstnavatelů. Data UP Nový Jičín.

**Ad 3)** Třetí skupinu metod prostorové analýzy představují **modelovací metody**. Jejich základem je vytvoření vhodného modelu, ověření jeho vhodnosti pro sledovaný účel (např. vysvětlení vlivu jistých faktorů) následně využití parametrů modelu pro interpretaci jevů, vztahů nebo využití modelu

pro zkoumání následků jistých změn parametrů, časového vývoje apod. V oblasti modelování může k nejdůležitějším úlohám patřit prediktivní modelování, využívající často lokalizační a alokační úlohy.

### 1.4.3 Rozdělení podle typu prostorové reprezentace

Dělení metod podle **typu prostorové reprezentace** odráží rozdělení prostorové reprezentace. Metody lze tedy dělit na metody vhodné pro kontinuální reprezentaci a pro diskrétní reprezentaci, kterou lze dále dělit podle typu geometrických primitiv použitých k reprezentaci.

### 1.4.4 Statistické prostorové analýzy (s.s.)

Statistické prostorové analýzy zahrnují metody založené na stochastické (náhodné) povaze uspořádání a vztahů. Nabízejí poměrně široké spektrum aplikací. Někdy jsou přímo označovány jako prostorová statistika.

**Podle počtu současně zkoumaných charakteristik** můžeme používané statistické metody rozdělit na **monovariační** (jednorozměrné) a **multivariační** (analýzy vícerozměrných objektů/událostí). Používání alternativních názvů, uvedených v závorce, není příliš vhodné vzhledem k jejich snadné záměně za počet dimenzí u geometrické reprezentace objektů.

Monovariační statistické analýzy pracují současně pouze s jednou charakteristikou objektu, zatímco multivariační statistiky využívají více charakteristik současně. U multivariačních technik lze ještě specifikovat, zda studují více charakteristik současně jen vizuálně (a pak analýza a interpretace probíhá vizuálně a mentálně), nebo zda jde skutečně o multivariační metody s.s.

Dělení **podle povahy statistických technik** využívá analogie z tradičního rozdělení statických technik na popisné (typicky výpočet statistických charakteristik) a indukční (na základě studia výběru usuzujeme na vlastnosti celku, provádíme testování hypotéz). Podle tohoto principu se dělí techniky statistické prostorové analýzy na:

- a.) **popisné** (deskriptivní, centrografické techniky) – především kvantitativní měření charakteristiky polohy a charakteristiky rozptýlení
- b.) **inferenční** (analýzy textury) – ty určují, zda distribuce je nebo není náhodná, popisují vztahy mezi 2 a více veličinami

Pokud hovoříme o distribuci geoprvků, máme na mysli texturu (vzor), kterou vytváří geoprvky svým rozmístěním ve sledované části prostoru.

#### 1.4.4.1 Tři základní typy prostorové distribuce

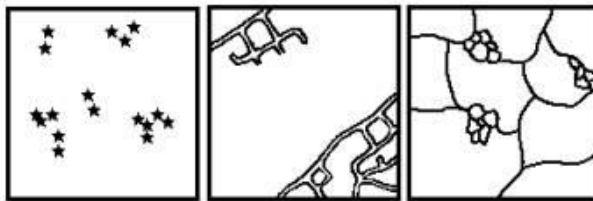
**Z hlediska statistického uspořádání geoprvků** rozlišujeme 3 základní typy prostorové distribuce:

- a.) **shluková** (clustered), případně skupinová
- b.) **pravidelná** (regular)
- c.) **náhodná** (random)

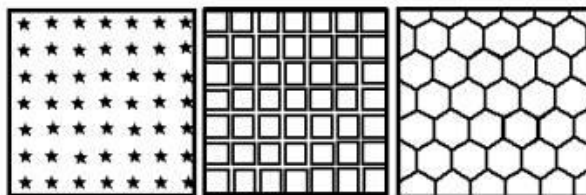
Někteří autoři ještě vymezují rovnoměrně rozmístěnou distribuci jevů, kterou lze zařadit mezi pravidelnou a náhodnou, Ivanička (1983) popisuje také aglomerizovanou a heterogenní distribuci. Většina autorů však vymezuje pouze uvedené 3 základní typy.

Uvedené dělení se využívá, pokud je naším cílem zjištění, zda (s jakou pravděpodobností) je prostorové rozložení geoprvků náhodné, resp. zda je statisticky průkazný jeho shlukový či pravidelný charakter. Skutečné distribuce mohou být testovány vůči těmto 3 ideálním typům, často však pouze vůči náhodné distribuci. Statisticky prokázaný výskyt shlukového či pravidelného vzorku může být

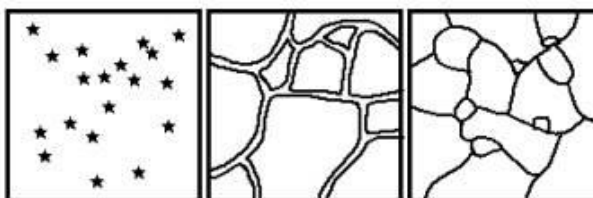
základem pro zjišťování příčin, které vedly k pozorovanému uspořádání. Např. prokázaný shlukový výskyt případů jisté choroby může (pokud byla data správně standardizována, a to i v časové ose) podpořit hypotézu infekčního charakteru (etiologie) choroby.



Obr.1-5 Shlukový typ distribuce pro 3 typy geoprvků: body, linie, areály



Obr.1-6 Pravidelný typ distribuce pro 3 typy geoprvků: body, linie, areály



Obr.1-7 Náhodný typ distribuce pro 3 typy geoprvků: body, linie, areály

#### 1.4.4.2 Podrobnější přehled základních technik pro provádění statistických prostorových analýz

Podrobnější přehled základních technik pro provádění statistických prostorových analýz zahrnuje:

1. Jednoduché deskriptivní analýzy, transformace dat a sumarizace
2. Metody nejbližších vzdáleností (*nearest neighbor*) a K-funkce
3. Kvadrantové, jádrové a Bayesovské vyhlazovací metody
4. Prostorová autokorelace a kovariační struktury
5. Geostatistické a prostorové ekonometrické modelování
6. Prostorové generalizované lineární modelování
7. Multivariační techniky
8. Prostorové interakční modely

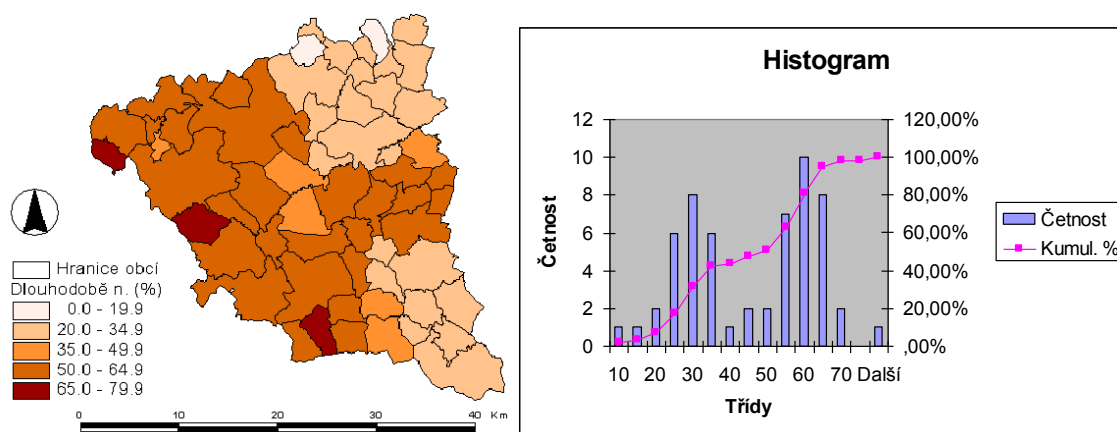
##### 1.1.1.1.1 Jednoduché deskriptivní analýzy, transformace dat a sumarizace

Tyto operace často nejsou samostatně uvedeny, avšak tvoří základní nástroje pro řadu dalších technik. Zahrnujeme do nich jednoduché grafické a numerické metody sumarizace dat a manipulace s daty (histogramy, vyrovnání histogramu jádrovým odhadem, ogiva, rankilové grafy, rozptylogramy, "vousaté" krabičky, projekce multivariačních dat do 2D zobrazení, výpočet základních statistických ukazatelů, zjišťování korelace, transformace dat). Popis těchto technik lze najít v základních statistických učebnicích nebo publikacích popisujících průzkumovou analýzu dat.

Zařazujeme zde i základní popisnou statistiku pro prostorové objekty a jevy (např. určení středu pro shluk bodů, tvorba elipsy disperze).

*Použití neprostorových statistických nástrojů k provedení prostorové analýzy je možné demonstrovat na 2 příkladech:*

- A. Výpočet statistiky dopravních nehod v různých částech města a její interpretace. Výběr dopravních nehod byl proveden pomocí dotazu na část města, do které přísluší ulice s registrovanou nehodou. Následně se vypočítají statistické ukazatele pro jednotlivé výběry dat a prověří se významnost těchto cílů. Cílem takové analýzy je zjistit a ověřit rozdíly v nehodovosti v různých částech města. Pokud bychom neprovedli rozdělení celého souboru měření (resp. nehod) na prostorovém základě, výsledkem by byla běžná statistická analýza nehod za celé město.
- B. Výpočet statistiky znečištění půdy Cd a porovnání znečištění u jednotlivých potenciálních znečišťovatelů. Provede se lokalizace vzorků půd se znečištěním Cd, vytvoří se obalové zóny kolem objektů jako jsou spalovny, silnice či továrny, následně se provede výběr bodů v polygonu a vypočítají se statistické ukazatele. V tomto případě je výpočet statistických ukazatelů součástí složitější analýzy, kdy jsou vyžadovány konstrukce nových grafických objektů a řešení geometrické úlohy „výběr bodů v polygonu“.



Obr.1-8 Podíl dlouhodobě nezam. z celkového počtu nez. (okr. Nový Jičín, 30.6.1999)

### 1.1.1.1.2 Metody nejbližších vzdáleností a K-funkce

Metody nejbližších vzdáleností (podrobný popis viz kapitola 3.2.2) a K funkce (podrobný popis viz kapitola 3.2.3) jsou určeny pro posouzení umístění událostí či objektů (především bodová reprezentace) a určení typu pozorované textury (náhodná, nenáhodná).

Metody nejbližších vzdáleností jsou založeny na grafickém srovnání pozorované distribuční funkce vzdáleností mezi objekty (nebo vzdáleností mezi náhodně umístěným bodem a pozorovaným objektem) s jinými pozorovanými distribučními funkcemi nebo s očekávanou (teoretickou) distribuční funkcí získanou z modelu vytvořeného z náhodných dat.

Funkce  $K(h)$  je definována jako očekávaný počet dalších objektů do vzdálenosti  $h$  od určitého objektu. Provede se grafické nebo statistické srovnání naměřené  $K$  funkce s  $K$  funkcí odvozenou z teoretických modelů a posuzuje se typ pozorovaného vzorku.

### 1.1.1.1.3 Kvadrantové, jádrové (kernelové) a Bayesovské vyhlazovací metody

Kvadrantové, jádrové a Bayesovské vyhlazovací metody jsou založeny na neparametrických technikách a slouží k transformaci dat z diskretní reprezentace do kontinuální, tedy k výpočtu hustoty událostí, k vyhlazení textury apod.

V prostorovém kontextu mohou být použity např. jako průzkumné metody pro identifikaci odlišných míst (tzv. *hot spots*, tedy více variabilních nebo naopak více homogenních míst), pro identifikaci vhodných modelů, pro analýzu shody modelů s naměřenými daty.

Kvadrantové metody představují nejjednodušší způsob transformace dat z diskretní reprezentace do kontinuální. V případě bodů se počítají výskyty bodů v jednotlivých buňkách překrývající mřížky.

Pokročilejší metody vycházejí z myšlenky **jádrového (kernelového) vyhlazování**. Vyhlazená hodnota v daném bodě je vypočtena jako vážený průměr z hodnot v okolních bodech, kde váhy jsou odvozeny z distribuce pravděpodobnosti se středem v příslušném bodě.

Jádrový odhad hustoty pracuje s lokalizačními daty a pak vyjadřuje prostorově vyhlazený odhad lokální intenzity výskytu objektů/událostí. Tuto lokální vyhlazenou intenzitu je možné chápat i jako povrch rizika výskytu těchto objektů/událostí.

Druhou možností je aplikace na atributová data a výpočet vyhlazeného odhadu sledovaných hodnot.

Podrobný popis kvadrantové metody pro testování náhodnosti viz kapitola 3.2.1, podrobný popis kvadrantové metody pro vyhlazování viz kapitola 3.4.1 a podrobný popis jádrového vyhlazení viz kapitola 3.4.2.

**Bayesovské vyhlazovací metody** jsou založeny na využití Bayesovy věty, často s aplikací Markovových řetězců.

#### **1.1.1.1.4 Prostorová autokorelace a kovariační struktury**

Při geostatistické analýze se považuje rozložení hodnot modelované veličiny (např. obsah Pb) za tzv. regionalizovanou proměnnou, která se vyjadřuje jako funkce souřadnic X,Y,Z. V každém bodě jistým způsobem vymezeného prostoru nabývá určité hodnoty. Je evidentní, že v případě přírodních objektů je hodnota v daném místě výsledkem řady procesů, z nichž některé mají výrazně náhodný charakter. Výsledkem je značná prostorová variabilita hodnot, jejich nespojitost a anizotropie.

Geostatistická analýza se snaží popsat chování této regionalizované proměnné. Jejím základním nástrojem jsou strukturální funkce.

Prostorové kovariační struktury se zjišťují v atributových datech a popisují závislost mezi rozptylem (resp. korelací) hodnot a vzdáleností měření. Uvádějí, zda a jak souvisí umístěním blízké hodnoty jedna s druhou.

Podrobný popis geostatistických metod viz kapitola 4.

#### **1.1.1.1.5 Geostatistické a prostorové ekonometrické modelování**

Geostatistické modelování je založeno především na provádění lokálních odhadů s využitím výsledků strukturální analýzy (aplikace interpolačních procedur).

Představují prostorové rozšíření standardních lineárních regresních modelů. Parametry jsou odhadovány pomocí funkce maximální věrohodnosti nebo zobecněnou metodou nejmenších čtverců.

Podrobný popis geostatistických metod viz kapitola 4.

#### **1.1.1.1.6 Prostorové generalizované lineární modelování**

Prostorové generalizované lineární modelování představuje zobecnění prostorových regresních modelů na případy, kdy modelovaná atributová data přísluší k výčtové doméně. Vycházejí z prostorového zobecnění myšlenek log-lineárního modelování kontingenčních tabulek a modelování Poissonových nebo binomických proměnných.

Podrobný popis prostorového generalizovaného lineárního modelování viz kapitola 6.7.3.

#### **1.1.1.1.7 Multivariační techniky**

Většina multivariačních technik není speciálně orientována na prostorově závislá data, ale i tak mohou být velmi užitečná jako nástroj pro redukci dat a pro identifikaci významné kombinace proměnných.

Metody se využívají např. v klasifikačních postupech při zpracování dat dálkového průzkumu Země.

Podrobný popis multivariačních postupů viz kapitola 6.6.

### **1.1.1.1.8 Prostorové interakční modely**

Prostorové interakční modely jsou založeny na modelování pozorovaného toku ze sady zdrojů do soustavy cílů. Pro soustavu zdrojů se definují požadavky, pro sadu cílů se popisuje atraktivnost. Měřítkem vzdálenosti může být i čas nebo náklady. Modely jsou běžně odvozovány ze zobecněných gravitačních modelů, zaměřených např. na minimalizaci procestovaných podmínek nebo na optimalizačních problémech - minimalizace procestované vzdálenosti, maximalizace entropie). Je zde možné zařadit i lokační a alokační úlohy a tvorbu servisních území.

Podrobný popis prostorových interakčních metod viz kapitola 5.

## **1.5 Vztah GIS a prostorové analýzy**

Geografické informační systémy dosáhly v poslední době značného rozšíření a staly se základním nástrojem pro správu a zpracování prostorových dat i prostředníkem pro poskytování a využívání prostorových informací. Jedním z prioritních cílů vytvářených geografických informačních systémů (GIS) je podpora uživatelů při rozhodování, ke kterému využívají zprostředkované prostorových informací ve vhodné formě, často jako výsledek prostorových analýz ve smyslu společné analýzy geometrické a tématické (atributové) složky dat. Již Aronoff (1989) vyzdvihuje na GIS jako jejich nejcennější rys právě provádění prostorových analýz. Kořeny prostorových analýz leží v různých disciplínách, jejich další rozvoj a především úspěšná implementace je však již dnes pevně spojována s GIS.

V těchto systémech bývá k dispozici dobrá sada nástrojů pro realizaci řady metod používaných v prostorových analýzách (výběry na základě dotazování, logické operace prováděné na základě atributů geoprvků, překryvné operace, měření vzdálenosti a spojitosti, charakteristiky okolí).

Přitom je zřejmé, že některé metody našly velmi rychlé uplatnění v GIS (typicky metody mapové algebry), jiné se však obtížně uplatňují. Důvodem může být jejich nesnadná algoritmizace či obtížnost použití (spojená např. s požadavkem odborné erudice při provádění příslušné analýzy). Příkladem mohou být některé statistické metody, regresní modely, multivariační postupy či geostatistické modelování.

Někteří odborníci doporučují místo implementace klasických technik v geografických informačních systémech raději vývoj nových postupů, které by plně využily koncepce a metod GIS. Přebírané metody by měly být dostatečně inteligentní, aby využily více realistické reprezentace prostoru v GIS, dále by měly být vhodné pro co nejširší použití (různé aplikace), snadno algoritmizovatelné a výpočetně jednoduché.

Uvedená situace je vcelku logická a odpovídá běžnému vztahu mezi vývojem metod, aplikací specifických postupů vhodných pro určité řešení, chápaných jako rozvoj základní a aplikované vědy, a na druhé straně nutné zcela jiné zaměření GIS chápaných jako informační systémy (IS). Připomeňme, že hlavním dogmatem informačních systémů je jednoduchost a standardizace. Proto i GIS v roli IS budou využívat především jednodušší, osvědčené a snadno aplikovatelné postupy. Ve formě flexibilního prostředí také geografické informační systémy poskytují prostor pro provádění náročných analýz a modelování. Výsledky jsou často ukládány v GIS, avšak využívány v kombinaci s jinými informačními zdroji a prezentovány pro uživatele IS.

Významný pokrok také přinesly GIS v oblasti vizualizace jak získaných dat tak i výsledků zpracování. Nástroje počítačové grafiky značně usnadnily některé nejběžnější typy výstupů, především tvorbu statistických map. Rovněž se začínají rozšiřovat i netradiční způsoby vizualizace dat.

## 2 Typy dat a vztah k prostorovým analýzám

Správné využití prostorových analýz není možné bez dostatečného poznání té části světa (reality), kterou popisujeme, zkoumáme, analyzujeme nebo modelujeme. Realita je popisována na základě určitého konceptuálního modelu pomocí dat, která představují formalizovaný záznam reality (přesný význam definuje ČSN 369001). Proto je pro nás zásadní důkladné seznámení s daty, jejich správný výběr a výběr i vhodných metod pro jejich zpracování.

Prostorová data představují popis geoprůvku, který zahrnuje geometrickou složku, atributovou (popisnou) složku, časovou, funkční a vztahovou složku.

### 2.1 Grafické atributy

Významná část prostorových dat vyjadřuje svou geometrickou složku popisu (tedy svou lokalizaci a prostorovou reprezentaci) pomocí grafického popisu, grafickým vyjádřením lokalizace.

**Z hlediska typu vazby** mezi grafickou a atributovou složkou dat (a způsobu zpracování dat) můžeme data rozdělit na:

a) **lokalizační data** (*locational data, event data, point process*). Tato data obsahují jen lokalizaci objektů či událostí. Předpokládáme, že sledujeme výskyt pouze jednoho typu objektu nebo události, o kterých kromě jejich lokalizace nic dalšího nevíme. Zajímá nás tedy, kde se to stalo (kde se vyskytly události), ne co se stalo (nezajímá nás druh či popis události). Příkladem může být lokalizace případů choroby, lokalizace dopravních nehod apod.

Pokud při zpracování využíváme údaje o rozdílných typech objektů/událostí, musíme použít multivariační techniky, resp. analýzu vícenásobných typů událostí.

Lokalizační data se vztahují k jednomu místu a nepotřebují popisné atributy. Prakticky se jedná o pouze grafická data. Vzhledem k charakteristice dat se uplatňují i postupy známé z počítačové grafiky. Kardinalita a parcialita vazby mezi grafickou a atributovou složkou je zpravidla 1:0.

b) **atributová data** (*attribute data*). Taková data obsahují hodnoty (atributy) spojené s nějakými objekty či událostmi. Lokalizace bývá vyjádřena body, buňkami pravidelného gridu nebo polygony. Jako příklad můžeme uvést půdní vzorky lokalizované bodem (protože mají z geografického hlediska zanedbatelný objem) a popsány např. výsledky geomechanických zkoušek či chemických analýz. Nebo jde o data z dálkového průzkumu Země uložená v pravidelné mřížce, která (po radiometrických a geometrických korekcích) odpovídají velikosti odrazivosti a emisivity příslušné části zemského povrchu. Poslední příklad se vztahuje k polygonům (resp. k areálům) - úmrtnost v jednotlivých okresech je reprezentována polygony.

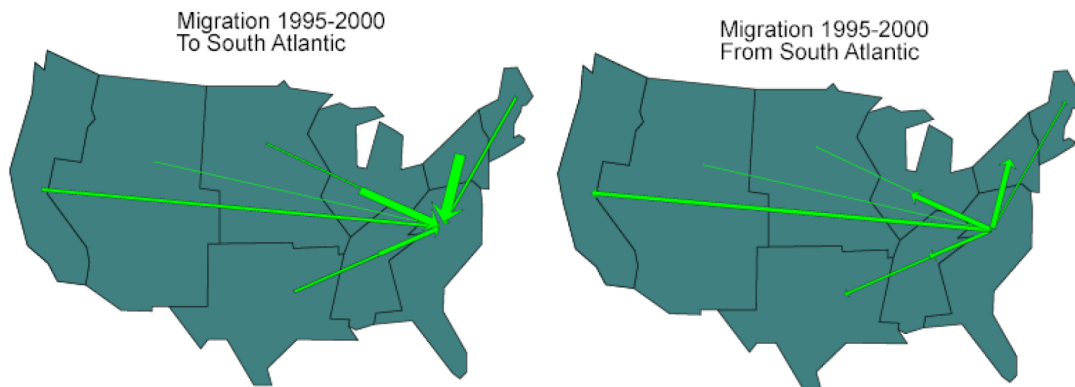
Opět se tedy jedná o data vztažená k jednomu místu, ke kterému se váží jisté atributy. Z hlediska kardinality vazby mezi grafickou a atributovou složkou jde u monovariačních dat o vztah 1:1 (případně N:1), u multivariačních pak běžně o M:N.

Multivariační postupy je nutné použít v případech, kdy sledujeme současně více atributů u každého objektu. Jedním z nich může být i čas.

V řadě případů je možné využít konceptu regionalizované proměnné, tj. považovat data za geostatistická.

c) **interakční data** (*interaction*). Typicky jsou to kvantitativní měření spojená s liniemi nebo páry míst (často 2 body, mohou to být ale i sady polygonů). Příkladem může být tok spotřebního zboží z míst uskladnění do obchodů nebo migrace obyvatel mezi územními celky (obr. 2-1).

Typicky jde o data vztažená ke 2 místům, proto bychom mohli hovořit o charakteristické vazbě mezi grafickou a atributovou složkou dat 2:1 či obecněji M:N.



Obr.2-1 Velikosti migračních toků v letech 1995 až 2000 směřujících do a z amerického jihovýchodu (<http://www.csiss.org/clearinghouse/FlowMapper/FlowTutorial.pdf>)

## 2.2 Popisné atributy

K tématickému popisu objektů používáme atributy, přesněji popisné atributy. Atributy nám slouží k záznamu i vyjádření vlastností objektů.

Postupem času (v průběhu rozvoje informatiky) se ukázalo, že je účelné přes velmi pestrou škálu informací, se kterými pracujeme, vymezit základní způsoby kódování (formalizace) dat, což má bezprostřední vliv na způsob uložení dat. Vymezení těchto základních tříd kódování vede k definici domén.

Doménu je možné charakterizovat jako potenční množinu dat, ze kterých je vybírána hodnota atributu. Základní typy domén z hlediska způsobu uložení dat v databázích jsou běžně vymezeny ve specifikaci jazyka SQL, který považujeme za základní nástroj standardizace v databázových systémech vedle nastupujícího jazyka XML.

Toto vymezení domén je založeno na **způsobu uložení dat**, i když určitým způsobem orientuje i dominantní způsob zpracování dat.

Jinou možností je vymezení typů domén podle **způsobu zpracování dat**. Vycházíme z toho, že z hlediska uložení dat pro nás může být výhodné použít např. doménu přirozených čísel, avšak interpretace těchto dat (jejich sémantika) může být velmi odlišná a může zásadně determinovat smysluplné možnosti zpracování (manipulace) dat. S číselnou doménou (zvláště racionálních čísel) je možné z matematického hlediska provádět plnou škálu matematických či statistických operací (aritmetické, goniometrické, statistické operace), při vědomí významu dat však může být množina operací omezena.

*Příklad - typ pokryvu krajiny může být kódován pomocí číselného kódu. 1 pak může označovat např. lesní porost, 2 zástavbu a 3 např. vodní plochy. Předpokládejme, že jsou data zaznamenána pomocí rastrového datového modelu a že chceme provést generalizaci pokryvu, tedy jisté vyhlazení (zanedbání malých odlišných oblastí). Generalizaci můžeme provádět na základě posouzení výskytu hodnot v širší oblasti (zpravidla posuvné okno, konvoluce). Přestože z hlediska uložení dat jde o číselné hodnoty, nelze použít pro charakteristiku oblasti např. vážený průměr (lesní a vodní plochy ukazují v „průměru“ na zástavbu!). Jedinou možností je zde použít nejčtenější hodnoty. Tedy převažuje-li v regionu les, budeme celý region klasifikovat jako les.*

Velmi jednoduché dělení nabízí rozdělení dat na **kvalitativní** a **kvantitativní**, vhodnější je ale následující přesnější vymezení čtyř tříd.

Z hlediska **zpracování dat** tedy rozlišujeme následující typy domén:



- **poměrová** (*ratio*) - odpovídá hodnotě na kalibrované, lineární škále ve vztahu k pevnému bodu. Takové hodnoty lze libovolně zpracovávat matematickými funkcemi. Patří sem věk, četnost, vzdálenost, cena apod. Z hlediska způsobu uložení dat jde pouze o číselnou doménu.
- **intervalová** (*interval*) - používá hodnoty s pozicí na kalibrované, lineární škále, která však nemá vztah k fixnímu bodu. Takové hodnoty lze porovnávat, ale ne např. násobit nebo poměřovat. Typickým příkladem je teplota: -5 stupňů C je o 10 stupňů méně než +5 stupňů C, ale nemá smysl vyjádřit jejich poměr (podobně i stupnice F versus C). Často charakterizují relativní pozici v prostoru, čase nebo velikost - zeměpisná šířka, nadmořská výška (problém různé srovnávací hladiny), délka, směr kompasu, čas během dne, normalizované skóre apod. Z hlediska způsobu uložení dat jde pouze o číselnou doménu.
- **pořadová** (*ordinal*) - vyjadřuje hodnotu na nekalibrované, lineární škále. Lze získat pouze kvalitativní rozdíl, ne kvantitativní. Hodnoty se tedy dají rozlišovat pořadím, ale ne velikostí. Např. 1., 2. a 3. místo na závodech nebo hodnoty „tmavě šedá“, „šedá“, „světle šedá“, „šedobílá“. Pavlík, Kühnl (1981) označují jevy popisované pomocí pořadové škály jako topologické (na rozdíl od metrických pro poměrovou škálu).
- **nominální, výčtová** (*nominal*) - hodnoty nemají vztah k lineární škále. Nelze je řadit, lze je pouze porovnat na rovnost či nerovnost. Reprezentují kvalitativní hodnoty. Typickým příkladem může být doména: {jehličnatý, smíšený, listnatý}. Nejsou kvantifikovatelné, ale pouze klasifikovatelné.

Shrneme-li, přestože použijeme k uložení dat doménu např. přirozených čísel, je aplikace matematických funkcí smysluplná pouze pro poměrová data. U intervalových dat je možné použít i rozdíl hodnot. Funkce maximum, minimum lze aplikovat i na pořadové hodnoty. Na výčtové hodnoty lze uplatnit pouze takové funkce jako testování shody, zjišťování variability, modus apod.

Někteří autoři (Pavlík, Kühnl 1981) vyčleňují z nominálních (výčtových) dat ještě binární data (ve smyslu dvouhodnotové logiky). Je zjevné, že pro takovou doménu by bylo možné uplatnit pouze testování shody a zjišťování variability už je smysluplné pouze při uplatnění prostorového charakteru úlohy (např. iterační testy náhodnosti sekvence jako je 11000101...).

Kraak, Ormelling (1996) uvádí vazbu těchto domén ke kartografickým výrazovým prostředkům (velikost, barva, tvar, textura, hodnota) a k požadovanému efektu (odlišení, uspořádání, vzdálenost, proporce). Popisuje vztah k diskrétním a kontinuálním objektům a jakým způsobem se mapují. DeMers (1997) uvádí tabulku možností využití 4 základních typů domén ve vztahu k základním grafickým formám reprezentace objektů.

Přiřazení zkoumaných dat k jednomu z uvedených typů domén však představuje pouze 1. krok při jejich poznání nezbytném pro volbu správného způsobu zpracování dat (zpracování s.l. - od vstupu, uložení, analýzy, syntézy, modelování až po prezentaci)

Někteří autoři ještě vymezují různé typy indexů - density (které mají vztah k ploše), nepravé poměry, průměry a potenciály.

Podobnou situaci můžeme doložit u číselných hodnot, kdy má smysl z hlediska zpracování odlišovat veličiny **absolutní** a **relativní**.

**Absolutní** veličiny zahrnují např. počet nebo hmotnost, často primární měřená či zjištěná data.

**Relativní** veličiny jsou poměrná čísla (např. míra nezaměstnanosti). U relativních veličin nemá smysl provádět součet. Relativní veličiny můžeme dále rozdělit na:

- extenzitní (ukazatele struktury, představují podíl z celku),
- intenzitní (ukazatele intenzity), dále dělené na:
  - míry,
  - kvocienty,
- indexy, které vyjadřují časový nebo geografický vývoj.

V anglosaské literatuře se ve smyslu absolutní - relativní používá označení extenzivní - intenzivní, je třeba však upozornit, že jiní autoři používají toto označení v jiném smyslu. Extenzivní data podle

nich představují měřitelné vlastnosti a intenzivní data neměřitelné vlastnosti, často složené pojmy jako „hornatost“.

### 3 Body (lokalizační data)

Body představují nejčastější způsob reprezentace geografických fenoménů. Zpravidla jde ve skutečnosti o 3D objekty. Body jsou umístovány v těžišti objektů, které představuje gravitační centrum pro analýzy. Těžiště se konstruuje např. v místě křížení nejdelsí a nejkratší osy objektu (zpravidla plochy). U konvexních objektů se tak může těžiště dostat i mimo vlastní objekt. Pokud je to problém, řešení může být založeno na středu vepsané kružnice, který vždy bude ležet uvnitř objektu.

Při analýze bodového vzorku (rozmístění, distribuce bodů) popisujeme příbuznost vzorku k pravidelné, shlukové nebo náhodné prostorové distribuci.

Velmi jednoduchým, ale také značně hrubým kritériem popisu bodového vzorku je **hustota bodů v ploše** (počet/plocha =  $n/R$ ), které nepříjemně závisí na definici hranic oblasti (resp. na velikosti oblasti). Z tohoto důvodu jsou lepší charakteristiky založené na vzdálenosti mezi body nebo na relativních vzdálenostech jako je např.  $d_i/d_{max}$ .

Při výpočtech v relativně malých oblastech používáme euklidovskou geometrii, protože se v nich neprojeví zakřivení Země.

#### 3.1 Popisná statistika pro body

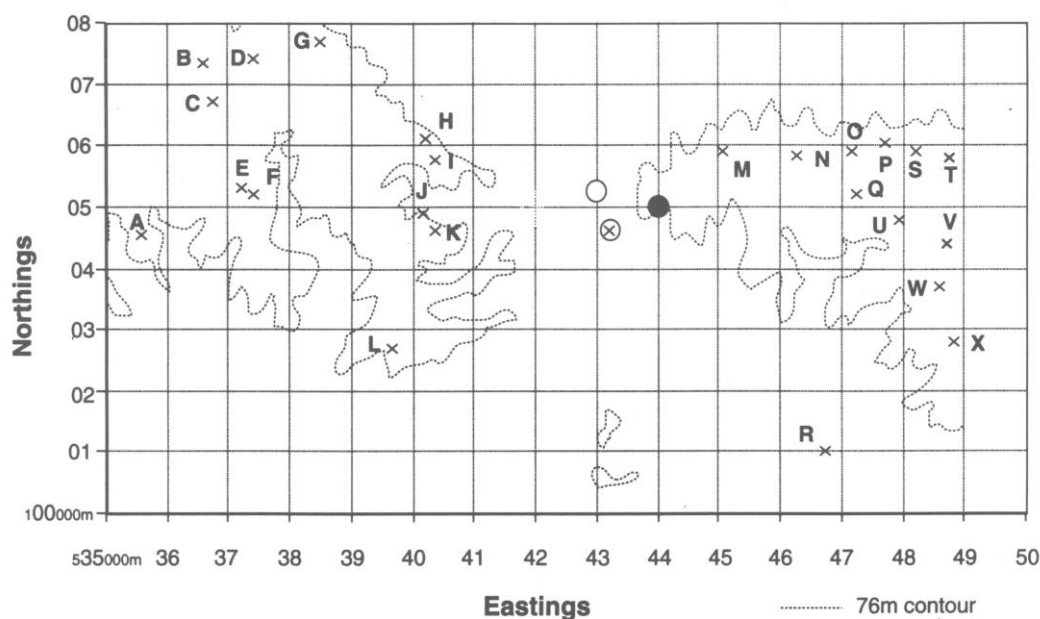
Popisné metody se zabývají určením charakteristiky polohy (např. určení geografického středu, mediánu) a charakteristiky rozptýlení (např. směrodatná vzdálenost nebo směrodatná elipsa). Popisují distribuci bodů pomocí základních statistických charakteristik. Používají se ke srovnání více bodových vzorků nebo ke sledování jejich vývoje v čase.

##### 3.1.1 Charakteristiky střední polohy

- **průměrný střed** (*mean centre*)  
Střed leží na průměru souřadnic X a Y.

$$\bar{X} = \frac{\sum x}{n} \quad \bar{Y} = \frac{\sum y}{n}$$

Příklad:



Tab. 3-1 Rozmístění mohyl ve východním Sussexu v Anglii a postup výpočtu středu (podle Walford 1995)

Místo	X	Y	Počet mohyl (W)	X*W	Y*W
A	3560	455	1	3560	455
B	3660	735	1	3660	735
C	3675	670	1	3675	670
D	3740	740	1	3740	740
E	3720	530	4	14880	2120
F	3740	520	2	7480	1040
G	3850	770	2	7700	1540
H	4020	610	1	4020	610
I	4035	575	1	4035	575
J	4015	490	1	4015	490
K	4035	460	1	4035	460
L	3965	270	1	3965	270
M	4505	590	1	4505	590
N	4625	585	1	4625	585
O	4715	590	2	9430	1180
P	4770	605	5	23850	3025
Q	4725	520	3	14175	1560
R	4675	100	4	18700	400
S	4820	590	2	9640	1180
T	4875	580	3	14625	1740
U	4795	480	1	4795	480
V	4870	440	2	9740	880
W	4860	370	1	4860	370
Y	4885	280	1	4885	280
	103135	12555	43	188595	21975

$$x = 103135/24 = 4297$$

$$y = 12555/24=523$$

Střed je vyznačen prázdným kroužkem.

- **vážený průměrný střed (weighted mean centre)**

Používá se v případě výskytu více událostí/objektů na stejném místě. Pak má každý bod váhu přímo úměrnou počtu událostí/objektů na tomto místě.

$$\bar{X} = \frac{\sum X_i * W_i}{\sum W_i} \quad \bar{Y} = \frac{\sum Y_i * W_i}{\sum W_i}$$

Příklad:

$$x = 188595/43 = 4386$$

$$y = 21975/43=511$$

- **agregovaný průměrný střed**

Alternativa váženého středu, kdy se nepoužívají původní souřadnice X,Y ale jen souřadnice čtverců s agregovaným počtem bodů uvnitř čtverce.

$$\bar{X} = \frac{\sum X_i * F_i}{n} \quad \bar{Y} = \frac{\sum Y_i * F_i}{n}$$

n je celkový počet objektů (mohyl)  
 F<sub>i</sub> frekvence objektů (mohyl) ve čtvercové buňce  
 X<sub>i</sub> a Y<sub>i</sub> souřadnice středů čtverc. buněk  
 i je od 1 do m, kde m je počet neprázdných čtverců

Průměrný střed má stejné problémy jako aritmetický průměr – je to především citlivost na extrémní hodnoty. Existují-li shluky, výsledek nereprezentuje dobře střed.

Příklad:

Tab. 3-2 Souřadnice středů čtverců pro lokalizaci mohyl ve východním Sussexu a postup výpočtu agregovaného středu

X	Y	Počet mohyl (W)	X*W	Y*W
3550	450	1	3550	450
3650	750	1	3650	750
3650	650	1	3650	650
3750	750	1	3750	750
3750	550	6	22500	3300
3850	750	2	7700	1500
3950	250	1	3950	250
4050	650	1	4050	650
4050	550	1	4050	550
4050	450	2	8100	900
4550	550	1	4550	550
4650	550	1	4650	550
4650	150	4	18600	600
4750	550	5	23750	2750
4750	650	5	23750	3250
4850	550	5	24250	2750
4850	450	3	14550	1350
4850	350	1	4850	350
4850	250	1	4850	250
			188750	22150

$$x = 188750/43 = 4390$$

$$y = 22150/43=515$$

Střed je vyznačen plným kroužkem na obr. 3-1.

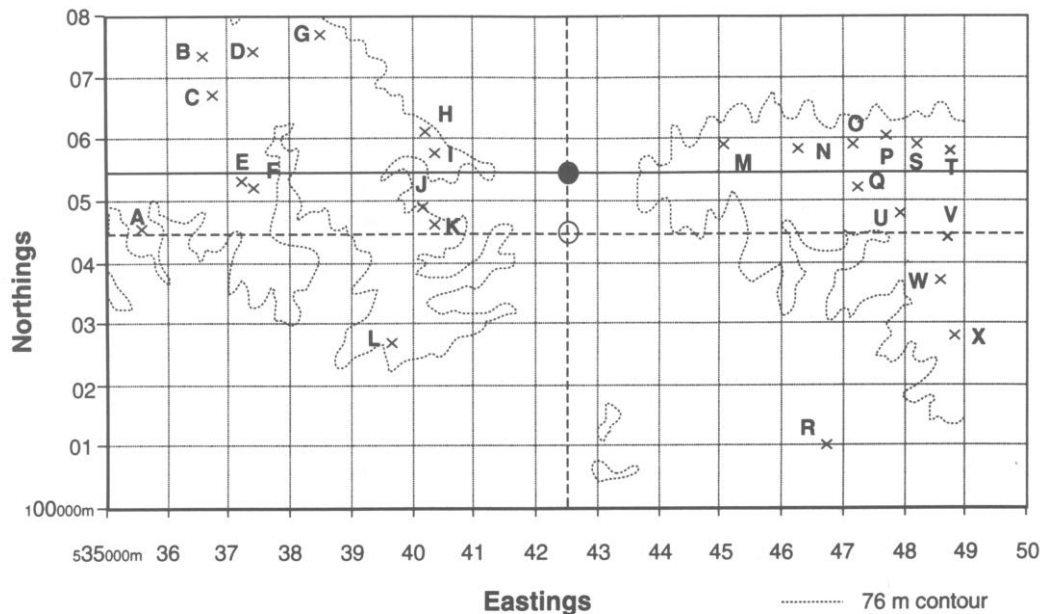
V tomto případě není agregovaný průměrný střed vhodný, protože jen málo buněk obsahuje body a všeobecně je málo bodů. Jen 19 ze 116 čtverců obsahuje body.

- **mediánový střed**

Jedná se o analogii mediánu, jsou ale rozdílné názory na způsob definice mediánového středu.

- a) najdeme medián na ose X a Y a vedeme z nich linie kolmé na směr osy. Takto definovaný „medián ze souřadnic“ ale nemusí odpovídat mediánu souboru bodů, protože distribuce nemusí být mezi kvadranty vyrovnaná.

Příklad:



Obr. 3-2 Rozmístění mohyla a určení mediánového středu (Walford 1995)

X=4235  
Y=553

Střed je vyznačen plným kroužkem.

- b) Rozdělíme počet bodů do 4 stejně početných skupin pomocí průměrů horizontální a vertikální linie. Problémem je, že je to nejednoznačná úloha, protože může existovat více řešení (více proložených linií).

Příklad:

X=4235  
Y=450

Střed je vyznačen prázdným kroužkem.

Na obr.3-2 je vykreslena izolinie 76 m. Většina bodů leží ve výše položených místech ve sledovaném území a tvoří 2 výrazné shluky. Vypočtené průměry však nepatří ani k jednomu bloku a dokonce spíše leží v údolí (pod uvedenou vrstevnicí).

- c) v severní Americe se prosazuje střed s nejvyšší dostupností (minim. vzdálenosti do všech bodů). Někdy se používá pro něj označení MAT (minimum aggregated travel). Tento mediánový střed o souřadnicích X a Y (obě s vlnkou) musí splňovat následující podmínku:

$$\min \sum_{i=1}^n \sqrt{(x_i - \bar{X})^2 + (y_i - \bar{Y})^2}$$

Praktický výpočet je založen na optimalizaci popsané v Smith et al. 2011 (s.119-120).

Vzdálenosti jsou zde uvažovány jako euklidovské, mohly by být ale použity i jiné metriky.

- **Geometrický střed**

Tato charakteristika střední polohy je méně citlivá na odlehlé hodnoty. Implementace je k dispozici v CrimeStat.

$$\text{Geometric Mean of X} = \text{GM}(X) = \sqrt[N]{\prod_{i=1}^N (X_i)^{1/N}}$$

$$\text{Geometric Mean of Y} = \text{GM}(Y) = \sqrt[N]{\prod_{i=1}^N (Y_i)^{1/N}}$$

- **Harmonický střed**

Tato charakteristika střední polohy je méně citlivá na odlehlé hodnoty. Pozor na hodnoty souřadnic blízké nule. Implementace je k dispozici v CrimeStat.

$$\text{Harmonic mean of X} = \text{HM}(X) = \frac{N}{\sum (1/X_i)}$$

$$\text{Harmonic mean of Y} = \text{HM}(Y) = \frac{N}{\sum (1/Y_i)}$$

- **Směrový střed**

Vypočítá se na základě střední úhlu a střední vzdálenosti od zvoleného počátku (např. bod o souřadnicích min. X a min.X, neboli levý dolní roh MBR). Jednotlivé body se spojí s tímto počátkem, a určí se délka spojnice d a úhel Theta. Střední směrová vzdálenost je jednoduše aritmetický průměr z délek d. Střední úhel (mean angle) se vypočte jako:

$$\text{Mean angle} = \bar{\theta} = \text{Abs} \left\{ \text{Arctan} \left[ \frac{\sum d_i \sin \theta_i}{\sum d_i \cos \theta_i} \right] \right\}$$

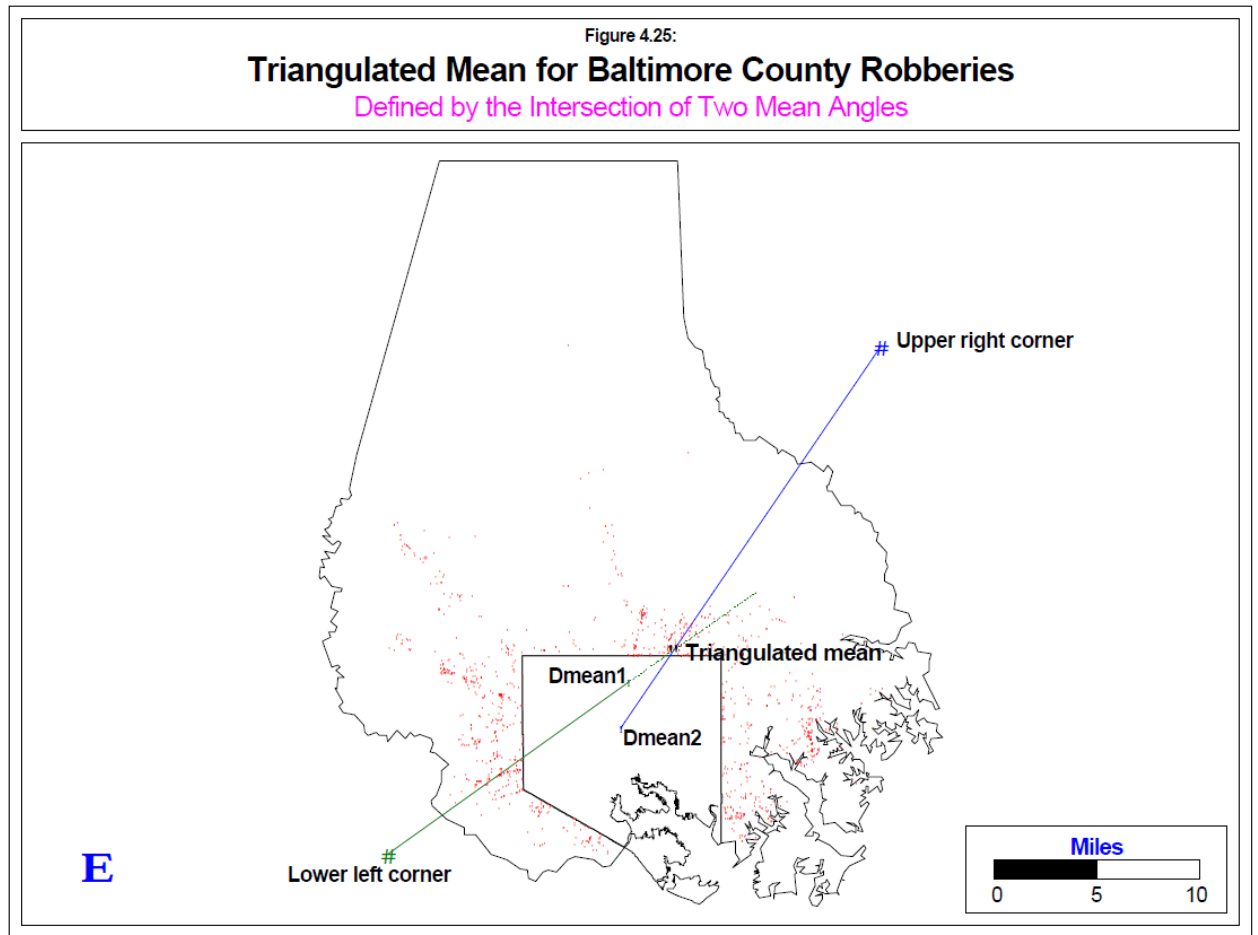
Což není nic jiného než – v čitateli je suma  $\Delta Y_i$ , ve jmenovateli suma  $\Delta X_i$ . Jejich poměr udává průměrný tangens, arctan tedy příslušný úhel.

Směrový střed se získá vynesemím úhlu z počátku a na něj příslušná střední vzdálenost.

- **Trojúhelníkový střed**

Vynesou se 2 polopřímky z počátku v levém dolním a pravém horním rohu MBR do příslušných směrových středů. Jejich průsečík udává trojúhelníkový střed.

Výpočet v CrimeStat, manuál kap. 4.



Obr. 3-3 Konstrukce trojúhelníkového středu (manuál k programu Crimestat, kap. 4)

### 3.1.2 Charakteristiky rozptýlení

- **Směrodatná vzdálenost** (*standard distance*)  
 Odpovídá měření rozptýlu kolem průměrného středu.

$$S_d = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n} + \frac{\sum (y_i - \bar{y})^2}{n}}$$

Příklad:

Tab. 3-3 Rozmístění mohyl ve východním Sussexu v Anglii a postup výpočtu směrodatné vzdálenosti (podle Walford 1995)

Místo	X	Y	Počet mohyl (W)	(X-x)*(X-x)	(Y-y)*(Y-y)
A	3560	455	1	543169	4624
B	3660	735	1	405769	44944
C	3675	670	1	386884	21609
D	3740	740	1	310249	47089
E	3720	530	4	332929	49
F	3740	520	2	310249	9
G	3850	770	2	199809	61009
H	4020	610	1	76729	7569
I	4035	575	1	68644	2704
J	4015	490	1	79524	1089



K	4035	460	1	68644	3969
L	3965	270	1	110224	64009
M	4505	590	1	43264	4489
N	4625	585	1	107584	3844
O	4715	590	2	174724	4489
P	4770	605	5	223729	6724
Q	4725	520	3	183184	9
R	4675	100	4	142884	178929
S	4820	590	2	273529	4489
T	4875	580	3	334084	3249
U	4795	480	1	248004	1849
V	4870	440	2	328329	6889
W	4860	370	1	316969	23409
X	4885	280	1	345744	59049
				5614849	556091

Poznámka: hodnoty jsou uvedeny bez desetinných míst; malé x a y znamenají hodnoty souřadnic průměrného středu.

$$S_d = \sqrt{\frac{5614849}{24} + \frac{556091}{24}} = 5.071km$$

- **koeficient relativního rozptýlení** (*coefficient of relative dispersion*)

Směrodatná vzdálenost / poloměr kruhu se stejnou plochou jakou má studovaná oblast

$$CRD = 100 * \frac{S_d}{A_k} = 100 * \frac{S_d}{\sqrt{\frac{R}{\pi}}} = 100 * S_d * \sqrt{\frac{\pi}{R}}$$

kde R je plocha oblasti

Je-li oblast různě velká (ohraničená), vznikají zavádějící hodnoty. Proto se někdy používá k získání relativního míry při studiu variability obyvatelstva poloměr země nebo státu místo poloměru kruhu se stejnou plochou jakou má studovaná oblast.

- **elipsa standardizované odchylky** (*standard deviation ellipse*)

Slouží k vyjádření směrové odchylky tam, kde je výrazná anizotropie v distribuci bodů (protažení apod.).

Je potřebné určit 3 prvky: úhel rotace, délku delší poloosy a délku kratší poloosy

Používá se rotace původních dat do nového souřadnicového systému.

Postup:

- 1) vypočtou se souřadnice průměrného středu X, Y.
- 2) pro každý bod se vypočtou transformované souřadnice

$$x'_i = x_i - \bar{X} \quad y'_i = y_i - \bar{Y}$$

- 3) vypočte se úhel rotace

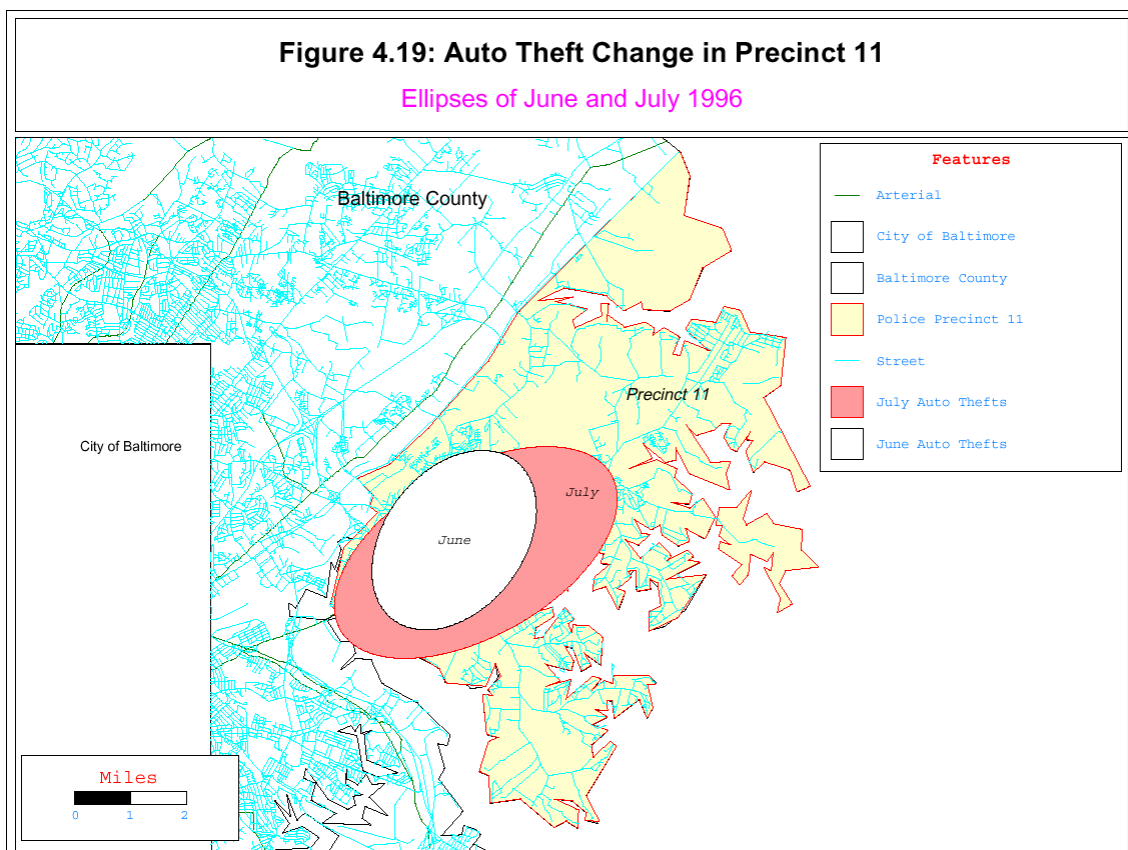
$$\tan \theta = \frac{\left( \sum_{i=1}^n x_i'^2 - \sum_{i=1}^n y_i'^2 \right) + \sqrt{\left( \sum_{i=1}^n x_i'^2 - \sum_{i=1}^n y_i'^2 \right)^2 + 4 \left( \sum_{i=1}^n x_i' \sum_{i=1}^n y_i' \right)^2}}{2 \sum_{i=1}^n x_i' \sum_{i=1}^n y_i'}$$

4) vypočtou se délky poloos elipsy

$$l_a = \sqrt{\frac{\sum_{i=1}^n (x_i' \cos \theta - y_i' \sin \theta)^2}{n}} \quad l_b = \sqrt{\frac{\sum_{i=1}^n (x_i' \sin \theta + y_i' \cos \theta)^2}{n}}$$

V CrimeStat a Smith et al. (2011) je to mírně jinak – s vysvětlením o správném výpočtu směrodatných odchylek ve směru obou os:

$$SD_x = \sqrt{\frac{2 \sum ((x_i - \bar{x}) \cos \theta - (y_i - \bar{y}) \sin \theta)^2}{n-2}} \quad SD_y = \sqrt{\frac{2 \sum ((x_i - \bar{x}) \sin \theta + (y_i - \bar{y}) \cos \theta)^2}{n-2}}$$



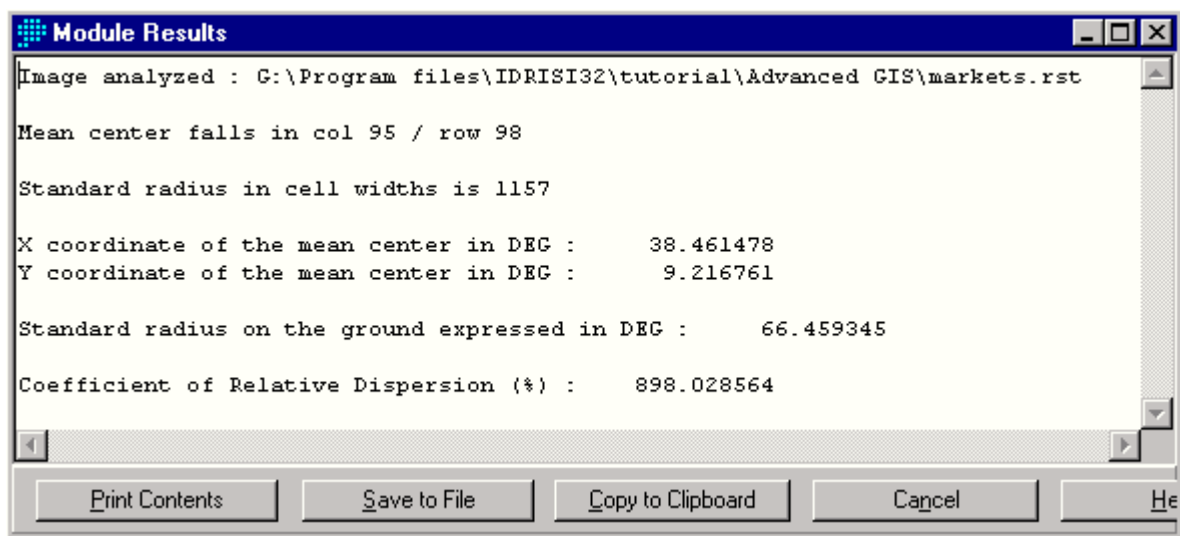
Obr. 3-4 Změny ve výskytu krádeží aut vyjádřené pomocí elips standardizované odchylky (© Levine and Associates)

### 3.1.3 Dostupné nástroje v programech

IDRISI má funkci CENTER, která vypočte vážený průměrný střed sady bodů, směrodatnou vzdálenost (*standard radius*) a koeficient relativního rozptýlení CRD. Funkce lze výhodně použít i při analýze časoprostorových změn (opakovaný výpočet pro různá období).

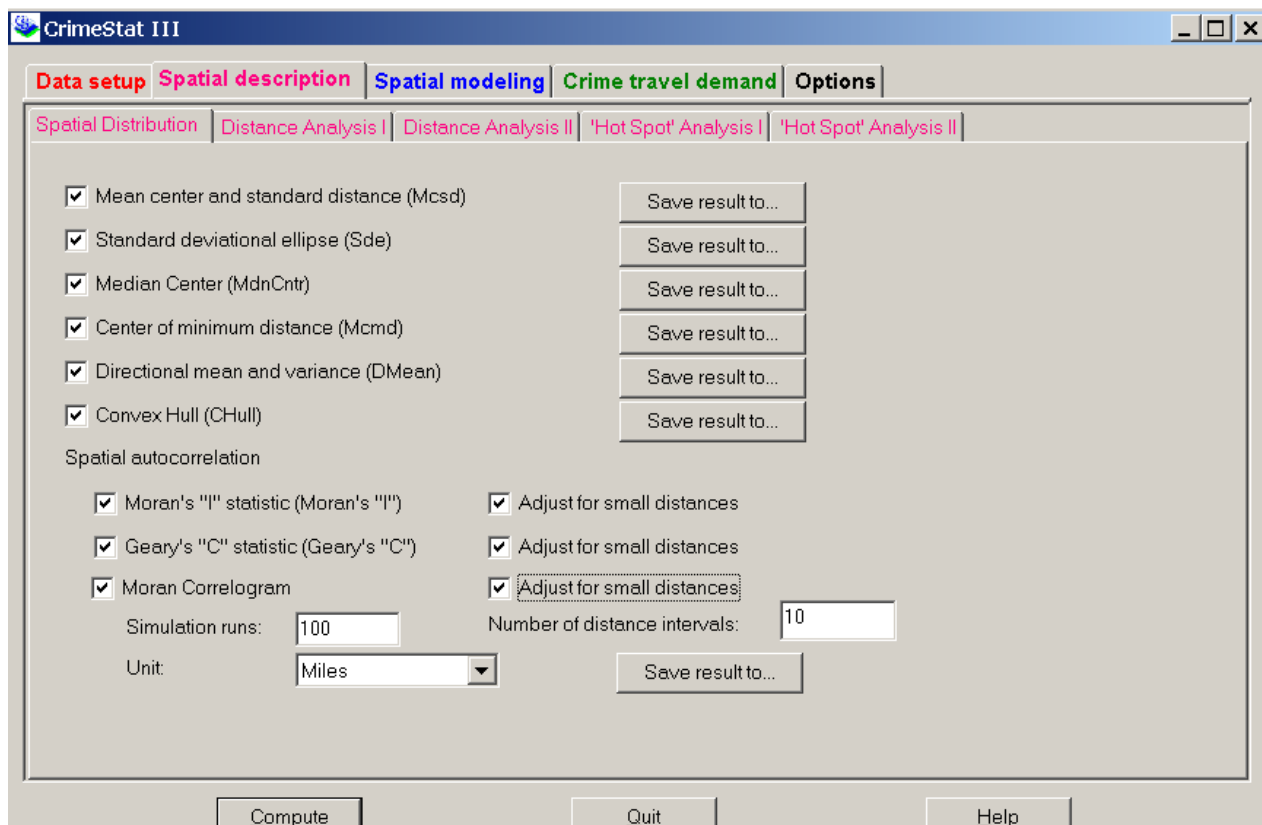


Obr. 3-5 Vstupní data pro výpočet váženého průměrného středu sady bodů, reprezentujících obchody



Obr. 3-6 Výsledek analýzy provedené v prostředí programového produktu IDRISI

Zřejmě největší sadu nástrojů nabízí CrimeStat.



Obr. 3-7 Nabídka funkcí pro popisnou statistiku sady bodů v programu CrimeStat III (manuál k CrimeStat, kap. 4)

## 3.2 Inferenční statistické testy pro body

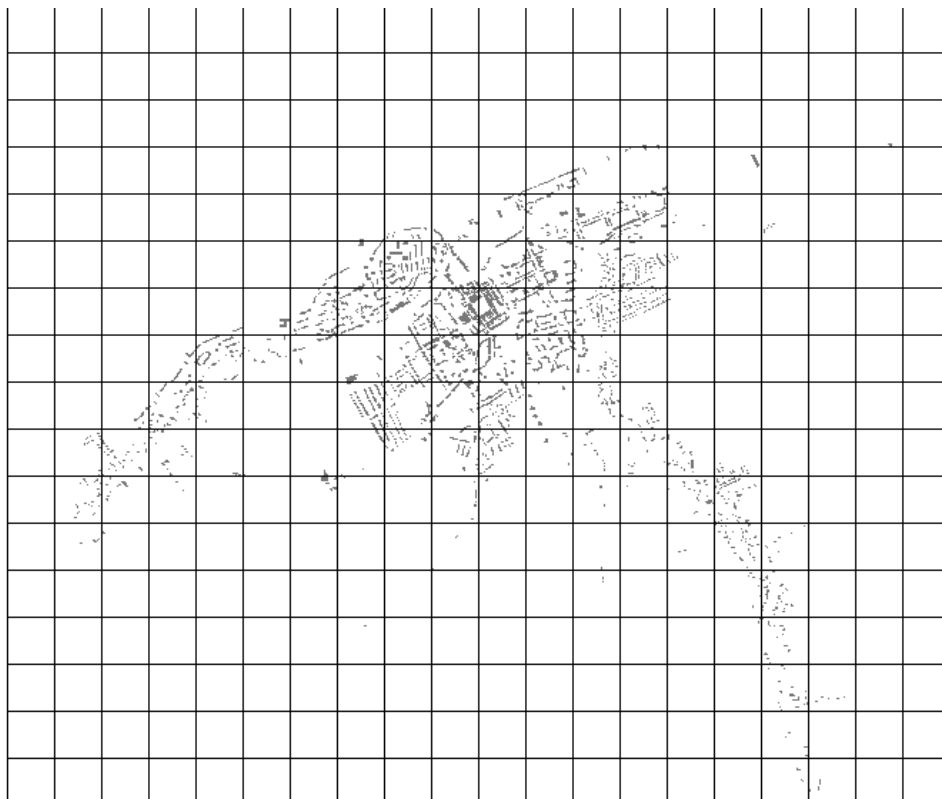
Určují pravděpodobnost, že určitá distribuce bodů vznikla náhodným procesem. Testy oceňují, zda rozdíl mezi očekávanou náhodnou distribucí a pozorovanou distribucí je významný.

K posouzení náhodnosti pozorovaného bodového vzorku lze použít řady technik. K nejběžnějším patří **kvadrantové testy náhodnosti**, **metoda nejbližších vzdáleností** a její modifikace a v poslední době nejvíce doporučované testy využívající **K-funkce**. Náhodnost bodového vzorku se zkoumá především v situacích, kdy body reprezentují místa určitých událostí (kriminální činy, výskyty nemocí, výskyty různých environmentálních či společenských jevů). Vzhledem ke skutečnosti, že potřebujeme odlišit terminologicky pozorované body od například náhodně rozmístěných bodů, je výhodné označit body reprezentující sledované jevy jako události a označení „body“ ponechat pro bodové objekty, které nerepresentují sledovaný jev. V této souvislosti se tedy hovoří o analýze náhodnosti událostí. Náhodnost událostí pak lze zkoumat z pohledu geografického jako testování náhodnosti rozmístění událostí (zkráceně náhodnosti textury událostí) nebo z pohledu časového jako testování náhodnosti výskytu událostí. Řada technik se pak zabývá současným časoprostorovým zkoumáním náhodnosti výskytu událostí.

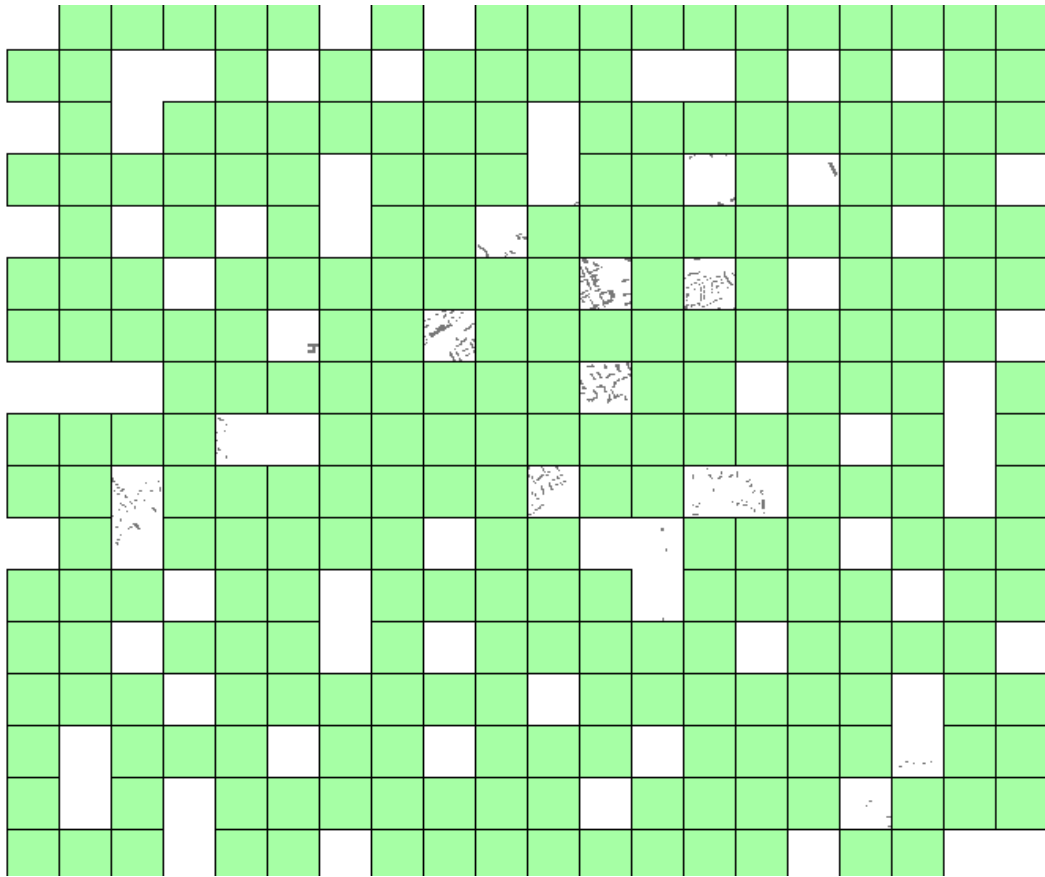
### 3.2.1 Kvadrantové testy náhodnosti

Základem kvadrantových testů náhodnosti je sledování četnosti událostí v uměle vymezených buňkách (kvadrantech). Předpokládáme, že četnost událostí má Poissonovu distribuci. Události musí být rozmístěny rovnoměrně v území, tvořit homogenní distribuci bez zjevného trendu (Smith et al. 2011, 251). Buňky mohou být různého tvaru (zpravidla však kruhové nebo čtvercové) a velikosti, v jedné oblasti  $\mathcal{R}$  se však používá konstantní tvar a velikost buněk. Mohou být rozmístěny pravidelně (obr. 3-8) nebo náhodně (obr. 3-10) v oblasti  $\mathcal{R}$ . Možnost náhodného rozmístění buněk může být

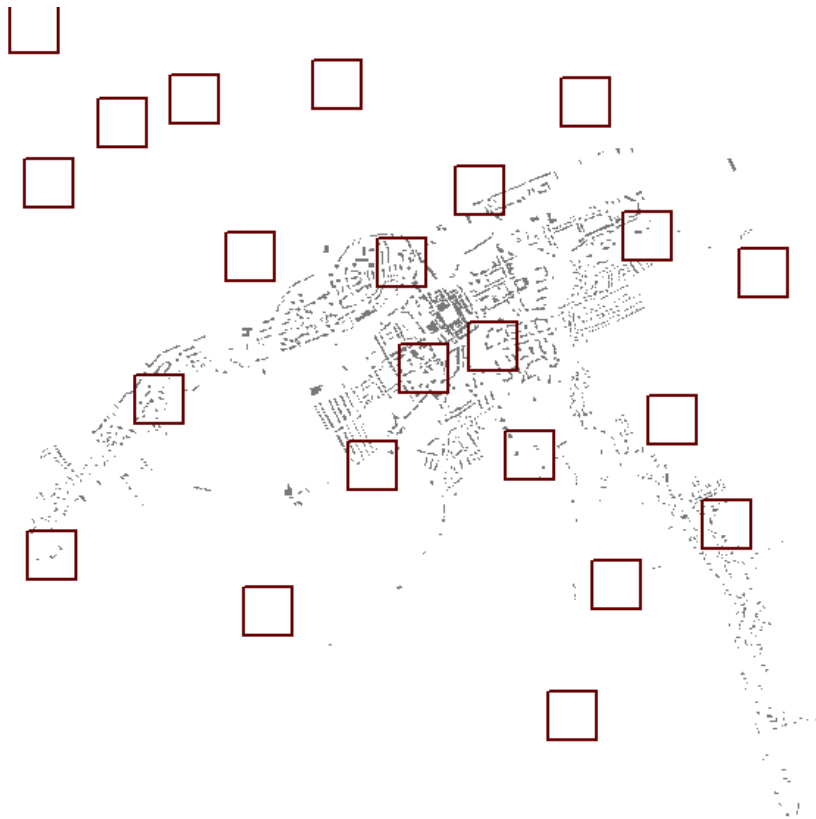
výhodná v situaci, kdy neznáme všechny výskyty událostí a považujeme je pouze za vzorek. Podobně i v případě pravidelného rozmístění buněk (často tvořících pravidelnou mřížku), lze náhodný vzorek simulovat náhodným výběrem z těchto buněk (obr. 3-9).



*Obr. 3-8 Pravidelná mřížka*



*Obr. 3-9 Pravidelná mřížka s náhodným výběrem*



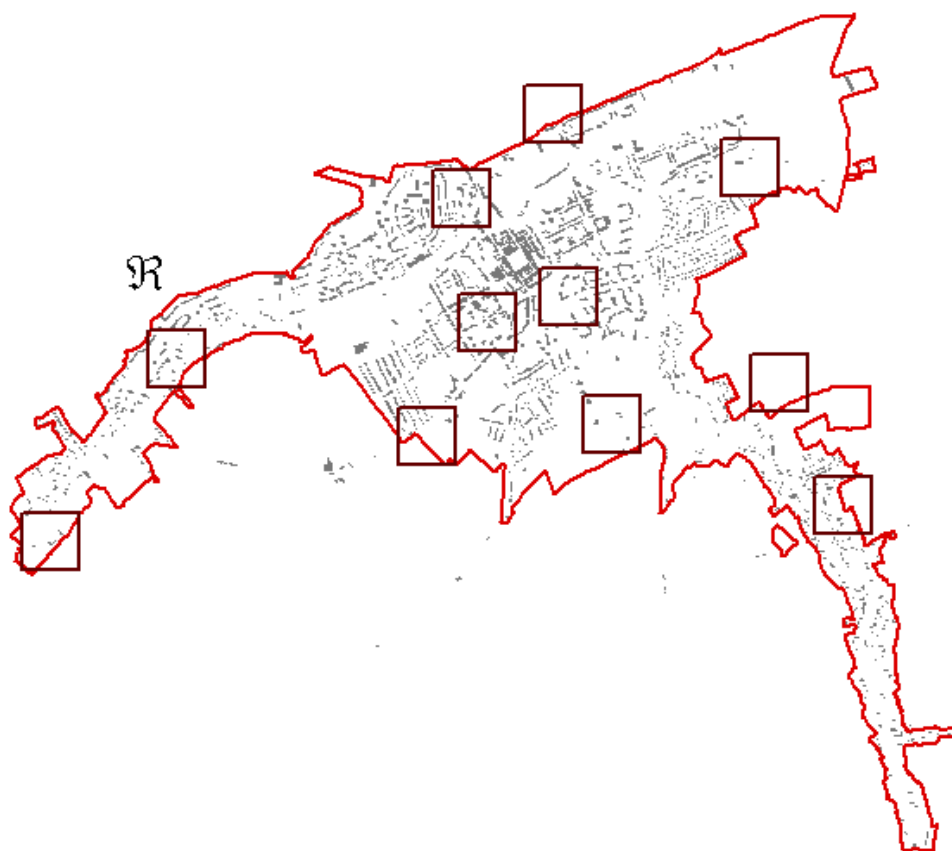
*Obr. 3-10 Náhodně rozmístěné buňky*

Označme  $m$  počet buněk vybraných k testování,  $x_i$  pozorovaný počet událostí v buňce  $i$  a tedy  $(x_1, x_2, \dots, x_m)$  je vektor počtu událostí v  $m$  buňkách v oblasti  $\mathfrak{R}$ .

Za předpokladu CSR (kapitola 3.3.1) je teoretickou pravděpodobnostní distribuci  $X^2$  rozdělení  $\chi^2_{m-1}$ , která dává dobré výsledky za předpokladu, že  $m > 6$  a střední hodnota počtu událostí  $E(x) > 1$ .

Náhodně rozmístěné buňky se mohou překrývat, a to způsobuje problémy v evidenci výskytů událostí, protože pak již počet  $x_i$  není nezávislý. Musíme proto přijmout takové vzorkovací schéma, abychom vyloučili překrývající se buňky.

Také překryv kvadrantu s hranicí  $\mathfrak{R}$  může způsobovat problémy (obr.3-11). Proto se doporučuje vytvářet ochranné pásmo uvnitř  $\mathfrak{R}$ , aby k tomuto překrytu nedošlo.



Obr. 3-11 Náhodně rozmístěné buňky, část z nich postižena hraničním problémem (hranici tvoří červená linie)

### 3.2.1.1 Test při náhodném rozmístění buněk

Náhodného rozmístění buněk v  $\mathfrak{R}$  se používá při odhadu intenzity událostí  $\lambda'$  v  $\mathfrak{R}$ . Musí ovšem platit předpoklad, že distribuce odpovídá homogennímu Poissonovu procesu. Pak

$$\lambda' = \frac{E(x)}{Q}$$

$E(x)$  aritmetický průměr počtu událostí v  $m$  buňkách  
 $Q$  plocha buňky (všechny jsou stejné)

Rozptyl lze přibližně odhadnout jako

$$VAR(\lambda) = \frac{\lambda'}{m * Q}$$

a interval spolehlivosti může být odvozen např. pro 95% hladinu významnosti jako

$$\lambda' \pm 2 \sqrt{\frac{\lambda'}{m * Q}}$$

### 3.2.1.2 Index disperze

Místo náhodného rozmístění buněk se ale zpravidla využívá metoda pravidelné mřížky a tu využívá i většina testů náhodnosti rozmístění událostí. Jednoduché testy využívají základní vlastnosti Poissonovy distribuce, že aritmetický průměr se rovná rozptylu (parametr  $\lambda$ ). Vypočteme tedy  $E(x)$  jako aritmetický průměr a  $VAR(x)$  jako rozptyl počtu událostí v  $m$  buňkách v oblasti  $\mathfrak{R}$ .

Potom poměr VMR se nazývá **indexem disperze** a ICS **indexem velikosti shluků**.

$$VMR = \frac{VAR(x)}{E(x)} \quad ICS = VMR - 1$$

Pokud  $VMR = 1$ , resp.  $ICS = 0$ , jde o **náhodný vzorek**, protože distribuce odpovídá požadavku Poissonovy distribuce. Pokud je  $VMR > 1$  (resp.  $ICS > 0$ ), indikuje se **shlukování událostí**. V případě, že  $VMR < 1$  (resp.  $ICS < 0$ ), naznačuje to existenci **pravidelného vzorku**.

Pozor při hodnocení na autokorelační efekt viz kapitola 4, který může výsledky deformovat.

Testování náhodnosti je založeno na prokázání významnosti odchylky jedním či druhým směrem od náhodného stavu ( $VMR=1$ ,  $ICS=0$ ). K testování se používá  $t$  test nebo  $\chi^2$  test dobré shody. Ve všech případech ( $t$  test i  $\chi^2$  test) bývá nulová hypotéza zamítnuta, je-li pravděpodobnost testované charakteristiky menší než kritická hodnota (běžně 0.05).

#### 1.1.1.1.9 $\chi^2$ test pro index disperze

$$\chi^2 = (m - 1) * VMR = \frac{(m - 1) * VAR(x)}{E(x)} = \frac{\sum_{i=1}^m (x_i - E(x))^2}{E(x)}$$

kde  $m$  je počet kvadrantů,  $E(x)$  je průměr pozorovaných událostí a  $VAR(x)$  jejich rozptyl.

#### 1.1.1.1.10 $t$ -test pro index disperze

$$t = \frac{VMR - 1}{S_{VMR}} = \frac{\frac{VAR(x)}{E(x)} - 1}{S_{VMR}} = \frac{(m - 1) * (VMR - 1)}{4}$$

s  $m-1$  stupni volnosti

Výpočet  $t$  testu zde vyžaduje standardizovanou chybu poměru VMR, odhadovanou jako:



$$S_{\text{VMR}} = \frac{4}{(m-1)}$$

### 3.2.1.3 $\chi^2$ test

Je možné se setkat i s následující variantou  $\chi^2$  testu pro zjišťování náhodnosti u kvadrantové metody.

Používá se  $\chi^2$  test ve tvaru, kde  $O_i$  je četnost pozorovaná a  $E_i$  četnost očekávaná s k-2 stupni volnosti (k je počet tříd).

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$$

Spočítáme počet událostí ve čtvercích, každý čtverec je klasifikován podle pozorované četnosti a tak získáme hodnoty  $O_i$ . Aby byl výpočet dostatečně věrohodný, měla by mít každá třída alespoň 5 čtverců, jinak se třídy spojují.

Očekávané četnosti  $E_i$  vypočteme z Poissonovy distribuce:

$$p(x) = \frac{\lambda^x}{x!} e^{-\lambda}$$

kde x bude četnost bodů ve čtverci a  $\lambda$  očekávaná střední hodnota (intenzita).

$$E_i = p(x) * m$$

kde m je celkový počet čtverců

Příklad:

počet událostí ve čtverci	počet čtverců $O_i$	k	p (výskytu n událostí)	očekávaný počet čtverců $E_i$
0 událostí	10	1	$p(0)=0,195$	$0,195*52=10$
1 událost	17	2	$p(1)=0,319$	$0,319*52=17$
2 události	12	3	$p(2)=0,261$	$0,261*52=14$
3 události	8	4	$p(3)=0,142$	$0,142*52=7$
4 události	5	5	$p(4)=0,058$	$0,058*52=3$
suma	52			

$$\lambda = 85/52=1,6346$$

$$p(0) = \frac{1,6346^0}{0!} e^{-1,6346}$$

$$\chi^2 = (10-10)^2/10 + (17-17)^2/17 + (12-14)^2/14 + (8-7)^2/7 + (5-3)^2/3=0,29+0,14+1,33=1,76$$

### 3.2.1.4 Problémy při kvadrantových metodách

Při kvadrantové analýze je citlivou volba sítě (při pravidelné mřížce) - počátek a velikost buněk. Pokud změním velikost kroku sítě, změní se i VMR. Na základě empirických zkušeností se doporučuje, aby střední hodnota počtu událostí v kvadrantu byla kolem 1.6).

Obecným nedostatkem kvadrantových metod je, že se nebere ohled na relativní polohu kvadrantů a relativní polohu událostí v kvadrantu. Pro základní situaci, kdy máme zjištěny všechny události v  $\mathcal{R}$  a jsou k dispozici počty událostí ve velkém počtu navazujících buněk, je možné uvažovat o využití informace o relativní pozici buňky nejenom informace o počtu událostí v buňce.

Někdy se namísto testování  $X^2$  distribuce používá Kolmogorov-Smirnovův test jako silnější a pružnější přístup (Smith et al., 2011, 251)

Problém volby velikosti buňky sítě pomáhá eliminovat Greig-Smithova procedura, která přináší i možnost částečně využít informaci o relativní poloze buňky.

**Greig-Smith procedura** vypočítává rozptyl počtu událostí v buňkách z originální mřížky a pak rozptyl pro další odvozené mřížky, které jsou vytvářeny postupným spojováním sousedních buněk v originální mřížce do bloků větší velikosti. Odhady rozptylu jsou vykresleny do grafu proti velikosti bloku (buňky) a opět vrcholy a poklesy na křivce jsou interpretovány jako přítomnost shlukování či pravidelného vzoru. Odkazy na použití této metody:

<http://www.jstor.org/pss/2530919>

[http://www.apsnet.org/phyto/PDFS/1993/Phyto83n04\\_419.pdf](http://www.apsnet.org/phyto/PDFS/1993/Phyto83n04_419.pdf)

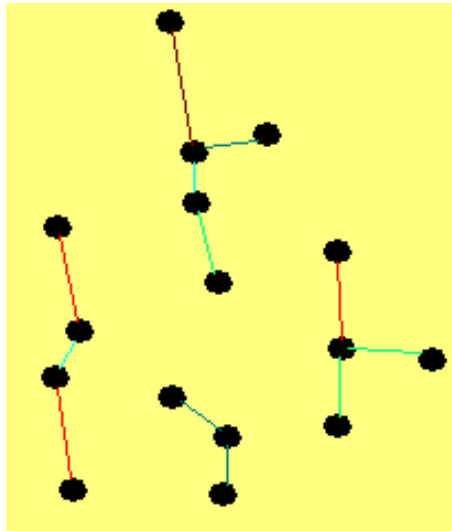
Existují i průzkumné metody založené na „řezu“ (tedy zjišťování počtu událostí v řádku). To může být užitečné pro hledání vzorkovacího schématu při terénních pracích.

### 3.2.2 Metoda nejbližších vzdáleností

Analýza nejbližších vzdáleností (*nearest neighbour distances*) (někdy také analýza nejbližšího souseda) studuje vzdálenosti mezi body, a to především mezi nejbližšími sousedními body. Neposuzuje tedy celkový vzorek a tak nepředchází základnímu nedostatku kvadrantového přístupu.

Při aplikaci této metody je nutné nejdříve rozhodnout, jaký typ vzdáleností se bude používat. Můžeme definovat různé druhy vzdáleností mezi událostmi. Budeme popisovat především vzdálenosti zjišťované podle principu nejbližšího souseda. Zvláště se zajímáme o vzdálenosti událost – událost  $W$ , které představují vzdálenosti mezi náhodně vybranou událostí a nejbližší událostí. Jinou možnost představují vzdálenosti bod – událost  $X$ , tedy vzdálenosti mezi náhodně vybraným bodem ve studované oblasti a nejbližší pozorovanou událostí. Obě tyto míry mohou být použity, pokud jsou k dispozici všechny události v  $\mathcal{R}$ .  $W$  je **nedefinované**, pokud nemůžeme provést náhodný výběr z událostí nebo pokud nejsou v oblasti všechny události. Naproti tomu  $X$  je užitečné i při testování vzorkování – můžeme vybrat náhodný bod v  $\mathcal{R}$  a pak nalézt nejbližší událost a měřit vzdálenost k ní.

Metoda nejbližších vzdáleností je založena na zkoumání pozorované distribuce jedné nebo dvou z těchto nejbližších vzdáleností ( $W$  nebo  $X$ ). Mějme na paměti, že nejbližší vzdálenosti poskytují informaci jenom o interakcích mezi událostmi ve velkém měřítku, tedy do relativně malých vzdáleností. Často ale můžeme využít i studia variací ve větších vzdálenostech bodového vzorku v celém  $\mathcal{R}$ .



Obr. 3-12 Nejblíže vzdálenosti  $W$  mezi událostmi (barevně rozlišená délka spojnice)

Nejčastěji se studuje kumulativní relativní křivka výskytu zjištěných vzdáleností (kumulativní křivka pravděpodobnosti), tedy distribuční funkce  $G(w)$  pro  $W$ , respektive  $F(x)$  pro  $X$ :

$$G'_{(w)} = \frac{COUNT(w_i \leq w)}{n}$$

$$F'_{(x)} = \frac{COUNT(x_i \leq x)}{m}$$

Tytéž funkce vyjádřené pomocí indikátorové funkce  $I$ :

$$G'_{(w)} = \frac{\sum_{i=1}^n I_i(w)}{n} \quad I_i = \begin{cases} 0 & \dots w_i > w \\ 1 & \dots w_i \leq w \end{cases}$$

$$F'_{(x)} = \frac{\sum_{i=1}^m I_i(x)}{m} \quad I_i = \begin{cases} 0 & \dots x_i > x \\ 1 & \dots x_i \leq x \end{cases}$$

### 3.2.2.1 Teoretický model nejblíže vzdáleností pro CSR

Pro distribuci  $W$  a  $X$  můžeme odvodit jisté teoretické výsledky, pokud budeme hraniční efekt zanedbávat.

Nechť je  $\lambda$  střední intenzita událostí na plošnou jednotku. Potom za předpokladu CSR (kapitola 3.3.1) jsou události nezávislé a počet událostí v libovolné ploše odpovídá Poissonově distribuci. Pak pravděpodobnost toho, že žádná událost nebude ležet v kruhu s poloměrem  $x$  kolem libovolně vybraného bodu, je dána jako  $\exp(-\lambda\pi x^2)$ .

Odvození:

Očekávaný počet událostí je:

$$E = \lambda * plocha = \lambda * \pi * x^2$$

Pravděpodobnost pro jev, kdy se v ploše nevyskytuje žádná událost:

$$P(0) = \frac{E^0}{0!} e^{-E} = e^{-\lambda\pi x^2}$$

Potom distribuční funkce  $F_{(x)}$  pro CSR je:

$$F_{(x)} = P(X < x) = 1 - e^{-\lambda\pi x^2}$$

pro  $x \geq 0$

Jde o doplněk pravděpodobnosti, že do vzdálenosti  $x$  leží aspoň jeden bod.

To znamená, že  $\pi x^2$  odpovídá exponenciální distribuci s parametrem  $\lambda$ , nebo také že  $2\pi\lambda X^2$  odpovídá distribuci  $\chi^2_2$ . Z toho můžeme odvodit:

$$E_{(x)} = \frac{1}{2\sqrt{\lambda}} \quad \text{Var}(x) = \frac{(4 - \pi)}{4\lambda\pi}$$

To také znamená, že pokud  $X_1, X_2, \dots, X_n$  jsou nezávislé na nejbližší vzdálenosti bod-událost, pak  $2\pi\lambda \sum X_i^2$  má distribuci  $\chi^2_{2n}$ .

Úplně stejné argumenty a postup bychom mohli uplatnit pro teoretický model CSR pro nejbližší vzdálenosti událost-událost  $W$ .

$$G(w) = \Pr(W \leq w) = 1 - e^{-\lambda\pi w^2} \quad w \geq 0$$

$E(w)$  a  $VAR(w)$  jsou stejné jako pro  $X$ .

Jiný přístup k odvození průměru a rozptylu pro teoretický model, založený na derivaci přírůstku kruhu a následné integraci, uvádí Smith et al. (2011, 261).

### 3.2.2.2 Varianty vyhodnocení metody

V zásadě můžeme postupy vyhodnocení distribuce vzdáleností rozdělit na **grafické** (zkoumají průběh grafu distribuce vzdáleností) a **numerické** (které vypočítají jistou číselnou charakteristiku a prověřují její náhodnost).

#### Grafické postupy

1. Graf závislosti  $G'(w)$  na  $W$
2. Graf závislosti  $F'(x)$  na  $X$
3. Graf závislosti  $G'(w)$  na  $F'(x)$
4. Graf závislosti  $G'(w)$  na teoretické  $G(w)$  pro kompletně vyčíslená pole
5. Graf závislosti  $F'(x)$  na teoretické  $F(x)$  pro kompletně vyčíslená pole
6. Graf závislosti  $G'(w)$  na průměru simulací  $\underline{G}(w)$  a pás spolehlivosti (vykreslení horního a spodního limitu spolehlivosti)

#### Numerické postupy

7. Test NNI
8. Clark-Evansův test

9. Hopkinsonův test
10. Byth-Ripleyův test

U některých **grafických postupů** je potřebné odvodit teoretickou funkci  $G(w)$  resp.  $F(x)$ , tedy funkce odpovídající distribuci CSR (kapitola 3.3.1).

Znalost teoretické distribuce  $W$  a  $X$  při CSR nám dovoluje odvodit (přinejmenším přibližně) distribuci statistických charakteristik z pozorovaných nejbližších vzdáleností. Ty můžeme použít jako základ pro testování CSR. Měli bychom poznamenat, že distribuční teorie pro většinu takových testů předpokládá, že nejbližší vzdálenosti použité k výpočtu sumární statistiky jsou nezávisle vzorkované z  $\mathfrak{R}$  a že žádná z nich není vychýlena hraničním faktorem.

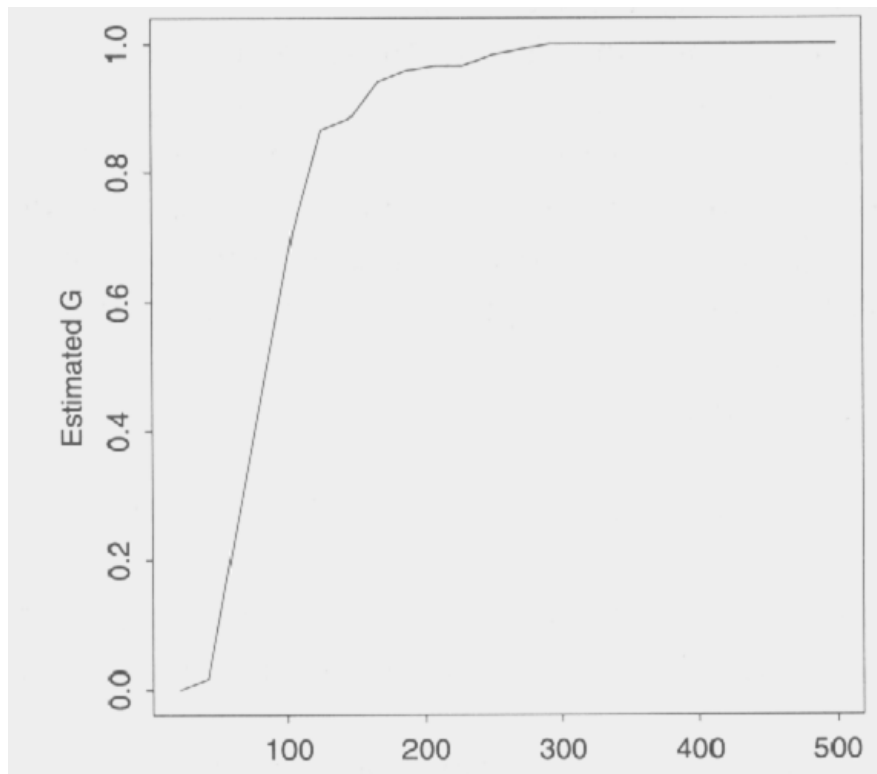
Předpoklad nezávislosti je pravděpodobně neplatný, pokud je celkový počet událostí ve studované oblasti malý a podíl událostí, který je vybrán do vzorku pro nejbližší vzdálenosti, je velký; představte si, že byste měli v oblasti pouze 2 události, pak 2 nejbližší vzdálenosti by byly identické. Jedním z možných doporučení je, aby počet  $m$  nejbližších vzdáleností byl  $m \leq 0,1 * n$ , kde  $n$  je celkový počet událostí. Obecný efekt nedostatku nezávislosti se projeví tak, že testovaná statistika bude mít větší rozptyl než teoretické hodnoty za předpokladu nezávislosti, což znamená, že standardní test může být označen za významně odchýlený od modelu CSR. Zvláště by měl být uživatel upozorněn, aby nepoužíval test nejbližších vzdáleností bez vhodných korekcí na sadu všech nejbližších vzdáleností událost-událost, pokud je k dispozici kompletní mapa událostí v  $\mathfrak{R}$ . Tento test totiž vyžaduje náhodný vzorek událostí, proto je samozřejmě výhodné, pokud jsou bodová pole náhodně vzorkovaná.

Další problém je oprava **hraničního efektu** (kapitola 3.2.2.3). Nejbližší vzdálenosti pro události blízké hranici  $\mathfrak{R}$  budou vychýlené a budou mít tendenci být větší než vzdálenosti pro události uvnitř regionu. Některé testy mají hraniční korekci zabudovanou, ale obecně se to dá aplikovat jen na omezenou sadu pravidelných tvarů oblasti.

Základem **numerických postupů** je redukce komplexního bodového vzorku do jednoduché sumární statistiky nejbližších vzdáleností pro testování CSR. To vede přece jen k významné ztrátě informace. Testy pouze indikují odchylky od CSR; bohužel je málo známo o chování těchto statistik, pokud CSR není dodrženo, rovněž nejsou k dispozici žádné informace pro působení jiných alternativních modelů, pokud je CSR odmítnuta. Obecně nemůžeme doporučit tento přístup, pokud je k dispozici celé vymapované pole, snad jedině pro případ předběžné analýzy.

#### **1.1.1.1.11 Graf závislosti $G'(w)$ na $W$ , graf závislosti $F'(x)$ na $X$**

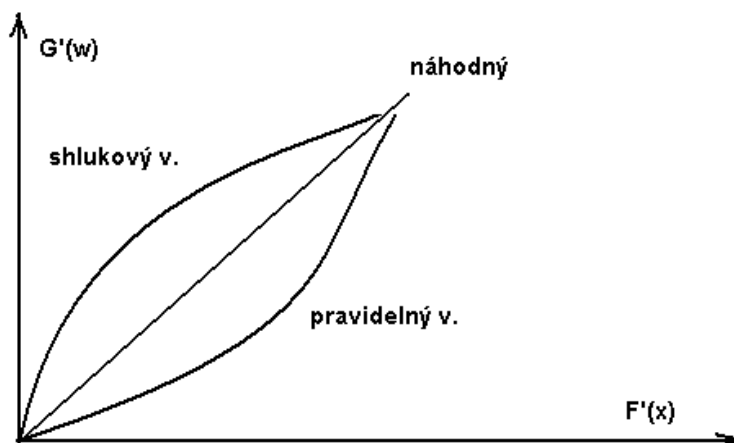
Výsledná empirická distribuční funkce  $G'(w)$  nebo  $F'(x)$  je vynesena do grafu v závislosti na  $W$  resp.  $X$  a interpretována pouze vizuálně. Např. pokud distribuční funkce roste příkře na začátku rozsahu a pak se vyrovná, naznačuje to, že je zastoupeno hodně krátkých vzdáleností na rozdíl od dlouhých, a že tedy zřejmě půjde o vzorek se shluky událostí. Naopak, pokud křivka roste až při konci rozsahu vzdáleností, naznačuje to jisté opakování nebo pravidelnost. Posuzování je značně subjektivní.



Obr. 3-13 Distribuční funkce nejbližších vzdáleností  $G(w)$  pro lokalizaci sopek v Ugandě (Bailey, Gatrell 1995)

#### 1.1.1.1.12 Graf závislosti $G'(w)$ na $F'(x)$

Vykresluje se  $G'(w)$  proti  $F'(x)$  do grafu. Pokud se neprojevuje interakce, budou si obě distribuce podobné a měli bychom očekávat přímou linii v grafu. V případě pozitivní interakce (shlukování) budou vzdálenosti bod-událost  $x_i$  větší než vzdálenosti událost-událost  $w_i$  a můžeme očekávat, že  $G'(w)$  překročí  $F'(x)$ . Opačný případ nastane pro pravidelný vzorek. Posuzování je značně subjektivní.



Obr. 3-14 Graf závislosti  $G'(w)$  na  $F'(x)$

### 1.1.1.1.13 Graf závislosti $G'(w)$ na teoretické $G(w)$ , graf závislosti $F'(x)$ na teoretické $F(x)$

Do grafu se vykresluje závislost  $\underline{G}(w)$  nebo  $\underline{F}(x)$  přímo vůči jejich teoretickým ekvivalentům pro CSR, tedy:

$$G(w) = 1 - e^{-\lambda \pi w^{**2}} \quad \text{nebo} \quad F(x) = 1 - e^{-\lambda \pi x^{**2}}$$

Interpretace takového grafu je určitě méně subjektivní než v předchozích případech, kde se kreslily grafy  $G'(w)$  nebo  $F'(x)$  vůči  $W$  resp.  $X$ , a zkoumaly se pouze s ohledem na jejich obecný tvar. Minimálně je možné srovnávat se známým teoretickým tvarem pro CSR. Stále však chybí formální způsob, jak ocenit významnost rozdílů v grafu. Žádný takový způsob není k dispozici, díky komplexním problémům vzniklým při odhadování hraniční korekce použité při odhadu  $G(w)$  nebo  $F(x)$ .

### 1.1.1.1.14 Graf závislosti $G'(w)$ na průměru simulací $\underline{G}(w)$ a pás spolehlivosti

Snad nejspokojivějším přístupem (i když výpočetně nejintenzivnějším) je srovnání odhadované distribuční funkce (bez hraniční korekce) se simulačním odhadem jejich teoretických distribučních funkcí při CSR (kapitola 3.3.1) za přítomnosti určitého hraničního efektu vzniklého v dané studijní oblasti  $\mathfrak{R}$ .

Představíme tuto metodu na příkladu  $W$ , analogický přístup lze použít pro  $X$ . Simulační odhad pro  $G'(w)$  při CSR je vypočten jako

$$\hat{G}(w) = \sum \frac{\hat{G}_i(w)}{m}$$

kde  $G_i(w)$  pro  $i = 1, 2, \dots, m$  jsou empirické distribuční funkce, z nichž každá je odhadována (odvozována) bez hraniční korekce z 1 z  $m$  nezávislých simulací  $n$  událostí při CSR v  $\mathfrak{R}$ , tedy  $n$  událostí nezávislých a stejnoměrně distribuovaných v  $\mathfrak{R}$ .

Pro účely ocenění významnosti odchylek mezi simulovanou CSR distribucí  $G(w)$  a tou, která byla skutečně pozorována  $G'(w)$  se také definuje horní a dolní simulační obálka (limit):

$$U(w) = \max_{i=1, \dots, m} \{G_i(w)\}$$

$$L(w) = \min_{i=1, \dots, m} \{G_i(w)\}$$

Nyní vykreslíme  $G'(w)$  (odhadované bez hraniční korekce) proti  $G(w)$  a do grafu přidáme  $U(w)$  a  $L(w)$ . Pokud data odpovídají CSR, bude závislost  $G'(w)$  na  $G(w)$  přibližně lineární pod úhlem  $45^\circ$ . Pokud se projevuje shlukování, bude křivka závislosti ležet nad ideální linií ( $45^\circ$ ) a opačně pro pravidelný vzor.  $U(w)$  a  $L(w)$  nám pomohou ocenit význam odchylek od ideální linie  $45^\circ$ , protože mají vlastnost:

$$\Pr(\hat{G}(w) \succ U(w)) = \Pr(\hat{G}(w) \prec L(w)) = \frac{1}{m+1}$$

To také ilustruje, jakou hodnotu bychom měli použít pro  $m$ , tedy kolik simulací bychom měli provést, abychom byli schopni rozeznat odchylky na zvolené hladině významnosti.

### 1.1.1.1.15 TEST NNI

Test NNI vychází z výpočtu poměru mezi pozorovanou a očekávanou střední hodnotou minimální vzdálenosti mezi událostmi NNI (*near neighbour index*). Doporučuje se mít alespoň 100 událostí (Smith et al., 2011, 261).

Platí, že průměrné vzdálenosti mezi událostmi u shlukového vzorku jsou menší než u náhodně rozptýlených událostí a ty jsou opět menší než u pravidelně rozmístěných událostí.

Platí, že:

NNI  $\geq 0$  AND NNI  $\leq 2.1491$

NNI = 0 ... shlukový vzorek

NNI = 1 ... čistě náhodný vzorek

maximální hodnota NNI – pravidelně rozptýlený vzorek

Postup:

a) výpočet očekávané hodnoty (pro CSR)  $r_e$

$$r_e = 0,5 * \sqrt{\frac{R}{n}}$$

kde  $R$  je plocha oblasti  $\mathfrak{R}$  a  $n$  počet událostí.

b) výpočet průměrné hodnoty  $r_o$  z měření

Zjištěn nejbližší soused pro každou událost, evidována příslušná vzdálenost  $w$  a nakonec vypočtena průměrná hodnota této minimální vzdálenosti.

$$r_o = \frac{\sum w_i}{n}$$

c) výpočet NNI

$$NNI = \frac{r_o}{r_e}$$

NNI se někdy používá jako popisná statistika ke srovnání distribuce různých fenoménů v téže oblasti (např. vzdálenosti mezi stromy různých druhů v téže zalesněné oblasti)

d) test významnosti

Provádí se statistický test významnosti rozdílu mezi  $r_o$  a  $r_e$ .

Nulová hypotéza: „NNI je odlišná od 1 jen jako výsledek vzorkovací chyby“ (NNI = 1 pro náhodné rozmístění událostí, odchylka od 1 je náhodná).

Používá se Z-test, založený na směrodatné odchylce. (Z jako normovaná veličina pro normální distribuci, tedy tzv. Z-skóre.)

$$Z = \frac{|r_o - r_e|}{S_d} \quad S_d = \frac{0,26136}{\sqrt{\frac{n * n}{R}}} = \frac{0,26136}{n} * \sqrt{R}$$



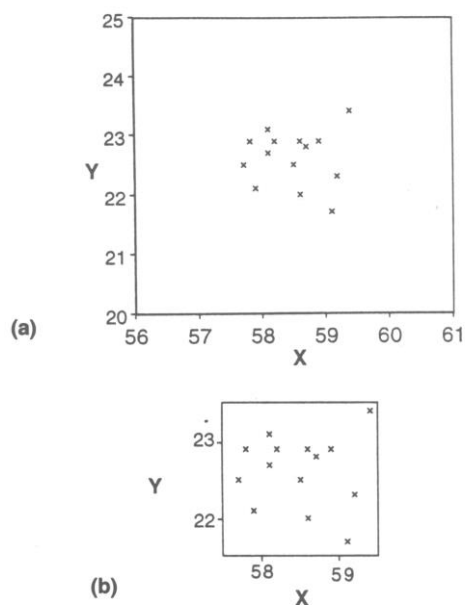
kde  $R$  je plocha oblasti  $\mathfrak{R}$  a  $n$  počet událostí.

Je-li  $Z >$  kritická hodnota (pro hladinu významnosti 0.05 nebo 0.01) pro normovanou normální distribuci, pak se zamítá nulová hypotéza, tedy konstatujeme, že vzorek je významně odlišný od náhodného.

Interval spolehlivosti pro hladinu 0.05 :  $NNI - 1.96 S_d$  až  $NNI + 1.96 S_d$

S analýzou nejbližšího souseda jsou spojeny 2 hlavní problémy:

- a) Charakteristika NNI a směrodatná odchylka silně závisí na velikosti studované oblasti.



Obr. 3-15. Stejný vzorek ve 2 různě velkých oblastech

- b) Někdy je minimální vzdálenost mezi událostmi nevhodná pro výpočet NNI např. jsou-li události v pravidelně rozmístěných shlucích (každý s podobným počtem bodů a vzdálenostmi mezi nimi). Pak totiž analýza NNI indikuje shlukovaný vzorek, i když jsou shluky rozptýleny. Někdy se proto používá 2., 3. nebo 4. nejbližší vzdálenost. Např. distribuce živočišných druhů může být charakterizována "rodinami" zabydlujícími určité teritorium (stádo jelení zvěře, smečka ...). Vzdálenosti uvnitř skupiny jsou malé. Proto, i když jsou vzdálenosti mezi skupinami velké, bude NNI odpovídat jednomu shluku.

### 1.1.1.16 Clark-Evansův test

Vypočte se charakteristika

$$\bar{w} = \frac{\sum w_i}{m}$$

a ta se srovnává s hodnotami normální distribuce s parametry:

$$N\left(\frac{1}{2\sqrt{\lambda}}, \frac{(4-\pi)}{4\lambda\pi m}\right)$$

Test je založen na vzdálenostech mezi událostmi a proto vyžaduje, aby bylo k dispozici kompletně vyčíslené bodové pole, ze kterého jsou náhodně vybírány jednotlivé události, pro které se určí jejich nejbližší vzdálenosti. To je posilováno faktem, že  $\lambda$  je neznámé a musíme ho nahradit vhodným odhadem, zřejmě

$$\hat{\lambda} = \frac{n}{R}$$

kde  $n$  je počet událostí a  $R$  plocha  $\mathfrak{R}$ .

Bylo by vhodné namísto náhodného výběru vzorku  $m$  událostí použít všech vzdáleností událost-událost (pro  $n$  událostí). Proto byly navrženy vhodné korekce pro  $E(w)$  a  $VAR(w)$ , které dovolují použít všechny nejbližší vzdálenosti místo vzorku ( $m=n$ ). Pak:

$$E(\bar{w}) = 0,5\sqrt{\frac{R}{n}} + 0,051\frac{P}{n} + 0,041\frac{P}{n^{3/2}}$$

$$VAR(\bar{w}) = 0,070\frac{R}{n^2} + 0,037P\sqrt{\frac{R}{n^5}}$$

kde  $P$  je obvod oblasti  $\mathfrak{R}$ , která má plochu  $R$ . Tyto aproximace neplatí pro velmi konvolutní oblasti  $\mathfrak{R}$ .

#### 1.1.1.1.17 Hopkinsův test

Hopkinsův test srovnává  $\sum x_i^2 / \sum w_i^2$  s procentními body distribuce  $F_{2m,2m}$ . Podstata testu spočívá v tom, že pro shlukový vzorek jsou vzdálenosti bod-událost ( $x_i$ ) relativně větší než vzdálenosti událost-událost ( $w_i$ ), a opačně v případě pravidelného vzorku. Protože test používá  $w_i$ , opět vyžaduje mít k dispozici kompletní sadu  $n$  událostí v oblasti  $\mathfrak{R}$  tak, aby mohly být náhodně vybrány vzdálenosti mezi událostmi. Navrhuje se, aby byl výběr událostí prováděn společně s výběrem bodových vzorků, pokud je použito „semisystematické“ výběrové schéma, zatímco je přes studovanou oblast přeložena pravidelná mřížka studovaných bodů. Alternativně jsou studijní body použity jako body pro výpočet vzdáleností bod-událost  $x_i$ . Kolem každého ze zbývajících studijních bodů se vytvoří (předpokládá) malá kruhová oblast, uvnitř které jsou vyčísleny všechny události (stačí tak velká oblast, aby obsahovala aspoň 5 událostí); výběr náhodných událostí pro výpočet vzdáleností událost-událost  $w_i$  je založen na populaci událostí uvnitř kolekce těchto úplně vyčíslených oblastí.

#### 1.1.1.1.18 Byth & Ripleyův test

Byth & Ripleyův test srovnává

$$\frac{1}{m} \sum \frac{x_i^2}{(x_i^2 + w_i^2)}$$

s procentními body distribuce  $N(1/2, 1/12m)$ ,

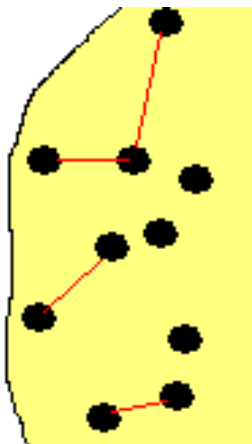
kde  $x_i$  hodnoty jsou náhodně párovány s  $w_i$  hodnotami.

Podmínky pro tento test jsou velmi podobné Hopkinsonovu testu.

### 3.2.2.3 Hraniční efekt

U hraničního efektu lze očekávat, že nejbližší vzdálenosti pro události ležící blízko hranic  $\mathfrak{R}$  jsou relativně větší než uvnitř  $\mathfrak{R}$ , protože událost poblíž hranice nemá možnost sousedit s událostí vně oblasti (obr. 3-16). To může být vážný problém ve všech výpočtech a hodnoceních. Při většině vyhodnocení, zvláště u relativně malých oblastí, je nutné se pokusit hraniční efekt eliminovat.

Vyhodnocování distribučních křivek  $G_{(w)}$  a  $F_{(x)}$  je proto komplikováno skutečností, že nejsou známy teoretické tvary těchto funkcí (za předpokladu CSR) pro konkrétní tvar oblasti  $\mathfrak{R}$ . Díky tzv. **hraničnímu efektu** bude průběh křivky vždy komplexním způsobem záviset na tvaru oblasti  $\mathfrak{R}$ .



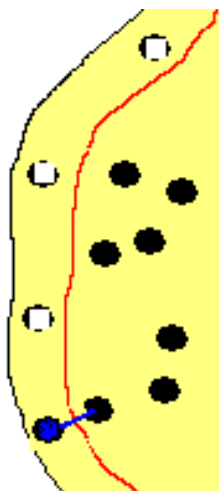
Obr. 3-16. Vybrané 4 události při hranici a jejich vzdálenosti k nejbližšímu sousedovi

Možnosti řešení hraničního efektu:

- a. Fixní ochranné pásmo
- b. Toroidální korekce
- c. Adaptivní ochranné pásmo
- d. Simulační postupy

Ad a) Fixní ochranné pásmo

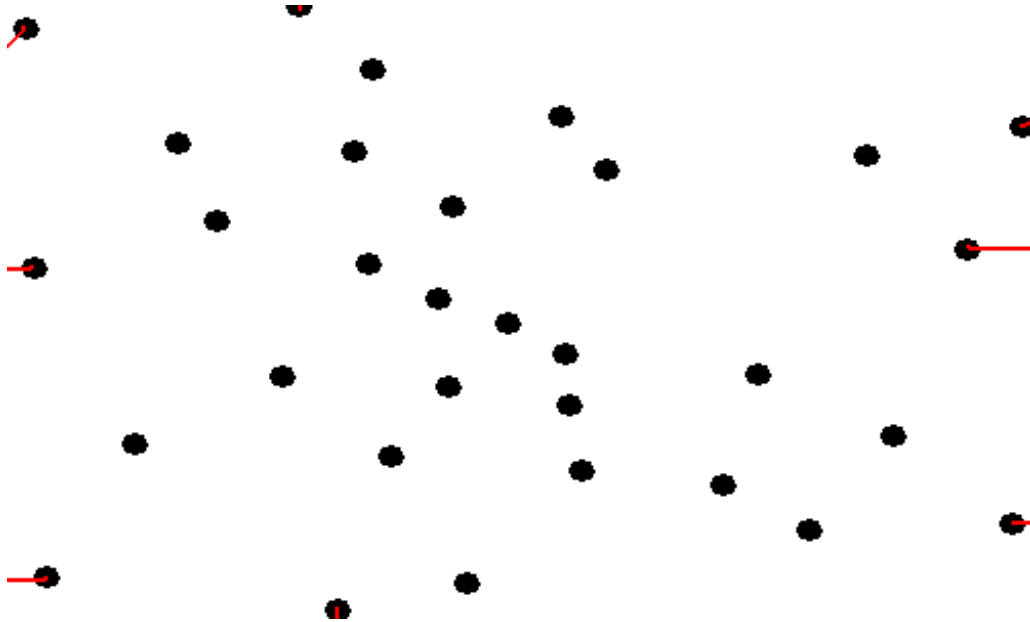
Problém může být řešen konstrukcí ochranného pásma (ohradníku) uvnitř obvodu  $\mathfrak{R}$ . Nejbližší vzdálenosti nejsou zjišťovány pro události (nebo body pro  $X_i$ ) ležící v této hraniční zóně, ale události v zóně se mohou stát nejbližšími sousedy jiných událostí ve zbytku zóny (nebo bodů). Nevýhodou je stanovení konstantního ochranného pásma, což může vést k vyloučení řady událostí z hodnocení.



Obr. 3-17. Vybrané 4 události při hranici, které jsou uvnitř ochranného pásma

Ad b) Toroidální korekce

Pokud je oblast pravoúhlá, lze uplatnit tzv. toroidální hraniční korekci. Pak předpokládáme, že horní hrana oblasti je spojena se spodní hranou, levá hrana zase s pravou. De facto je nyní studována oblast o rozměru 3 x 3 původních oblastí. Připojené kopie oblasti dovolují hledat bližší události i za hranicí.



Obr. 3-18. Červeně vyznačené spojnice k nejbližšímu sousedovi přes okraj oblasti

Ad c) Adaptivní ochranné pásmo

Výpočet s hraniční korekcí lze vyjádřit buď

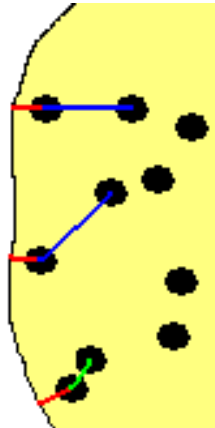
$$\hat{G}_{(w)} = \frac{\#(b_i \succ w \geq w_i)}{\#(b_i \succ w)}$$

kde # představuje počet a  $b_i$  je vzdálenost události  $i$  od nejbližšího bodu na hranici, nebo lze výraz zapsat pomocí indikátorových funkcí  $I$  a  $J$

$$G'_{(w)} = \frac{\sum_{i=1}^n I_i(w)}{\sum_{i=1}^n J_i(w)} \quad I_i = \begin{cases} 0 \dots (w_i \succ w) \text{ OR } (b_i \leq w) \\ 1 \dots (w_i \leq w) \text{ AND } (b_i \succ w) \end{cases} \quad J_i = \begin{cases} 0 \dots b_i \leq w \\ 1 \dots b_i \succ w \end{cases}$$

kde  $b_i$  je vzdálenost události  $i$  od nejbližšího bodu na hranici.

Postup ignoruje hodnoty  $w_i$  pro události blízké k okraji. Podobně lze vyjádřit  $F'(x)$ . Tímto postupem vzniká adaptivní šířka ochranného pásma, kde jsou z výpočtu vyloučeny jen ty události (resp. vzdálenosti k nim), které jsou příliš blízko hranice oblasti a může se u nich projevit hraniční efekt.



Obr. 3-19. Porovnání vzdáleností k nejbližšímu sousedovi a k hranici, akceptována pouze událost

#### Ad d) Simulační postupy

Pravděpodobně nejspolehlivější postup nabízí stochastická simulace. Analogicky s použitím simulací v následující kapitole se provádí náhodná simulace bodové textury (model CSR), z ní se vypočítá příslušná křivka, která odpovídá jedné realizaci funkce pro náhodnou texturu. Pokud provedeme dostatečně množství simulací (typicky 100), sada křivek nám vymezuje prostor, kde se může vyskytovat náhoda dle CSR. Jestliže empirická křivka vybočuje z pásma náhodných realizací, je zřejmé, že pozorovaná textura je nenáhodná.

### 3.2.3 K funkce

Redukovaná míra druhého momentu neboli K funkce poskytuje efektivní přehled prostorové závislosti událostí pro široký rozsah měřítek. Je těsně korelovaná s variabilitou intenzity 2.řádu, která nemůže být z pozorovaných událostí přímo odvozena.

Pro odvozování K-funkce a její interpretaci musí platit 2 základní předpoklady – distribuce událostí v  $\mathfrak{R}$  odpovídá homogennímu a izotropnímu procesu pro sledovanou měřítku (stacionarita). Nelze totiž rozlišit změny ve variacích intenzity 1.řádu (typu trend) a 2.řádu (typu kovariance). Pokud existují změny 1.řádu v celé oblasti  $\mathfrak{R}$ , lze někdy vybrat menší části oblasti a v nich studovat změny 2.řádu za předpokladu, že v malé oblasti budou změny 1.řádu zanedbatelné.

#### 3.2.3.1 Definice K-funkce

K funkce odpovídá standardizovanému průměrnému počtu událostí do vzdálenosti  $h$  od libovolné události. Standardizace se provádí pomocí intenzity  $\lambda$ .

Tedy

$$K_h = \frac{1}{\lambda} * E(N_h) \qquad K_h = \frac{1}{\lambda} * \frac{1}{n} * \sum N_h$$

- $K_h$  hodnota K funkce ve vzdálenosti  $h$ ,
- $N_h$  počet událostí do vzdálenosti  $h$  od libovolné události,
- $E()$  označuje střední hodnotu.
- $\lambda$  intenzita neboli střední počet událostí v plošné jednotce, který by měl být konstantní v celé  $\mathfrak{R}$

Hodnota  $K$  funkce představuje vypočtenou hustotu standardizovanou intenzitou  $\lambda$ , která je nezbytná k eliminování závislosti na konkrétní hustotě událostí v určité oblasti, neboť střední hodnota  $K$  funkce bude na celkovém počtu událostí a intenzitě událostí samozřejmě záviset.

Význam  $K(h)$  jako sumárního měření efektu 2.řádu je v tom, že jsme schopni odvodit z pozorovaného bodového vzorku jeho odhad  $K'(h)$ , zatímco nejsme schopni odhadnout přímo variability intenzity 2.řádu.

Je-li  $R$  plocha oblasti  $\mathfrak{R}$ , pak očekávaný počet událostí v  $\mathfrak{R}$  je  $n=\lambda \cdot R$ . Z definice  $K$ -funkce plyne, že očekávaný počet uspořádaných párů událostí ve vzdálenosti do  $h$  je  $\lambda^2 \cdot R \cdot K_{(h)}$  (tento způsob výpočtu je samozřejmě výhodnější než realizace vyplývající z definice - tj. navštívit každý bod a počítat vzdálenosti k okolním bodům). Proto vhodným odhadem  $K_{(h)}$  bude:

$$k'_{(h)} = \frac{1}{\lambda^2 R} \sum_{i \neq j} \sum I_h(d_{ij}) \quad \text{resp.} \quad k'_{(h)} = \frac{R}{n^2} \sum_{i \neq j} \sum I_h(d_{ij})$$

$d_{ij}$  vzdálenost mezi  $i$ -tou a  $j$ -tou pozorovanou událostí v  $\mathfrak{R}$

$$I_h(d_{ij}) = \begin{cases} 0 & \dots h < d_{ij} \\ 1 & \dots h \geq d_{ij} \end{cases}$$

Význam  $K$ -funkce se nejlépe pochopí z demonstrace. Kolem vybrané události se vytvoří sada koncentrických kruhů a zjišťuje se počet událostí v každém kruhu (větší kruh obsahuje automaticky i menší kruhy, jde o kumulativní funkci). Postupně jsou navštíveny všechny další události a celkové součty jsou přepočítány faktorem  $R/n^2$  (samozřejmě při zanedbání okrajového efektu).

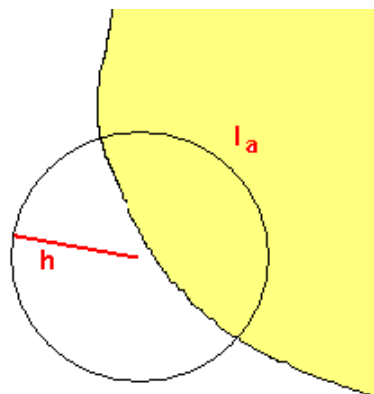
Stejně jako u analýzy nejbližších vzdáleností je nutno uvažovat o vlivu hranic  $\mathfrak{R}$  na náš odhad a o korekci této chyby.

### 3.2.3.2 Hraniční korekce pro $K$ funkci

Uvažujme kružnici se středem v události  $i$ , která prochází událostí  $j$ , pak  $w_{ij}$  bude podíl obvodu kružnice, který leží uvnitř  $\mathfrak{R}$  (obr. 3-20). Potom  $w_{ij}$  je podmíněnou pravděpodobností, že událost je pozorována v  $\mathfrak{R}$  v závislosti na vzdálenosti  $d_{ij}$  od  $i$ -té události. Výsledný tvar:

$$k'_{(h)} = \frac{1}{\lambda^2 R} \sum_{i \neq j} \sum \frac{I_h(d_{ij})}{w_{ij}}$$

$$w_i = \frac{l_a}{2\pi h}$$



Obr. 3-20. Odvození  $w$  pro hraniční korekci

Nyní ještě nahradíme neznámou intenzitu  $\lambda$  za její odhad. Nejjednodušší odhad je

$$\hat{\lambda} = \frac{n}{R} \quad k'_{(h)} = \frac{R}{n^2} \sum_{i \neq j} \sum \frac{I_n(d_{ij})}{w_{ij}}$$

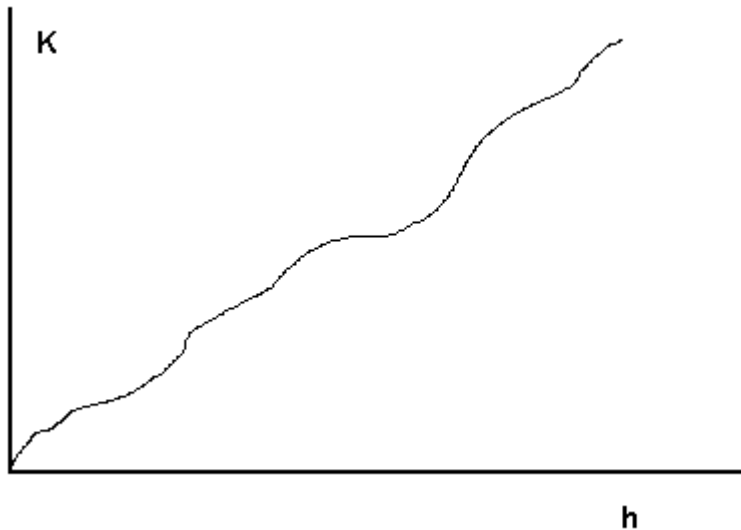
Okrajová korekce odhadu K-funkce je smysluplná, pokud  $h$  není příliš velké vzhledem k  $\mathfrak{R}$ . Je tedy nezbytné omezení  $h$ , jinak  $w_{ij}$  klesá absurdně až k mezní situaci dělení 0).

V praxi to není závažný problém, protože nás obvykle zajímají jen malé hodnoty  $h$  – není realistické zkoumat efekt 2.řádu na týchž rozměrech jako má  $\mathfrak{R}$ .

V praxi není výpočet  $K(h)$  jednoduchý zvláště pro libovolné tvary  $\mathfrak{R}$ , kde je výpočet  $w_{ij}$  problematický. Přesný vzorec pro  $w_{ij}$  lze zapsat jen pro jednoduché tvary jako je pravoúhlá nebo kruhová oblast  $\mathfrak{R}$ , v jiných případech vyžaduje odvození  $w_{ij}$  plnou algoritmizaci úlohy (posloupnost řady kroků, nemožnost použít jeden vzorec) a je výpočetně náročný.

### 3.2.3.3 Vyhodnocení K-funkce

Pokud vypočteme  $K(h)$ , vykreslíme ho do grafu v závislosti na hodnotách  $h$ , abychom ocenili prostorové závislosti sledovaných událostí. Avšak na rozdíl od odhadované distribuce nejbližších vzdáleností, kde bylo možné alespoň částečně interpretovat graf, u K-funkce nejsme schopni interpretaci provést, protože nevíme, jak by měl vypadat graf v případě nulové prostorové závislosti (tj. plné náhodnosti výskytu bodových vzorků).



Obr. 3-21. Ukázka K funkce

Musíme proto transformovat K funkci na jinou funkci, kterou lze již hodnotit. Odvozuje se  $K_{(h)}$  pro náhodný proces (tedy výskyt událostí v libovolném bodě v oblasti  $\mathfrak{R}$  je nezávislý od všech dalších událostí a je stejný pro celou oblast  $\mathfrak{R}$ ). Tedy pro náhodný proces bude očekávaný počet událostí do vzdálenosti  $h$  roven  $n = \lambda * R = \lambda \pi r^2 = \lambda \pi h^2$ . To znamená, že podle definice K-funkce bychom měli očekávat  $K_{(h)} = \pi h^2$  pro homogenní proces bez prostorových závislostí.

Pravidelný vzorek bude mít  $K_{(h)}$  (tedy počet událostí) menší než je očekávaných  $\pi h^2$ , vzorek se shluky bude mít  $K_{(h)}$  větší než  $\pi h^2$ . Proto běžně srovnáváme odhad K-funkce s hodnotou  $\pi h^2$ . Běžně se vykresluje graf závislosti hodnoty  $L(h)$  na  $h$ , kde:

$$L(h) = \sqrt{\frac{k_{(h)}}{\pi}} - h$$

V tomto grafu vrcholy s kladnou hodnotou indikují shlukování („přitažlivost“) událostí a poklesy se zápornou hodnotou indikují pravidelnost. Souhrnně se tedy v grafu ukazuje pulsace událostí pro odpovídající hodnotu  $h$ .

Alternativa k použití odmocninové transformace je využití logaritmické transformace a vykreslování  $l(h)$  proti  $h$ , kde:

$$\hat{l}_{(h)} = \frac{1}{2} \log\left(\frac{\hat{k}_{(h)}}{\pi}\right) - \log h$$

Opět vrcholy indikují shlukování a poklesy pravidelnost pro určité hodnoty  $h$ .

Řešení lze zjednodušit vykreslováním hodnoty " $K(h) - \pi h^2$ " proti  $h$ .

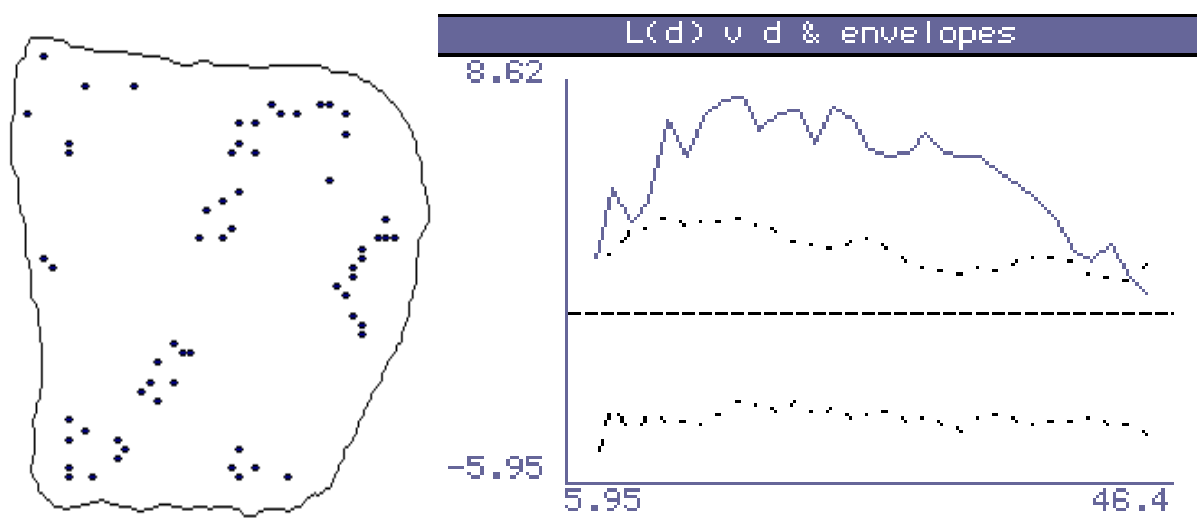
Formální ocenění významnosti těchto vrcholů/poklesů vyžaduje znalost distribuce  $L(h)$  a tedy  $K(h)$  pro CSR. Ta je neznámá a problematicky zjistitelná díky hraničním korekcím zabudovaným (obsaženým) v  $K(h)$ . Je však možné použít analogický přístup jako pro nejbližší vzdálenosti  $h$  získáním simulačního odhadu pozorované distribuce. Je potřebné zkonstruovat horní a dolní simulační obálky (limity):

$$U(h) = \max_{i=1, \dots, m} \{L_i(h)\}$$

$$L(h) = \min_{i=1, \dots, m} \{L_i(h)\}$$

z  $m$  nezávislých simulací  $n$  událostí v  $\mathcal{R}$  při CSR, při kterých je pozorována funkce  $L_i(h)$ .

Tyto limity  $U(h)$  a  $L(h)$  jsou potom vykreslovány v grafu aktuálně pozorované  $L(h)$  proti  $h$  (obr. 3-22).





Obr. 3-22. Shluková textura bodových událostí a odpovídající L funkce, probíhající nad horním limitem

Významnost vrcholů a poklesů je oceněna na základě vztahu:

$$\Pr(\hat{L}(h) > U(h)) = \Pr(\hat{L}(h) < L(h)) = \frac{1}{m+1}$$

Stejně jako u nejbližších vzdáleností je význam  $m$  dán počtem simulací nezbytných pro dosažení potřebné hladiny významnosti.

Jako metoda sumarizace a průzkumu dat má K-funkce několik výhod: poskytuje informace o prostorovém vzorku v různých měřítkách, zahrnuje používání přesné lokalizace událostí a využívá všechny vzdálenosti mezi událostmi, nejenom nejbližší vzdálenost. Navíc je znám teoretický tvar funkce  $K_{(h)}$  pro různé prostorové bodové modely.  $K(h)$  tedy nemusí být použit jenom k průzkumu prostorových závislostí, ale také k návrhu modelů vhodných k reprezentaci a odhadu parametrů modelů.

ESRI implementuje Ripley K funkci ve Spatial Analyst.

Přitom srovnáváme hodnotu vypočtené L funkce pro pozorování a L funkce očekávané, získané na základě vztahů (Kukuliač 2011).

Vypočtená K funkce:

$$K_h = \frac{A}{n \cdot (n-1)} \cdot \sum_{i=1}^n \sum_{j=1, i \neq j}^n k(i, j) \quad k_{i,j} = \begin{cases} 0 & d_{i,j} > d \\ 1 & d_{i,j} \leq d \end{cases}$$

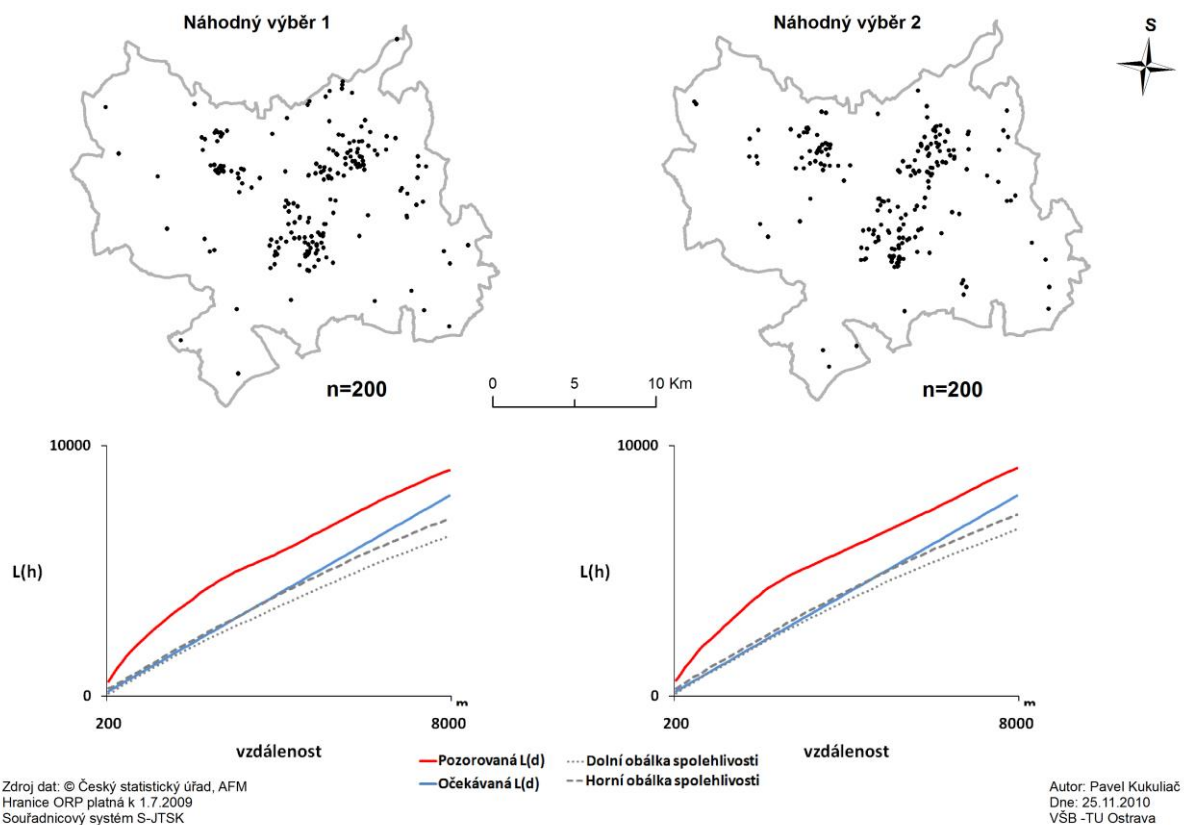
kde  $d_{i,j}$  je vzdálenost mezi  $i$ -tou a  $j$ -tou událostí v R. Je zřejmé, že jediný rozdíl je ve jmenovateli zlomku, kde se dělí  $n \cdot (n-1)$  místo  $n^2$ . Vzhledem k počtu párů při kombinaci bodů je vhodnější  $n \cdot (n-1)$ , i když praktické rozdíly jsou zanedbatelné.

Funkce L se vyjádří jako:

$$L_h = \sqrt{\frac{K_h}{\pi}}$$

Z toho vyplývá, že je rostoucí a že v případě náhodné distribuce se rovná vzdálenosti  $h$  (a tedy stoupá v grafu pod úhlem 45 st.).

## TESTOVÁNÍ K-FUNKCE PRO NÁHODNÝ VÝBĚR ADRESNÝCH BODŮ PODNIKŮ NA ÚZEMÍ ORP OSTRAVA PRO ROK 2009



*Obr. 3-23. Shlukování podniků na území Ostravy (L funkce probíhá výrazně nad pásmem spolehlivosti) pro 2 náhodné výběry*

### 3.3 Modelování prostorového uspořádání událostí (bodů)

Mezi základní úlohy v rámci interferenčních metod patří zkoumání náhodnosti pozorované distribuce bodového vzorku, případně přiřazení k shlukové či k pravidelné distribuci.

Nejčastěji je pozorovaná distribuce srovnávána s modelem úplné prostorové náhodnosti CSR (kapitola 3.3.1) (*complete spatial randomness*).

Často se ale vyskytují situace, kdy CSR není vhodným modelem. Při průzkumové analýze je diskutován typ struktury, který je přítomen v datech a obecně indikuje shlukování nebo pravidelnost, které mohou být zřejmé a priori díky povaze dat. Pokud chceme „vysvětlit“ určitou povahu pravidelnosti nebo shlukování, potřebujeme pro to formální modely, které jsou vztažené k jiným modelům než CSR

CSR používá model založený na homogenním procesu Poissonovy distribuce ve studované oblasti  $\mathfrak{R}$ . Pro model jsou podstatné 2 aspekty

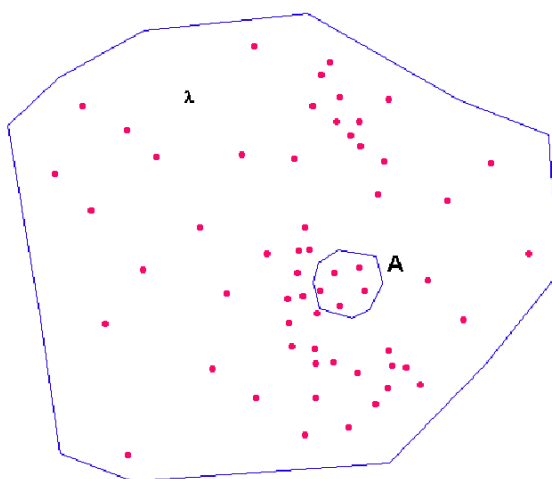
1. Je homogenní, což znamená, že  $\lambda$  je konstantní v  $\mathfrak{R}$ , není zde přítom žádný efekt 1.řádu (např. trend)
2. Má Poissonovu distribuci, což znamená prostorovou nezávislost ve výskytu událostí. Počet událostí, který se vyskytuje ve dvou sousedících regionech, není korelován, není zde přítom žádný efekt 2.řádu.

Pro obecný prostorový proces víme, že je docela běžné, že  $\lambda$  se mění v  $\mathfrak{R}$  (heterogenita) a/nebo že počet výskytů událostí v sousedních regionech je korelován (prostorová závislost).

Pro simulaci shlukové distribuce lze použít heterogenní Poissonův proces (kapitola 3.3.2), Coxův proces (kapitola 3.3.3), Poissonův shlukovací proces (kapitola 3.3.4), mohou být ale použity i jednoduché inhibiční procesy pro modelování pravidelného bodového pole. Markovův bodový proces (kapitola 3.3.5) může modelovat jak shlukování tak pravidelnost, od pravidelnosti v malém měřítku až po shlukování ve velkém.

#### 3.3.1 CSR - homogenní Poissonův proces

Standardní model pro CSR předpokládá, že události odpovídají homogennímu Poissonovu procesu v celé studované oblasti.



Obr. 3-24. Distribuce událostí a podoblast A

$\mathcal{A}$  je podoblast  $\mathcal{R}$ ,  $Y(\mathcal{A})$  je počet událostí v oblasti  $\mathcal{A}$ . Je-li prostorový bodový proces považován za sadu náhodných proměnných  $\{Y(\mathcal{A}), \mathcal{A} \in \mathcal{R}\}$ , pak to znamená, že  $Y(\mathcal{A}_i)$  a  $Y(\mathcal{A}_j)$  jsou nezávislé pro jakýkoliv výběr  $\mathcal{A}_i$  a  $\mathcal{A}_j$  a dále, že distribuce pravděpodobnosti  $Y(\mathcal{A})$  je Poissonovou distribucí

$$f_{Y(\mathcal{A})}(y) = \frac{(\lambda \cdot \mathcal{A})^y}{y!} e^{-\lambda \mathcal{A}}$$

$\mathcal{A}$  je plocha  $\mathcal{A}$

$\lambda$  je konstanta, intenzita neboli průměrný počet událostí na jednotku plochy.

Průměrná hodnota je  $\lambda \cdot \mathcal{A}$ .

Z toho vyplývá, že  $n$  událostí je nezávisle a jednotně distribuováno v  $\mathcal{R}$ . Jinak řečeno – libovolná událost má stejnou pravděpodobnost výskytu v libovolném místě  $\mathcal{R}$  a umístění události nezávisí na poloze jiné události, tedy mezi událostmi nejsou interakce.

Můžeme simulovat  $n$  událostí z takového procesu uzavřením  $\mathcal{R}$  do obdélníku  $\{(x,y): x_1 \leq x \leq x_2; y_1 \leq y \leq y_2\}$ , konkrétně generováním událostí se souřadnicí  $x$  z jednoduché distribuce na intervalu  $(x_1, x_2)$  a  $y$  souřadnicí z jednoduché distribuce na intervalu  $(y_1, y_2)$  a odmítnutím těch událostí, které neleží v  $\mathcal{R}$ . Události jsou generovány až do dosažení požadovaného počtu.

Základní testovanou hypotézou je srovnání pozorovaného vzorku s pravidelnou, shlukovou a náhodnou distribucí.

### 3.3.2 Heterogenní Poissonův proces

Heterogenní Poissonův proces je asi nejjednodušší alternativou k CSR. Konstantní intenzita  $\lambda$  pro CSR je nahrazena proměnnou funkcí intenzity  $\lambda(s)$ , ale výskyt každé události zůstává nezávislý na jiné. Výsledný proces je jednoduchým typem nestacionárního bodového pole s projevem vlivu pouze efektu 1.řádu. Nejjednodušší cestou k simulaci procesu je simulovat CSR na  $\mathcal{R}$  s intenzitou  $\lambda_{\max} = \max(\lambda(s))$  a pak vypustit v oblasti  $S$  tolik bodů, aby zbylý počet v místě odpovídal hodnotě  $\lambda(s)$ , tedy nezávisle ponechat události v oblasti  $S$  s pravděpodobností  $\lambda(s)/\lambda_{\max}$ .

### 3.3.3 Coxův proces

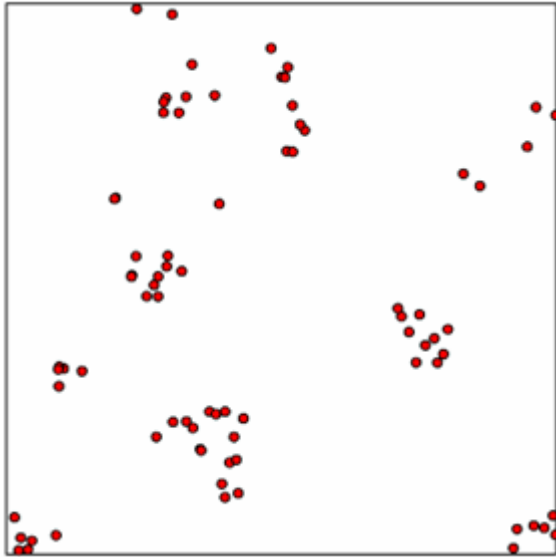
Coxův proces je přirozeným rozšířením heterogenního Poissonova procesu – intenzita  $\lambda(s)$  kolísá náhodně a ne deterministicky. Takový proces je často popisován jako dvojnásobně náhodný (stochastický).  $\lambda(s)$  je odvozena z pravděpodobnostní funkce přes  $\mathcal{R}$ , a pak podmíněně na hodnotě  $\lambda(s)$  události tvoří heterogenní Poissonův proces s intenzitou  $\lambda(s)$ . Výsledný proces může být stacionární i nestacionární. Bude stacionární jen tehdy, pokud distribuce pravděpodobnosti, ze které je intenzita generována, je stacionární. V principu může být takový proces simulován nejdříve simulací distribuce pravděpodobnosti pro  $\lambda(s)$  přes  $\mathcal{R}$  (tak se stanoví maximální počet  $\lambda(s)$  v daném místě) a pak použít vylučovací techniku popsanou u heterogenního Poissonovu procesu.

### 3.3.4 Poissonův shlukovací proces

Poissonův shlukovací proces vzniká z explicitního zahrnutí prostorového shlukovacího mechanismu přímo do modelu. Rodičovské události tvoří CSR proces a každý rodič produkuje náhodný počet potomků realizovaný nezávisle a stejně pro každého rodiče podle distribuce pravděpodobnosti  $f(\cdot)$ . Nakonec jsou pozice potomků k rodičům nezávislé a stejně distribuované podle určité bivariační (2 rozměrné) hustotní funkce  $g(\cdot)$ . Konečný proces zahrnuje jen potomky. Tento proces je stacionární, je také izotropní, pokud je funkce  $g(\cdot)$  radiálně symetrická. Metoda simulace takového procesu vychází přímo z jeho definice. Nejdříve je simulován rodičovský proces, aby byla získána místa umístění rodičovských událostí. Pak je pro každého rodiče nezávisle simulován počet

potomků podle funkce  $f(\cdot)$ . Nakonec je každý z nich umístěn kolem jeho rodiče podle funkce  $g(\cdot)$ . Potomci potom tvoří realizaci procesu. Abychom se vyhnuli hraničním problémům, musí být rodiče simulováni pro větší oblast než je  $\mathfrak{R}$ , aby nebyly ztraceny potomci, kteří leží v oblasti  $\mathfrak{R}$  a jejichž rodiče leží vně  $\mathfrak{R}$ .

Simulování za uvedených podmínek někteří autoři označují jako Thomasův proces (obr. 3-25). V případě použití rovnoměrné distribuce pro rozmístění potomků kolem rodičů jde o Matérnův proces. Obecnější modelem je Neyman-Scottův proces (Smith et al. 2011, 265).



Obr. 3-25. Thomasův PSP proces s parametry 20 shluků, směrodatná odchylka 0,03, průměr 5 (Smith et al, 2011, 265)

### 3.3.5 Markovovy bodové procesy

Markovovy bodové procesy umožňují vytvořit více obecnou skupinu pravděpodobnostních modelů pro bodové procesy. Zvláště někteří členové této rodiny poskytují flexibilnější rámec pro modelování pravidelnosti než je jednoduchý inhibiční proces. Takové modely např. dovolují realizovat případ, který může být nepravděpodobný, ale ne nemožný - 2 události se vyskytnou v těsné blízkosti a přitom jde o pravidelný vzorek. Markovovy procesy jsou teoreticky poněkud komplexní modely. Mezi jednodušší příklady patří třídy „párových interakčních procesů“ (*pairwise interaction processes*), například Strausův proces (Strauss). Zde jsou body v  $\mathfrak{R}$  považovány za sousedy, pokud jsou vzdálené méně než  $\delta$ . Spojená hustotní funkce pro  $\mathbf{n}$  bodových umístění  $(s_1, s_2, \dots, s_n)$  v  $\mathfrak{R}$ , která obsahuje  $\mathbf{m}$  různých párů sousedů, je určena jako:

$$f(s_1, \dots, s_n) = \alpha \beta^n \gamma^m \quad \beta > 0 \quad 0 \leq \gamma \leq 1$$

kde  $\alpha$  je normalizační konstanta,  $\beta$  sleduje intenzitu procesu a  $\gamma$  popisuje interakce mezi sousedy. Například při  $\gamma = 1$  získáváme CSR proces; při  $\gamma = 0$  získáváme hard core inhibiční proces s mezibodovou vzdáleností  $\delta$ ; střední hodnoty  $\gamma$  reprezentují formu mírné inhibice.

Obecná simulace Markovových procesů může být dosti komplexní a výpočetně náročná a nebudeme zde popisovat její detaily.

### 3.3.6 Srovnání jiných modelů než CSR s prostorovým bodovým procesem

Určité Poissonovy shlukovací procesy a určité Coxovy procesy mohou být statisticky nerozlišitelné. Bailey, Gatrell (1995) uvádí, že pokud je distribuce potomků  $f(\cdot)$  opět Poissonovou distribucí, je možné dokázat, že se vždy najde distribuce pravděpodobnosti pro  $\lambda(s)$  v Coxově procesu, která bude produkovat úplně stejný efekt jako ve shlukovacím procesu pro libovolné  $g(\cdot)$ . To je zajímavý a jistým způsobem nešťastný výsledek. V takových případech nemáme žádnou metodu statistické analýzy, abychom rozlišili oba procesy (data vzniklá z obou procesů), ačkoliv interpretace těchto modelů je zásadně odlišná (první je projevem heterogenity 1.řádu a prostorové nezávislosti, druhý reprezentuje stacionární efekt 2.řádu). Tak například při studiu výskytu vzácné nemoci nemůžeme být nakonec schopni najít správnou příčinu - jde o důsledek vystavení různému riziku (např. v důsledku faktorů životního prostředí) nebo se potvrzuje infekční charakter šíření choroby v důsledku existence kontaktů mezi jednotlivými případy.

Při výběru možného modelu můžeme použít průzkumné metody. Hlavním problémem je určit stupeň variability pozorované intenzity událostí v  $\mathfrak{R}$  a do jaké míry je dán stupněm variability pozorované intenzity událostí v  $\mathfrak{R}$  a do jaké míry je to způsobeno heterogenitou procesu nebo prostorovou závislostí 2.řádu. Možnou cestou odhadu intenzity v  $\mathfrak{R}$  je použití nějakého typu jádrových metod nebo redukovanou mírou 2.momentu nebo K-funkce, které sledují prostorové závislosti mezi regiony za předpokladu, že proces je izotropní. Informace získané z těchto průzkumných analýz mohou být použity k navržení vhodného typu modelu pro bodový proces - např. heterogenní Poissonův proces nebo Coxův proces.

Je nutné řešit situaci, jak odhadnout neznámé parametry v takovém modelu a jak formálně ocenit shodu modelu s daty. To je obtížná oblast a jedna z těch, kde zvolený přístup bude záviset na volbě konkrétního modelu. Např. v případě čistého efektu 1.řádu, jako je tomu u heterogenního Poissonova procesu, se můžeme pokusit použít jádrové metody k odhadu parametrů jednoduchých trendových modelů pro variace intenzity 1.řádu.

V jiných případech, kde model zahrnuje izotropickou prostorovou závislost, můžeme zkusit srovnat teoretickou K-funkci různých izotropických bodových procesů s odhadovaným pro pozorované pole. Takový přístup může být použit pro odhad parametrů takového modelu. Předpokládejme například, že náš navrhovaný model zahrnuje vektor parametrů  $\theta$ . Tedy teoretická K-funkce pro tento model by měla rovněž zahrnovat  $\theta$  a můžeme psát  $K(h, \theta)$ . Rozumnou cestou k odhadu  $\theta$  může být vybrat hodnotu  $\theta^*$ , která minimalizuje výraz:

$$\int_0^{h_0} ([k'(h)]^c - [k(h, \theta)]^c) dh$$

kde  $K'(h)$  je odhadovaná K-funkce pro pozorované bodové pole a  $h_0$  a  $c$  jsou „vytlačovací“ konstanty vybrané k získání požadovaných odhadovaných vlastností.

Jako ilustraci této myšlenky uvažujeme o aplikaci určitého typu Poissonova shlukovacího procesu na pozorované bodové pole. Předpokládejme, že vezmeme určitý příklad Poissonova shlukovacího procesu s Poissonovým počtem potomků na 1 rodiče, kde  $g(\cdot)$  je radiálně symetrická normální distribuce s rozptylem  $\sigma^2$ . Pak:

$$k(h) = \pi h^2 + \frac{(1 - e^{-h^2/4\sigma^2})}{\rho}$$

kde  $\rho$  je intenzita CSR procesu rodičů.

zde  $\theta = (\rho, \delta)^T$  a může být odhadnut jako bylo výše uvedeno.

Je jasné, že tyto přístupy nejsou zrovna jednoduché a mohou zahrnovat jak komplexní teoretické úvahy, tak intenzivní výpočty. Často jsou technicky založeny na simulacích. Nemusíme vědět, např. jaká je teoretická K-funkce nebo jaká relevantní vlastnost by měla odpovídat realizaci určitého modelu v naší oblasti  $\mathfrak{R}$ . Často je to pro obtíže s hraničním efektem. V tomto případě bychom mohli opakovaně simulovat navrhovaný model v  $\mathfrak{R}$  a odhadnout zájmovou vlastnost z těchto simulací. To může být potom srovnáno s odpovídajícím měřením odvozeným z pozorovaného bodového pole a významnost odchylek oceněna pomocí intervalu spolehlivosti určeného ze simulací.

### 3.4 Transformace bodové textury do kontinuálního pole

Transformace bodové textury do kontinuálního pole slouží k převodu dat a reprezentaci jevu z bodové reprezentace do kontinuální. Je využívána např. pro vizualizaci dat původně lokalizovaných v bodech pomocí 2,5D povrchů.

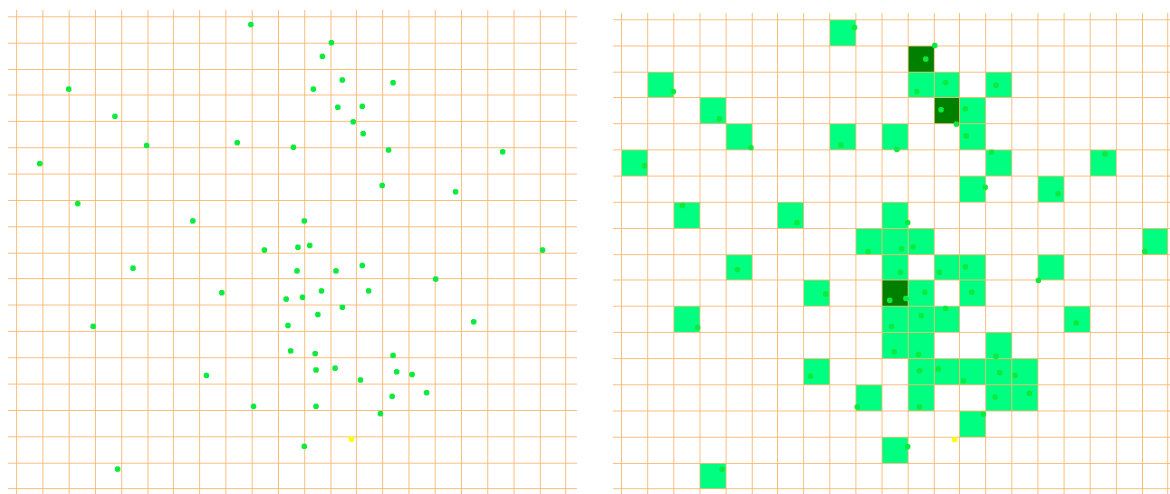
Používáme metody:

- Kvadrantová metoda
- Jádrové vyhlazení

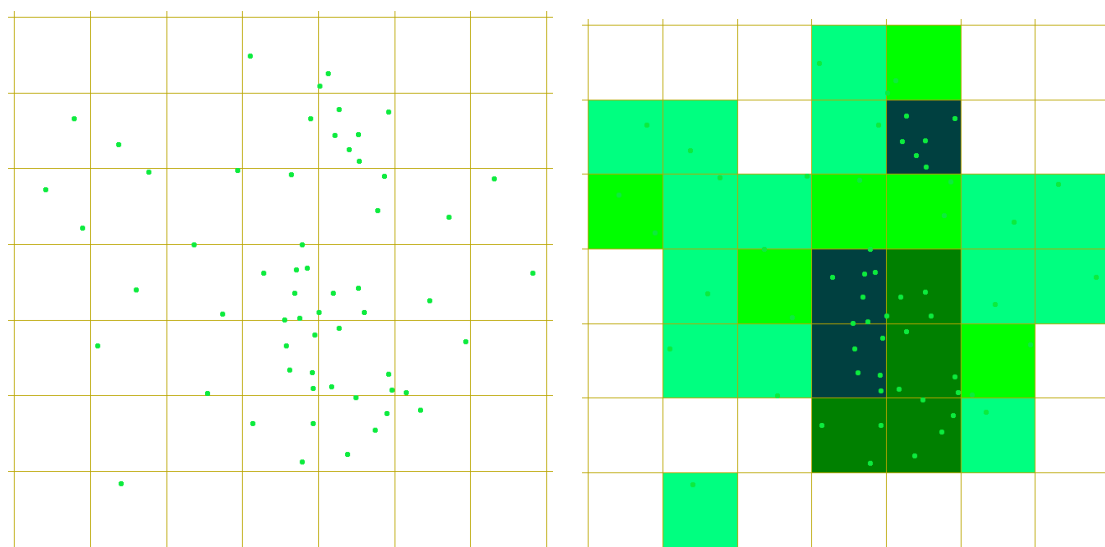
#### 3.4.1 Kvadrantová metoda

Základem kvadrantových metod je sledování **četnosti událostí** ve vymezených buňkách (kvadrantech). Pro transformaci se používá pravidelná mřížka. Počet událostí v buňce tedy udává hodnotu kontinuálního povrchu v daném místě.

Při kvadrantové analýze je citlivou volba sítě (při pravidelné mřížce), tj. stanovení počátku a velikosti buněk. Obecným nedostatkem kvadrantových metod je, že se nebere ohled na relativní polohu kvadrantů a relativní polohu událostí v kvadrantu.



Obr. 3-26. Příliš jemná síť, kde není vidět kontinuální změna intenzity událostí

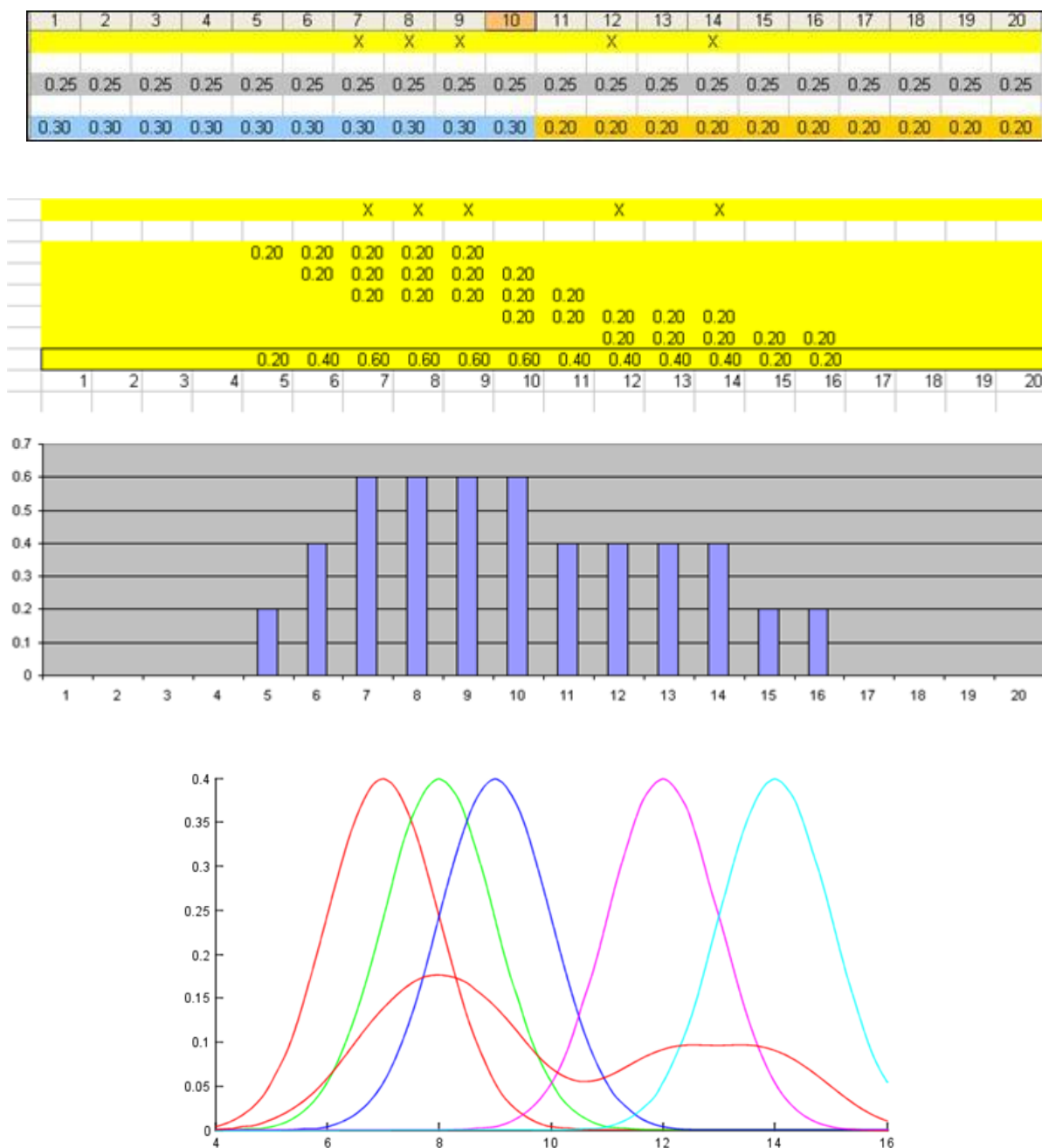


Obr. 3-27. Příliš hrubá síť s hrubým vymezením anomálních míst

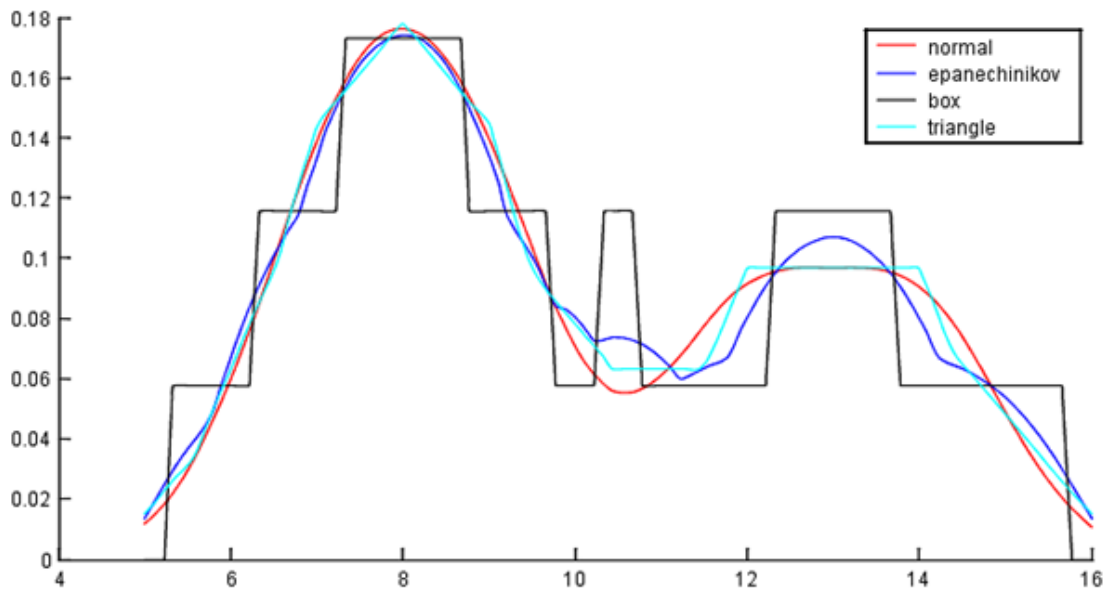


### 3.4.2 Jádrové vyhlazení

Jádrový odhad byl původně vyvinut pro získání vyhlazeného odhadu monovariační nebo multivariační hustoty pravděpodobnosti (křivky četnosti) získaného vzorku pozorování, tedy k vyhlazení histogramu.



Obr. 3-28. Postup odvození jádrového odhadu v případě jednorozměrných dat (histogram). Směrodatná odchylka je 2. Sečtou se plochy pod křivkami a podělí 5.



Obr. 3-29. Čtyři varianty jádrového odhadu podle typu jádrové funkce pro histogram.

Jedná se o neparametrickou metodu, protože neurčuje tvar funkční závislosti regresního vztahu, podobně jako klouzavé aritmetické průměry, oproti nim však představuje jisté zobecnění.

Odhad intenzity prostorového bodového vzorku je velmi podobný odhadu dvojrozměrné hustoty pravděpodobnosti, a proto dvojrozměrný jádrový odhad může být snadno upraven k odhadu intenzity. Jestliže  $\mathbf{S}$  reprezentuje obecně místo v  $\mathcal{R}$ , a  $\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_n$  místa  $n$  pozorovaných událostí, potom intenzita v bodě  $\mathbf{S}$  označená  $\lambda(\mathbf{s})$  může být odhadnuta jako:

$$\lambda'_\tau(\mathbf{s}) = \frac{1}{\delta_\tau(\mathbf{s})} \sum_{i=1}^n \frac{1}{\tau^2} k\left(\frac{\mathbf{s} - \mathbf{s}_i}{\tau}\right)$$

$k(\cdot)$  je vhodně vybraná funkce dvourozměrné hustoty pravděpodobnosti, známá jako kernel (jádro), která musí být symetrická kolem počátku. Parametr  $\tau > 0$  se označuje jako šířka pásma (bandwidth) a určuje stupeň vyhlazení – v podstatě je to poloměr kruhu se středem v  $\mathbf{S}$ , v kterém každý bod  $\mathbf{S}_i$  významně přispívá do  $\lambda'_\tau(\mathbf{s})$ .

Faktor

$$\delta_\tau(\mathbf{s}) = \int_{\mathcal{R}} \frac{1}{\tau^2} k\left(\frac{\mathbf{s} - \mathbf{u}}{\tau}\right) d\mathbf{u}$$

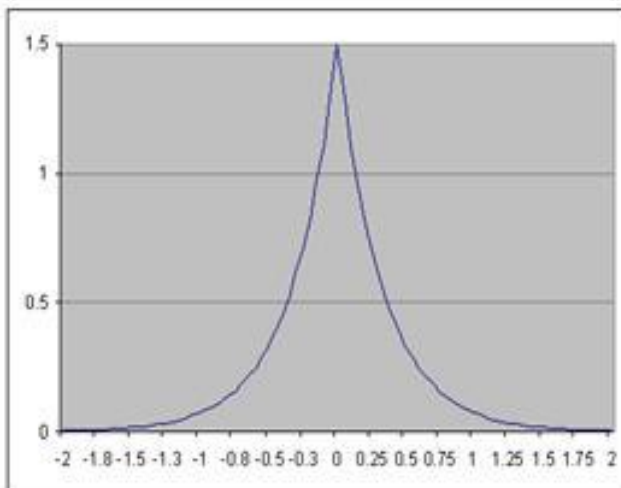
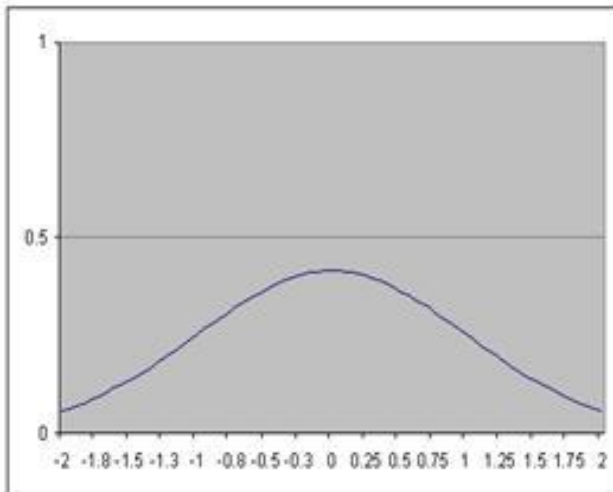
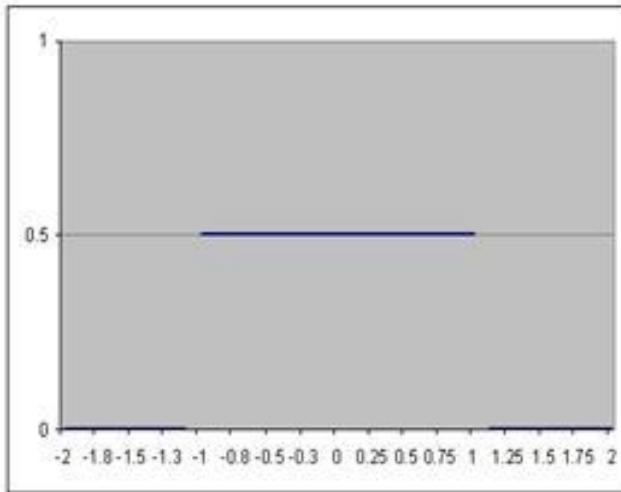
označuje okrajovou (hraniční) korekci – je to objem uzavřený pod kernelem se středem v  $\mathbf{S}$ , ležící uvnitř  $\mathcal{R}$ .

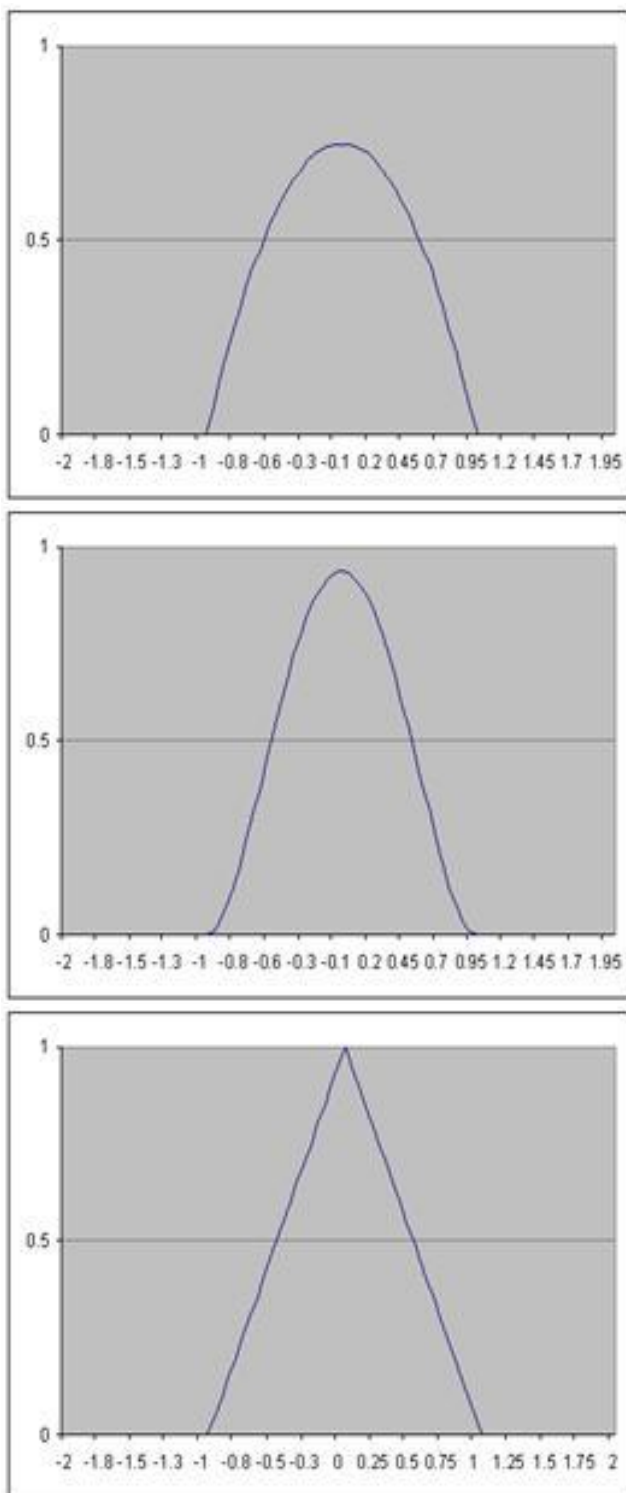
Hodnoty  $\lambda'_\tau(\mathbf{s})$  mohou být stanoveny pro každé místo ve vhodně vybrané jemné mřížce a mohou poskytovat užitečnou vizuální indikaci variability intenzity  $\lambda(\mathbf{s})$ .

Výběr vhodné funkce pro kernel  $K(\cdot)$  je relativně snadný, protože pro většinu vybraných funkcí požadovaných vlastností je jádrový odhad pro určitou šířku pásma velmi podobný. Typicky může být takovou funkcí kvartický (quartic) kernel. Pak, pokud zanedbáme hranovou korekci, dostaneme:

$$\lambda'_\tau(s) = \sum_{h_i \leq \tau} \frac{3}{\pi \cdot \tau^2} \left(1 - \frac{h_i^2}{\tau^2}\right)^2$$

$h_i$  vzdálenost mezi bodem  $\mathbf{S}$  a místem pozorované události  $\mathbf{S}_i$

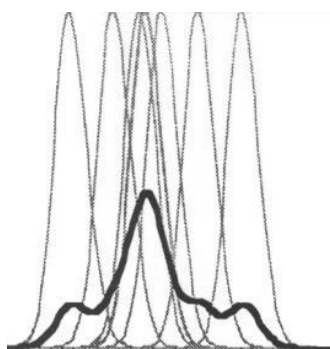




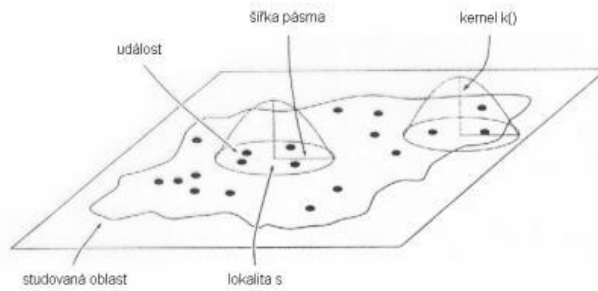
Obr. 3-30. Šest variant jádrové funkce - konstanta,  $N$ , exponenciální, kvadratický (Epanechikov, parabola), kvartický, trojúhelníkový (lineární) (manuál CrimeStat)

Kernel	Formula	Comments. Note $t=d_{ij}/h$ , $h$ is the bandwidth
Normal (or Gaussian)	$\frac{1}{2k} e^{-\frac{t^2}{2}}$	Unbounded, hence defined for all $t$ . The standard kernel in <a href="#">Crimestat</a> ; bandwidth $h$ is the standard deviation (and may be fixed or adaptive)
Quartic (spherical)	$\frac{3}{k} (1-t^2)^2, t \leq 1$ $= 0, t > 1$	Bounded. Approximates the Normal. $k$ is a constant
(Negative) Exponential	$Ae^{-k t },  t  \leq 1$ $= 0, t > 1$	Optionally bounded. $A$ is a constant (e.g. $A=3/2$ ) and $k$ is a parameter (e.g. $k=3$ ). Weights more heavily to the central point than other kernels
Triangular (conic)	$1- t ,  t  \leq 1$ $= 0, t > 1$	Bounded. Very simple linear decay with distance.
Uniform (flat)	$k,  t  \leq 1$ $= 0, t > 1$	Bounded. $k$ =a constant. No central weighting so function is like a uniform disk placed over each event point
Epanechnikov (paraboloid/quadratic)	$\frac{3}{4} (1-t^2),  t  \leq 1$ $= 0, t > 1$	Bounded; optimal smoothing function for some statistical applications; used as the smoothing function in the Geographical Analysis Machine ( <a href="#">GAM/K</a> ) and in ArcGIS

Sumace se provádí pouze pro  $h_i \leq \tau$ . Oblast vlivu, uvnitř které události přispívají do  $\lambda'_\tau(s)$  je tedy kruh o poloměru  $\tau$  kolem  $s$ . Na obr. 3-24 je nakreslen řez přes sadu radiálně symetrických funkcí  $K()$ , které odpovídají lokalizaci jednotlivých událostí. Funkce  $K()$  je posunuta počátkem do  $S$  a přepočítána faktorem  $\tau$  tak, aby poskytovala vhodné váhy pro události kolem  $S$ . V místě  $S$  (vzdálenost = 0) je váha nejvyšší ( $3/\pi\tau^2$ ) a klesá postupně až na nulu ve vzdálenosti  $\tau$ . Sumací příspěvků jednotlivých funkcí  $K()$  získáme výsledný odhad intenzity.



Obr. 3-31 Princip sumace příspěvků jednotlivých funkcí při tvorbě jádrového odhadu intenzity (Brunsdon 1995)



Obr. 3-32 Kernel jako 3D plovoucí funkce (podle Bailey, Gatrell 1995)

Pro lepší pochopení můžeme kernel považovat za 3D plovoucí funkci, která postupně navštíví každý bod  $S$  jemné mřížky (obr. 3-25). Vzdálenost ke každé pozorované události  $S_i$ , která leží uvnitř zóny vlivu (vzdálenost  $\tau$ ), je změřena a přispívá k výpočtu intenzity v místě  $S$ .

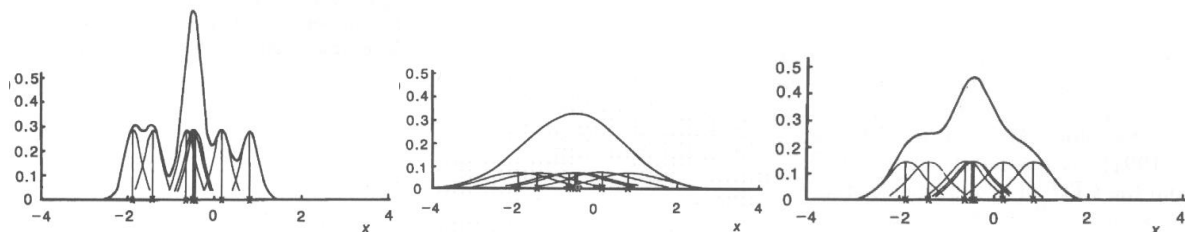
Nehledě na volbu funkce  $K(\cdot)$ , efekt zvětšení šířky pásma  $\tau$  vede k roztažení pásma kolem bodu  $S$ , ve kterém události ovlivňují odhad intenzity. Pro velké  $\tau$  je  $\lambda'_\tau(s)$  značně ploché a chybí lokální prvky. Pokud je  $\tau$  příliš malé, může se  $\lambda'_\tau(s)$  projevovat jako sada vrcholů s centry v místech událostí.

Vhodnou velikost pro  $\tau$  lze odhadnout jako:

$$\tau = 0,68 n^{-0,2} \cdot R$$

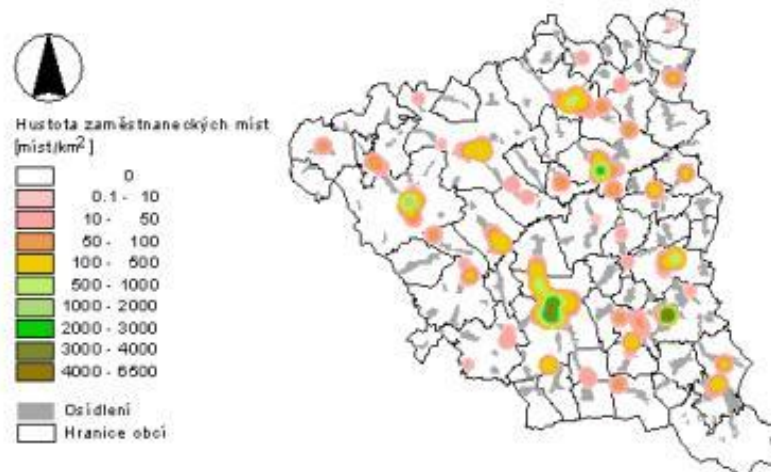
- R plocha oblasti  $\mathfrak{R}$
- n počet událostí v  $\mathfrak{R}$

V praxi se ale optimální hodnota  $\tau$  hledá zkoušením, zda výsledný obraz jádrového vyhlazení vyhovuje především z hlediska vhodného postizení variability pole.



Obr. 3-33 Příliš malá šířka pásma, příliš velká šířka pásma a optimální šířka pásma v řezu (citace)

Další informace především o praktické aplikaci, algoritmizaci a implementaci v ArcGIS Spatial Analyst najdete na <http://www.quantdec.com/SYSEN597/GTKAV/section9/density.htm>.



Obr. 3-34 Hustota zaměstnaneckých míst v okrese Nový Jičín k 31.3.2000 (jádrový odhad s šířkou pásma 1000 m)

### 3.4.2.1 Adaptivní jádrový odhad

Adaptivní jádrový odhad hustoty se liší tím, že parametr šířky pásma není konstantní, ale automaticky se mění zpravidla podle hustoty bodů (čím větší hustota bodů, tím menší šířka pásma, aby více vynikly lokální variace).

K odvození lokální distribuce bodů se používá technika pohyblivého jádrového odhadu hustoty (*moving kernel density estimation technique*). Při této analýze je umístěno okno na každý bod a jeho velikost se mění podle okolní hustoty bodů.

Pokud je  $\tau$  nahrazeno funkcí  $\tau(s_i)$ , pak platí (s vynecháním hraničního faktoru):

$$\lambda'_\tau(s) = \sum_{i=1}^n \frac{1}{\tau^2(s_i)} k\left(\frac{s-s_i}{\tau(s_i)}\right)$$

Potřebujeme specifikovat  $\tau(s_i)$ . Jednou z možností osvědčených v praxi je nejdříve provést neadaptivní jádrový odhad s jistou rozumnou hodnotou  $\tau_0$  a získat pilotní odhady  $\lambda'(s)$ . Pak je počítán geometrický průměr  $\lambda_g$  pilotních odhadů  $\lambda'(s_i)$  pro každé  $S_i$ . Pak

$$\tau(s_i) = \tau_0 \left( \frac{\tilde{\lambda}_g}{\tilde{\lambda}(s_i)} \right)^\alpha$$

$$0 \leq \alpha \leq 1 \quad \text{parametr citlivosti}$$

Pokud je  $\alpha = 0$ , jde o neadaptivní jádrový odhad (opět použijeme  $\tau_0$ ),  $\alpha = 1$  odpovídá maximálnímu lokálnímu přizpůsobení. Nejlepší výsledky podle Bailey, Gatrell (1995) dává hodnota  $\alpha = 0,5$ .

V modifikaci se používá i pro geograficky váženou regresi.

### 3.4.2.2 Problémy u jádrového odhadu

Jádrový odhad se zaměřuje na mapování intenzity bodové textury (zpravidla událostí), i když je možné ho uplatnit i pro jiné typy grafické reprezentace např. areály. Výskyt některých geografických jevů je však podmíněn distribucí jiného „základního“ jevu. Pokud budeme sledovat výskyt jednoho

druhu stromu v lese, není taková vazba (podmíněnost) zřejmá a výsledky můžeme interpretovat jako intenzitu výskytu tohoto druhu a zajímat se o příčiny pozorovaného uspořádání (textury). V řadě případů a zvláště u socioekonomických jevů tomu tak není. Např. pokud bychom sledovali výskyt bydlišť obyvatel s vyšším odborným vzděláním, získáme po vyhlazení vzor, který zřejmě bude především ukazovat distribuci obydlí, stěží budeme identifikovat podíl intenzity příslušející tomuto typu vzdělání. Ještě horší situace bude u jevů, které jsou podmíněny hierarchicky nebo paralelně více jevy.

Příkladem čisté hierarchické podmíněnosti může být např. počet uchazečů o zaměstnání v evidenci déle než 1 rok (tj. dlouhodobě nezaměstnaní). Zde se do výsledku postupně promítá distribuce bydlišť obyvatelstva, distribuce ekonomicky aktivních obyvatel mezi bydlicími, distribuce uchazečů mezi ekonomicky aktivními a distribuce dlouhodobě nezaměstnaných mezi uchazeči (a to ještě zanedbáváme faktor skryté, tj. neregistrované nezaměstnanosti). Je zjevné, že mapování intenzity výskytu dlouhodobě nezaměstnaných nelze samostatně správně interpretovat.

V případě jádrového vyhlazení bodové textury lze použít funkce relativní intenzity (Bithel 1996), která měří intenzitu v bodě relativně k průměrné intenzitě nezávislého jevu jako poměr obou funkcí (zpravidla sledovaná proměnná ku hodnotám pozadí). V některých případech se používá poměr logaritmů obou proměnných, případně rozdíl v intenzitách (bez nebo se standardizací), součet intenzit (bez nebo se standardizací). Používá se stejná šířka pásma. Problémy vznikají v případě, že se hodnota pozadí blíží nule.

V případě složitějších vazeb (více ovlivňujících jevů) se místo zahrnutí mnoha faktorů použijí kontrolní vzorky, které zabezpečí standardizaci (odvození kontrolních vzorků je vhodné realizovat stochastickou simulací).

Levine (2007, kap. 8) poskytuje vhodný přehled různých technik a jejich parametrů, včetně alternativních metod pro výběr šířky pásma, a uvádí příklady pro analýzu kriminality, zdravotnické analýzy, urbanistické a ekologické.



### 3.5 Analýza vícenásobných typů událostí

Zatím jsme předpokládali, že všechny sledované události jsou jednoho typu. Nyní předpokládejme, že máme 2 a více typů událostí.

Předchozí metody nám dovolily analyzovat události odděleně (pouze 1 typu) nebo všechny dohromady, ukázat shlukování nebo pravidelnost, avšak ne zodpovědět otázky typu:

*Je textura výskytu události 1 druhu ve vztahu k textuře výskytu události jiného typu?  
Vysvětluje distribuce 1 typu událostí distribuci druhého typu událostí?*

Obecně bychom měli hledat a ověřit existenci **nezávislosti** mezi typy událostí v protikladu k přitahování jednotlivých typů událostí nebo jejich odpuzování. Testování je tedy založeno na pozorované bivariační nebo multivariační bodové textuře, kde za základní situaci považujeme nezávislost mezi výskytem různých typů událostí.

Nezávislost znamená, že celková textura událostí je tvořena z nezávislých dílčích procesů, každý pro 1 typ události. Je nutné si uvědomit, že nezávislost neznamená, že by nějaký z dílčích procesů musel být typu homogenního Poissonova procesu (kapitola 3.1.1). Mohou být nezávislé např. i shlukovací proces a pravidelný proces. Obecně, každý z dílčích procesů může produkovat vzájemně odlišné textury nebo textury odlišné od celkové textury.

Při hledání teoretických modelů pro multivariační bodové procesy je možné generalizovat Poissonův, Poissonův shlukovací (kapitola 3.3.4) a Coxův procesy (kapitola 3.3.3) tak, aby byla zahrnuta myšlenka spojených procesů, které buď zobrazují přitahování nebo odpuzování – např. spojené Coxovy procesy nebo spojené párové procesy.

Testování nezávislosti vícenásobných typů událostí zpravidla využívá:

1.  $\chi^2$  testu pro kvadrantovou metodu (kapitola 3.5.1),
2. metodu nejbližších vzdáleností (kapitola 3.5.2),
3. K funkce.

#### 3.5.1 Kvadrantová metoda s $\chi^2$ testem

Nejjednodušší přístup k testování nezávislosti prostorové distribuce 2 typů událostí je založen na sčítání počtu událostí v buňkách ve sledované oblasti  $\mathcal{R}$ . Sčítání se provádí v buňkách buď náhodně rozptýlených nebo pravidelně umístěných v mřížce v případě, kdy máme k dispozici úplný bodový vzorek.

Běžně se prezentují výsledky v podobě tabulky 2 x 2 („přítomnost – nepřítomnost“), kde je zapsán počet kvadrantů  $c_{ij}$ , které obsahují buď oba typy událostí ( $c_{12}$  nebo  $c_{21}$ ) nebo pouze 1 typ ( $c_{11}$ ,  $c_{22}$ ).

Tab. 3-1 Sledování přítomnosti 2 typů událostí v buňkách

		Přítomnost událostí 2. typu v buňce	
		není	je
Přítomnost událostí 1. typu v buňce	není	$C_{11}$	$C_{12}$
	je	$C_{21}$	$C_{22}$

Výsledky mohou být testovány jednoduchým standardním  $\chi^2$  testem nezávislosti. Srovnáváme

$$\chi^2 = \frac{(c_{11}c_{22} - c_{12}c_{21})^2 \sum_i \sum_j c_{ij}}{\sum_j c_{1j} \sum_j c_{2j} \sum_i c_{i1} \sum_i c_{i2}}$$

s procentními body  $\chi^2$  distribuce. Překročení kritické hodnoty (zpravidla 0,05) znamená zamítnutí hypotézy nezávislosti mezi 2 texturami.

Problém u tohoto jednoduchého přístupu je stejný jako u jiných kvadrantových metod - velké množství informace o lokalizaci, které je obsaženo v pozorované textuře, je v testu ignorováno.

### 3.5.2 Metoda nejbližších vzdáleností

Distribuční funkce nejbližších vzdáleností typu událost-událost pro monovariační distribuci je nahrazena v multivariačním případě soustavou distribučních funkcí nejmenších vzdáleností  $G_{ij}(h)$ .  $G_{ij}(h)$  je pravděpodobnost, že vzdálenost náhodně vybrané události typu  $i$  k nejbližší události typu  $j$  je menší nebo rovna  $h$ .

$$G_{ij}(h) = \Pr(d_{ij} \leq h)$$

$F_j(h)$  je distribuční funkcí nejmenší vzdálenosti typu bod-událost a znamená pravděpodobnost toho, že vzdálenost náhodně vybraného bodu k nejbližší vzdálenosti typu  $j$  je menší nebo rovna  $h$ . Odhady těchto funkcí  $F'_j(h)$  a  $G'_{ij}(h)$  mohou být získány z pozorovaného bodového vzoru. V případě potřeby jsou do odhadů zahrnovány i okrajové korekce.

Pokud jsou jednotlivé (díleč) bodové procesy v multivariačním bodovém vzoru nezávislé, potom by měla být distribuce nejmenších vzdáleností k události typu  $j$  tatáž, ať už je počátkem náhodně vybraný bod nebo náhodně vybraná událost typu  $i$  ( $i \neq j$ ). Tedy:

$$G_{ij}(h) = F_j(h) \quad i \neq j$$

Vyhodnocení se provádí s pomocí grafu nebo se používá korelační analýza.

#### 3.5.2.1 Grafické hodnocení

Při **grafickém hodnocení** se zpravidla vykresluje  $F_j(h)$  a  $G_{ij}(h)$  do jednoho grafu, kde  $h$  se mění od 0 do maximální vzdálenosti nejbližších sousedů typu událost-událost nebo bod-událost. Graf je pak zkoumán, zda se objevuje nějaká odchylka mezi 2 odhadovanými distribucemi.

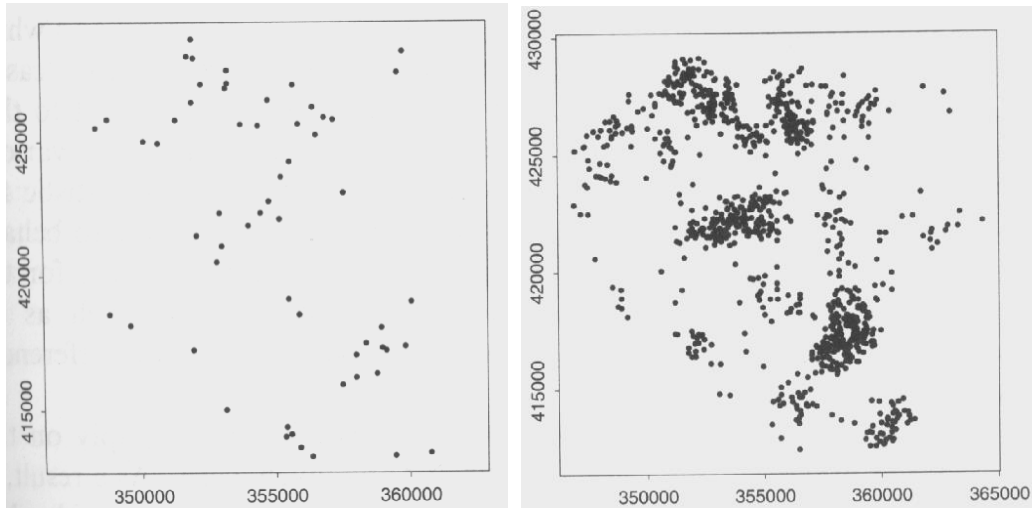
#### 3.5.2.2 Numerické hodnocení

Při **korelační analýze** se nemusí provádět plný odhad  $G_{ij}(h)$  a  $F_{ij}(h)$ . Nejjednodušší z hlediska implementace je použít náhodný vzorek bodů v  $\mathfrak{R}$  a měřit nejbližší vzdálenosti bod-událost typu  $i$  a nejbližší vzdálenosti bod-událost typu  $j$  pro tytéž body. Vzdálenosti v každém vzorku jsou pak nahrazeny jejich pořadím (ve variační řadě) v příslušném vzorku. Výsledkem je sada párů pořadí vzdáleností bod-událost. Nezávislost může být testována Spearmanovým nebo Kendallovým koeficientem korelace pro pořadí. Jednou z výhod takového testu je, že může být použit i na pouhý vzorek bodů stejně jako na zcela zmapovanou distribuci, další výhodou je nezávislost na typu distribuce. Doporučuje se zavést okrajovou korekci – zpravidla se vypouští ze zpracování obou vzorků vzdálenosti, kde je náhodně zvolený bod blíže k hranicím  $\mathfrak{R}$  než k nejbližší události jakéhokoliv typu (hraniční korekce typu adaptivní ochranné pásmo).

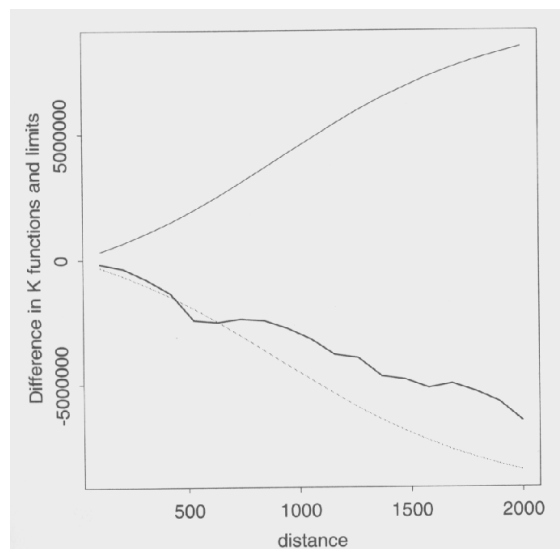
### 3.5.3 K funkce

Pro multivariační sledování se provádí úprava monovariačních K funkcí (kapitola 3.2.3) podobným způsobem jako u metody nejbližších vzdáleností.

Používá se  $K_{ij}(h)$  – charakteristika počtu událostí  $j$  do vzdálenosti  $h$  od události  $i$ .



Obr. 3-35 Výskyty rakoviny hrtanu a plic (Bailey, Gatrell 1995)



Obr. 3-36 Rozdíl v K funkcích obou případů rakoviny a simulační obálky (Bailey, Gatrell 1995)

## 3.6 Časoprostorové analýzy

Události/objekty jsou lokalizovány v bodech a čas je zaznamenán pouze 1 údajem.

### 3.6.1 Popisná statistika

Některé systémy také nabízejí tzv. analýzu korelované procházky (Correlated Walk Analysis) – např. Crimestat nebo Hawth nástroje pro ArcGIS. Cílem je předpověď pravděpodobného výskytu další události.

### 3.6.2 Inferenční statistika

K testování se používá několik jednoduchých testů:

1.  $\chi^2$  test (kapitola 3.6.3)
2. Knoxův test (kapitola 3.6.4)
3. Mantelův test (kapitola 3.6.5)
4. D funkce (kapitola 3.6.6)
5. STAM (prostor-čas-atributový stroj) (kapitola 3.6.7)

V některých případech se transformuje  $t_{ij}$  na časovou funkci  $w_{ij}$ , tedy funkci času mezi událostmi  $i$  a  $j$ .

Vhodným typem může být např.:

$$W_{ij} = \frac{1}{(a + |t_i - t_j|)} \quad \text{a je konstanta zajišťující, abychom nedělili 0 (např. a = 0.1)}$$

### 3.6.3 $\chi^2$ test

Zvolíme limity  $D$  a  $T$  pro vzdálenost a čas a rozdělíme události podle tabulky 3-2.

Tab. 3-2 Distribuce událostí v čase a prostoru

Čas	prostor	$d_{ij} < D$	$d_{ij} > D$
$t_{ij} \leq T$		Události, které se staly blízko u sebe a současně	daleko od sebe, současné
$t_{ij} > T$		Blízko, časově oddělené	daleko, časově oddělené

$t_{ij}$       rozdíl mezi časy  $t_i$  a  $t_j$  událostí  
 $d_{ij}$       vzdálenost mezi událostmi

K testování lze použít standardní  $\chi^2_1$  test nezávislosti v tabulce 2 x 2.

Celý časoprostor dat pak můžeme klasifikovat do 4 kategorií a můžeme najít shluky událostí.

Protože však uvedené rozdělení vyjadřuje vztahy pouze mezi 2 událostmi, je nutné v klasickém případě převést řešení na problém rozdělení do shluků.

### 3.6.4 Knoxův test

Tento test zavedl epidemiolog George Knox.

Z  $n$  událostí je možno získat  $n*(n-1)$  uspořádaných párů. Pro každý z párů měříme prostorovou vzdálenost a časový interval.

Spočítáme počet uspořádaných párů "blízkých v prostoru" (**S**).

Spočítáme počet uspořádaných párů "blízkých v čase" (**T**).

Čas a prostor je tedy rozdělen do dvou kategorií - blízké x neblízké.

Spočítáme počet uspořádaných párů současně blízkých v čase i prostoru, získáme **X**.

Pokud platí úvodní předpoklad nezávislosti prostoru a času, pak **X** odpovídá Poissonově distribuci s průměrným počtem událostí  $\lambda$

$$\lambda = \frac{S * T}{n * (n - 1)}$$

a vypočítáme pravděpodobnost jevu X:

$$p(x) = \frac{\lambda^x}{x!} e^{-\lambda}$$

Je-li p(x) menší než standardně používané limity pravděpodobnosti 0.05 resp. 0.01, nelze jev považovat za náhodný a musíme odmítnout nezávislost tvorby shluků v čase a prostoru.

Nevýhodou Knoxova testu je, že pracuje jen s logikou ano/ne (události jsou/nejsou blízké).

### 3.6.5 Mantelův test

Mantelův test nepoužívá Booleovskou logiku jako v případě  $\chi^2$  nebo Knoxova testu, ale sleduje skutečné velikosti intervalů (vzdálenosti v euklidovském i časovém prostoru). Používá se statistika (ukazatel)

$$\sum_{i \neq j} X_{ij} Y_{ij}$$

kde  $X_{ij}$  je vzdálenost události i a j;  
 $Y_{ij}$  je časový interval mezi událostmi i a j;

Ty lze vypočítat jako

$$X_{ij} = \frac{1}{(k_s + d_{ij})}$$

$$Y_{ij} = \frac{1}{(k_t + t_{ij})}$$

kde  $k_s$  a  $k_t$  jsou libovolné konstanty vybrané tak, aby všechny události měly identické prostorové a časové souřadnice

Pro výše uvedený statistický ukazatel lze vypočítat aritmetický průměr a rozptyl podle teoretické distribuce a ty se použijí k posouzení významu vypočteného ukazatele. Tak se prověří náhodnost distribuce událostí.

Mantelův test se používá často v ekologii ([http://en.wikipedia.org/wiki/Mantel\\_test](http://en.wikipedia.org/wiki/Mantel_test)), ovšem pro měření podobnosti hodnot. Používá se zpravidla ve standardizované formě (odečítá se průměr a dělí se odchylkou), jehož statistická pravděpodobnost se určuje permutačními testy (zaměňují se čísla v řádcích a sloupcích při velkém počtu iterací – aspoň 1000x). <http://www.nceas.ucsb.edu/scicomp/Dloads/SpatialAnalysisEcologists/SpatialEcologyMantelTest.pdf>

### 3.6.6 D funkce

Metoda využívá modifikace (rozšíření) K funkcí (kapitola 3.2.3). Funkce K(h,t) udává očekávaný počet událostí do vzdálenosti **h** a do časového intervalu **t**, standardizovaná počtem událostí na plochu a čas. Vycházíme z předpokladu, že pokud jsou události nezávislé, pak platí:

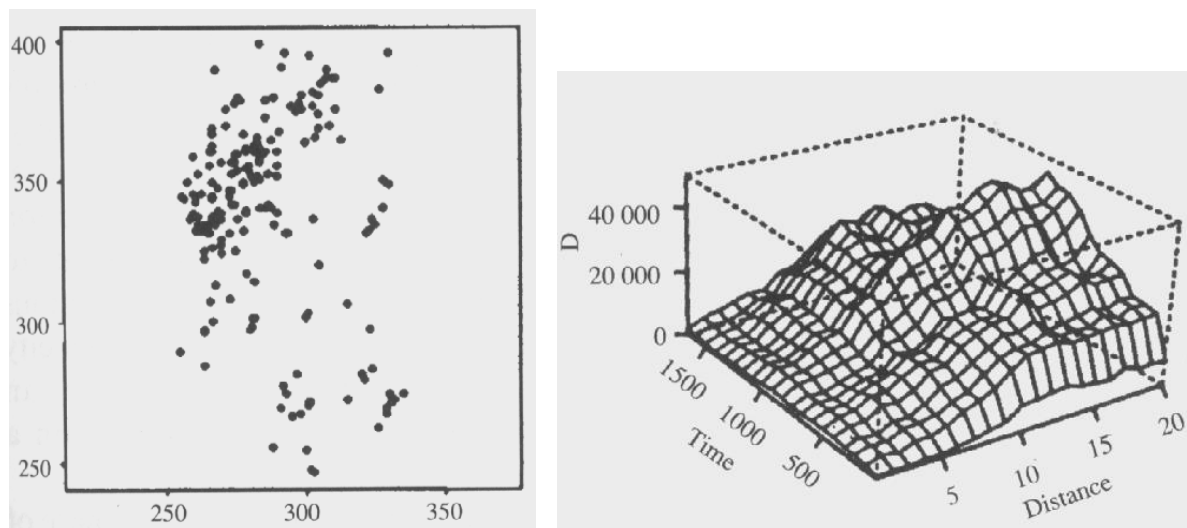
$$K(h,t) = K_S(h) * K_T(t)$$

$K_S(h)$  klasická K funkce pro pozorovaný vzorek událostí

$K_T(h)$  K funkce aplikovaná pro výskyt událostí v čase

Proto se sleduje rozdílová funkce  $D'(h,t) = K'(h,t) - K'_S(h) * K'_T(t)$

Výsledek se nejlépe interpretuje ve 3D grafu. Zvýšené hodnoty D funkce naznačují zvýšené interakce mezi událostmi. Např. na obr. 3-30 je nejvyšší odchylka při hodnotě vzdálenosti mezi událostmi 15 a časového rozdílu 1000. Validace je možná opět simulačními postupy.



Obr. 3-37 Výskyt případů Burkittova lymfomu a odpovídající  $D(h,t)$  (Bailey, Gatrell 1995)

### 3.6.7 Prostor-čas-atributový analytický stroj

Prostor-čas-atributový analytický stroj (*space-time-attribute analysis machine STAM*) je příkladem nasazení jednoduché analytické metody ve velkém měřítku. Prohledává všechny body v geografickém, časovém i atributovém „prostoru“.

Příkladem může být analýza výskytu nemoci, kde je cílem nalézt podobný vzorek na základě všech 3 „prostorů“.

Algoritmus:

- 1) definice rozsahu prohledání (stanovení limitů) ve všech 3 prostorech
- 2) výběr pozorování, která do limitů spadají
- 3) definuje se geografická prohledávaná oblast se středem v pozorovaném případě (místě) a s poloměrem  $g_r$  odvozeným z velikosti limitu
- 4) definuje se prohledávaná oblast v čase se středem v pozorovaném případě a poloměrem  $t_r$  odvozeným z velikosti limitu
- 5) definuje se atributová prohledávaná oblast se středem v pozorovaném případě (atributu) a s poloměrem  $a_r$  odvozeným z velikosti limitu
- 6) prohledá se databáze, aby se našly záznamy, které leží uvnitř prohledávaných oblastí (vymezených v krocích 3-5)
- 7) použije se testovací procedura Monte Carlo, která určí významnost pro získaný výsledek hledání
- 8) je-li náhodnost výsledku nízká (tedy nelze-li korelaci výskytu vysvětlit jako náhodu), uloží se identifikátor záznamu a prohledávací parametry
- 9) vyzkouší se všechny kombinace  $g$ ,  $t$  a  $a$  prohledávaných parametrů
- 10) tiskne se výsledek

11)změní se pozorovaný případ a opakují se body 3-11.

U atributů se měří shoda např. počtem shodných příznaků. Např. výskyty událostí jsou si podobné, jestliže se shodují alespoň v 5 atributech popisujících událost.

Geografický prostor se může např. prohledávat od 1 do 20 km s krokem 2 km.

## 4 Geostatistické metody pro kontinuální pole

**Prostorová (regionalizovaná) proměnná** představuje veličinu, která je funkcí polohy, tj. např. souřadnic  $X, Y, Z$ . Každému bodu ve vymezeném prostoru může být přiřazena hodnota veličiny, která je obecně složena ze systematické (strukturní) a náhodné části. Systematická složka je funkcí souřadnic  $U=f(X, Y, Z)$  a bývá předmětem popisu a interpretace. Náhodná složka představuje realizaci náhodných vlivů, šumu. K popisu prostorové proměnné lze využít teorie náhodných funkcí, kde pozorovanou hodnotu považujeme za jednu realizaci náhodné funkce v daném místě.

Přímé studium  $U(X, Y, Z)$  je prakticky vyloučené vzhledem k prostorové variabilitě a malému počtu známých realizací veličiny (tedy jejich měření či pozorování). Pro popis celkové struktury se někdy využívá aproximující funkce, která popisuje trend v poli.

Na základě funkční závislosti hodnoty sledované veličiny na poloze lze očekávat, že existuje i vzájemná závislost mezi jednotlivými hodnotami sledované veličiny. Ta se může projevit především při malých vzdálenostech mezi zkoumanými místy, kdy konstatujeme podobnost hodnot sledované veličiny. Jde tedy o korelaci hodnot téže veličiny a hovoříme o tzv. **autokorelaci**, která je závislá na vzdálenosti pozorování. Pokud závislost pozorujeme, hovoříme o tzv. kontinuálním poli. Pokud tato závislost nezávisí na směru zjišťování v poli, jedná se o izotropní pole, v opačném případě jde o anizotropní pole.

Předpoklady (řazeno od nejsilnějších k nejslabším požadavkům):

### **Stacionarita (striktní):**

- **distribuce veličiny je zcela nezávislá na místě, kde ji zkoumáme.**

### **Stacionarita 2.řádu:**

- **průměrná hodnota veličiny je konstantní v celé zkoumané oblasti**
- **hodnota kovariační funkce mezi 2 libovolnými místy závisí pouze na vzdálenosti a směru (vektor  $h$ )**

### **Intrinsikční hypotéza:**

- **Očekává se nulový rozdíl mezi hodnotami veličiny ve 2 místech oddělených vektorem  $h$**
- **Hodnota (semi)variogramu závisí pouze na vzdálenosti a směru (vektor  $h$ )**

**Geostatistická analýza** (někdy se konkretizuje jako strukturální analýza) se snaží popsat chování této prostorové proměnné. Jejím základním nástrojem jsou strukturální funkce (kapitola 4.1).

Prvním krokem při geostatistické analýze by mělo být ověření statistické distribuce sledované veličiny. Většina uváděných metod předpokládá normální distribuci známých hodnot a používá lineární metody. V případě jiného typu distribuce se provádí transformace hodnot (normalizace dat) nebo se používají nelineární, často neparametrické metody (např. indikátorové techniky).

Po geostatistické analýze se provádí zpravidla odhad hodnot v poli, tedy i v místech, která nebyla primárně měřena. K interpretaci se využívají geostacionární metody odhadu, které označujeme jako krigování (kapitola 4.2). Vedle této metody se stále více uplatňují po zkoumání lokálního odhadu simulační postupy (kapitola 4.2.12).

### 4.1 Strukturální funkce

Strukturální funkce slouží pro vizualizaci, modelování a průzkum prostorové autokorelace prostorově proměnné veličiny, tedy pro popis a analýzu variability proměnné ve studovaném poli.



Běžně používanou strukturální funkcí je **semivariogram**. Semivariogram vyjadřuje, jak se mění proměnná  $U$  mezi místem  $x$  a místem  $(x+h)$ , mezi nimiž je vzdálenost  $h$ . Jak název napovídá, jde o měření variability (změny 2.řádu):

$$\gamma(h) = \frac{1}{2n_h} \sum_{n_h} [u(x) - u(x+h)]^2$$

$n_h$  je počet párů při kroku  $h$

Při výpočtu semivariogramu je potřebné zvolit vhodný způsob výběru hodnot ve výpočtu. Praktické provádění se liší podle toho, zda zkoumáme hodnoty na linii (kapitola 4.1.1) nebo hodnoty v ploše (kapitola 4.1.5).

Jinými typy strukturálních funkcí jsou:

- kovariační (autokorelační) funkce

$$K(h) = \frac{1}{2n_h} \sum_{n_h} [u(x) - m(x)] * [u(x+h) - m(x+h)]$$

$m(x)$  je střední hodnota veličiny  $u$  v blízkosti místa  $x$

- korelogram (normovaná kovariační funkce)

$$K_N(h) = K(h) / K(0)$$

kde  $K(0)$  je kovariance při nulové vzdálenosti, zpravidla se rovná statistickému rozptylu  $s^2$

- relativní semivariogram - snaží se kompenzovat proporcionální efekt

$$\gamma_r(h) = \gamma(h) / m^2$$

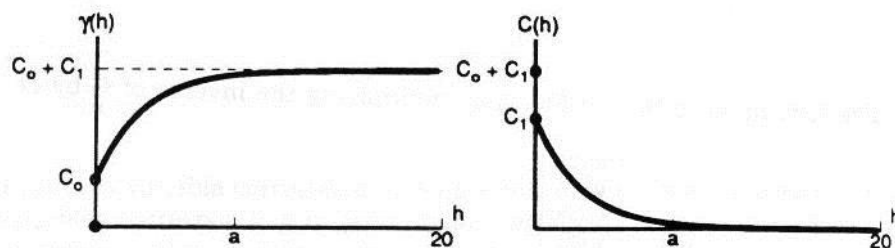
kde  $m^2$  je kvadrát střední hodnoty veličiny  $U$

- neergodický semivariogram - snaží se kompenzovat vliv trendu

$$\gamma_{ne}(h) = s^2 - K(h)$$

Vztah mezi kovariační funkcí a semivariogramem:

$$\gamma(h) = K(0) - K(h)$$



Obr. 4-1 Vztah semivariogramu (vlevo) a kovariační funkce (vpravo)

Existence **zbytkového rozptylu**  $c_0$  se projevuje tak, že model semivariogramu nevychází z počátku souřadných os, ale z výšky  $c_0$ . Je způsoben rozptylem v oblasti menší velikosti než je základní krok nebo nepřesností zjištěných hodnot (např. výsledky opakovaných měření na jednom místě se mohou lišit).

Při geostatistické (strukturální) analýze se provádí výpočet strukturálních funkcí, aplikace vhodného teoretického modelu (kapitola 4.1.1) a interpretace průběhu strukturální funkce, kde se studuje především charakter pole (kapitola 4.1.2).

Pro geostatistickou analýzu a následně i provádění lokálních či globálních odhadů si musíme být vědomi určitých omezujících předpokladů, na základě kterých budeme provádět výpočty. Zpravidla se u zkoumaného pole předpokládá, že průměrná hodnota proměnné je v poli konstantní (pokud se rozdělí celé pole na několik úseků a v nich se vypočítají průměrné hodnoty proměnné, neměli by být výsledky odlišné) a rovněž prostorová autokorelace je v poli konstantní (strukturální funkce, vypočtené odděleně v jednotlivých částech pole, budou shodné). Není-li možné výše uvedené předpoklady přijmout, musíme postupovat podle hypotézy univerzálního krigování, která předpokládá, že změnu průměrné hodnoty proměnné v poli lze popsat určitým trendem a k výpočtu strukturální funkce použijeme namísto hodnot proměnné rozdíly mezi hodnotou a očekávanou hodnotou proměnné.

Vývoj střední hodnoty v poli je možné popsat **trendem**. Trend (drift) se vyjadřuje lineární kombinací funkcí.

$$D(x) = \sum_{p=0}^q a_p * f^p(x)$$

např.  $D(x) = a_1X_1 + a_2X_2 + a_3X_1^2 + a_4X_1X_2 + a_5X_2^2$

V případě existence významného trendu v poli je potřebné používat místo původních strukturálních funkcí strukturální funkce reziduí.

Semivariogram reziduí:

$$\gamma(v_1, v_2) = \frac{1}{2n} (\sum [u(v_1) - u(v_2)]^2) - [m(v_1) - m(v_2)]^2 = \frac{1}{2n} (\sum [r(v_1) - r(v_2)]^2)$$

$v_1, v_2$  jsou 2 body, mezi nimiž se zjišťuje hodnota semivariogramu

$u(v_1)$  zjištěná hodnota v bodě  $v_1$

$m(v_1)$  střední hodnota veličiny  $u$  v blízkosti  $v_1$  určená zpravidla z trendu

$r(v_1)$  reziduum veličiny  $u$  v blízkosti  $v_1$

Pokud v poli sledované proměnné nejsou na sobě nezávislé, ale existuje mezi nimi korelace, doporučuje se použít strukturální funkce, které zahrnují např. 2 sledované proměnné. Takový postup se označuje jako koregionalizace dat (kapitola 4.1.6) (vícenásobné strukturální funkce).

Pro analýzu kvalitativních (předně nominálních) dat se používají neparametrické strukturální funkce (kapitola 4.1.7).

Při hodnocení se předpokládá nulový (nebo alespoň stejný) rozměr vzorků, pokud tento předpoklad není možné dodržet, používá se regularizace a deregularizace (kapitola 4.1.8).

### 4.1.1 Teoretické modely semivariogramu

Při interpretaci strukturálních funkcí se využívají jejich základní teoretické typy, označované jako modely. Jednotlivé modely si uvedeme spolu s jejich popisem pomocí semivariogramu.

#### Modely přechodového typu

U těchto klasických modelů je vyjádřena skutečnost, že při malých vzdálenostech je shoda mezi zjištěnými hodnotami vysoká (a tedy variabilita nízká), s rostoucí vzdáleností se "neshoda" zvyšuje až do určité vzdálenosti (=dosah), kde se úroveň neshody stabilizuje kolem hodnoty statistického rozptylu. Za touto vzdáleností se již neuplatňuje prostorová vazba mezi zkoumanými místy a variabilita je plně určována statistickým rozptylem.

Nejběžnější modely:

- Sférický model (kapitola 4.1.1.1)
- Kvadratický model (kapitola 4.1.1.2)
- Exponenciální model (kapitola 4.1.1.3)
- Gaussovský model (kapitola 4.1.1.4)
- Cirkulární model (kapitola 4.1.1.5)
- Pentasférický model (kapitola 4.1.1.6)
- Lineární model s prahem (kapitola 4.1.1.7)

#### Modely bez přechodu

Výskyt těchto modelů si lze zjednodušeně představit jako určitý extrémní případ klasického přechodového modelu. Představme si, že bychom u něho prováděli výpočet semivariogramu jen do vzdálenosti nepřesahující rozpětí  $d$ . Pak bychom při vynesení bodů nenašli žádnou oblast stabilizace křivky semivariogramu a daný případ bychom interpretovali jako model bez přechodu.

Často jsou to případy, kdy je variabilita ve zkoumaném poli výrazně nižší než je rozsah a detail našeho zkoumání.

Nejběžnější modely:

- Lineární (kapitola 4.1.1.8)
- Logaritmický (kapitola 4.1.1.9)
- Náhodný (kapitola 4.1.1.10)

#### Oscilační modely

Oscilační charakter (tj. s nehomogenním charakterem) má zkoumané pole nejčastěji v důsledku pravidelného střídání pásů s vyššími a nižšími hodnotami. Průměrná šířka pásů se dá odhadnout podle rozměru půlky periody vlny.

U těchto modelů se často projevuje nestabilita. Nepoužívají se pro odvození parametrů potřebných pro krigování (upřednostňují se robustní, jednoduché přechodové modely).

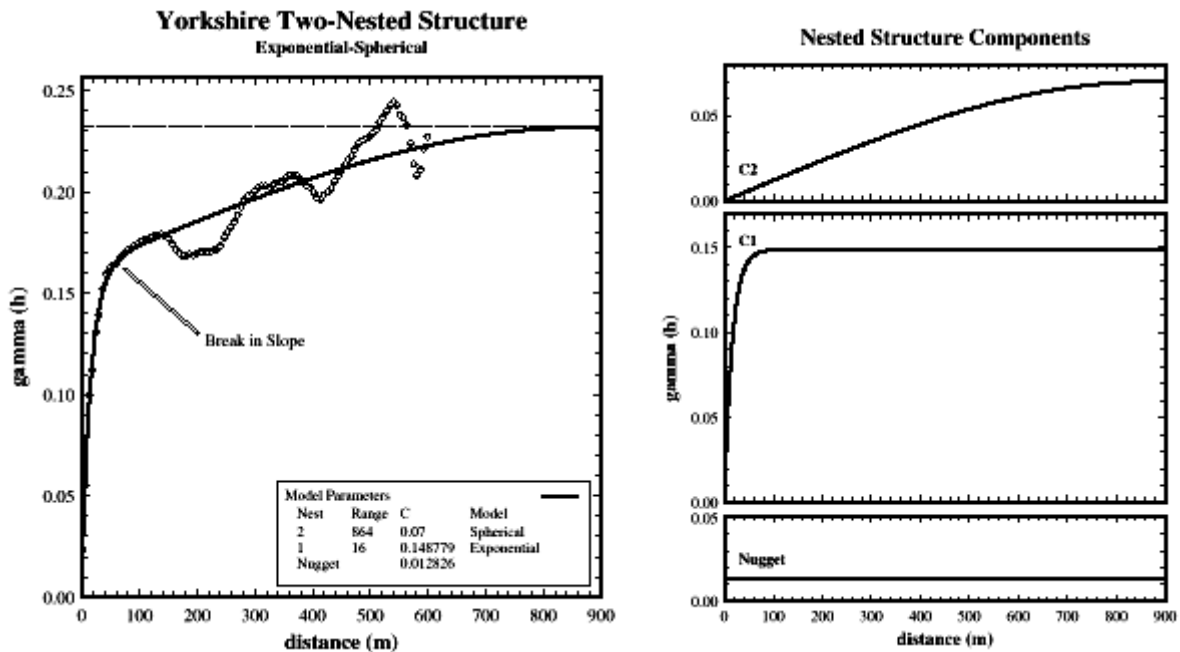
Nejběžnější modely:

- Sinový (kapitola 4.1.1.11)
- Kosinový (4.1.1.12)

#### Složené modely

Modely lze jednoduchým způsobem skládat do tzv. složených semivariogramů. Každému zdroji variability pak odpovídá vlastní semivariogram. Rovněž existenci zbytkového rozptylu lze vysvětlit pomocí složeného semivariogramu - jako existenci semivariogramu s relativně velmi malým dosahem (resp. náhodný semivariogram). Interpretaci více složených strukturálních funkcí lze doporučit pouze v případě věrohodného výpočtu semivariogramu (řádově stovky párů v každém kroku).

Vzorec složeného semivariogramu:  $\gamma^*(h) = \Sigma \gamma_i^*(h)$



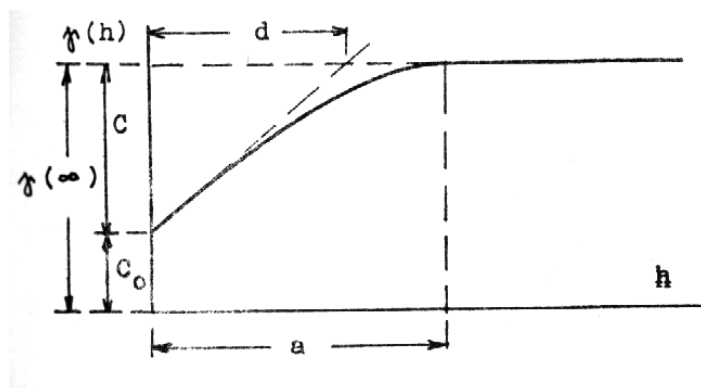
Obr.4-2 Složený model (vlevo průběh experimentálního a výsledného složeného modelu, vpravo dílčí semivariogramy) (Wingle, Poeter 1999)

#### 4.1.1.1 Sférický model

$$\gamma^*(h) = \begin{cases} C_0 + C * [1.5 * \frac{h}{a} - 0.5 * (\frac{h}{a})^3] & \dots\dots\dots h \leq a \\ C_0 + C & \dots\dots\dots h > a \end{cases}$$

$$a = 1.5 d$$

Často používaný model. Typicky se projevuje v případech, kdy v poli dominuje 1 zdroj variability.

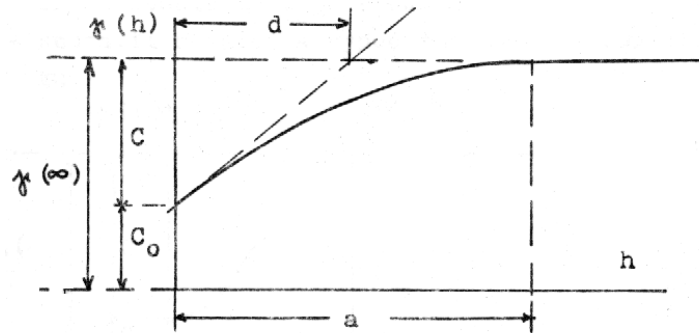


Obr. 4-3 Sférický model (Schejbal 1985)

#### 4.1.1.2 Kvadratický model

$$\gamma^*(h) = \begin{cases} C_0 + C^* \left[ 2 \cdot \frac{h}{a} - \left( \frac{h}{a} \right)^2 \right] & \dots\dots\dots h \leq a \\ C_0 + C & \dots\dots\dots h > a \end{cases}$$

$$a = 2d$$



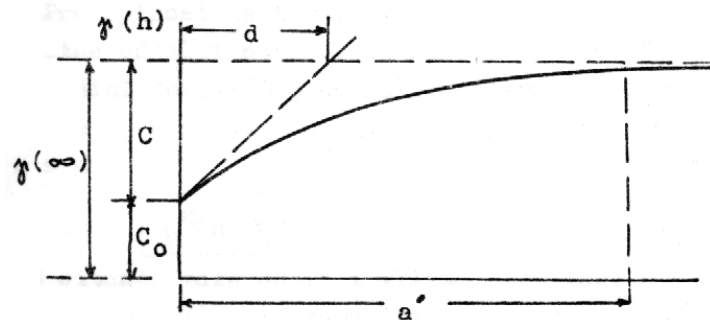
Obr. 4-4 Kvadratický model (Schejbal 1985)

#### 4.1.1.3 Exponenciální model

$$\gamma^*(h) = C_0 + C^* \left[ 1 - \exp\left(-\frac{h}{d}\right) \right]$$

$$a = 3d$$

Teoreticky tento model nemá práh ani dosah. Dosah se prakticky určuje dle místa, kde dosáhne křivka 95% maximální hodnoty. Model se objevuje např. u polí, kde působí více významných zdrojů variability.



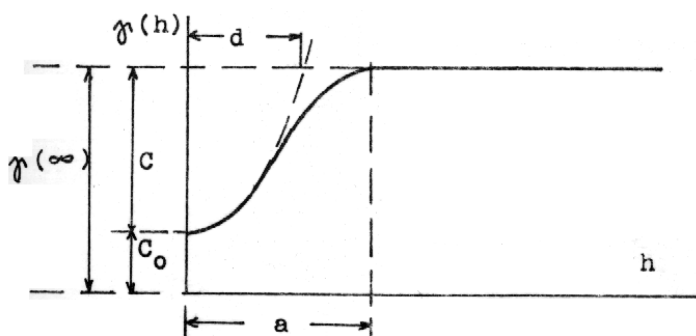
Obr. 4-5 Exponenciální model (Schejbal 1985)

#### 4.1.1.4 Gaussovský model

$$\gamma^*(h) = C_0 + C^* \left[ 1 - \exp\left(-\frac{h^2}{a^2}\right) \right]$$

$$a = d \cdot \sqrt{3}$$

Je dokumentován u dobře prozkoumaných polí. Naznačuje existenci plynulých změn hodnot. Setkáváme se s ním např. při modelování výškových dat. Posa (1989) však upozorňuje na častou nestabilitu tohoto modelu.



Obr. 4-6 Gaussovský model (Schejbal 1985)

#### 4.1.1.5 Cirkulární model

$$\gamma^*(h) = \begin{cases} C_0 + C * f & \dots\dots\dots h \leq a \\ C_0 + C & \dots\dots\dots h > a \end{cases}$$

kde

$$f = \left[ \frac{2h}{\pi a} \sqrt{1 - \left(\frac{h}{a}\right)^2} + \frac{2}{\pi} \arcsin \frac{h}{a} \right]$$

#### 4.1.1.6 Pentasférický model

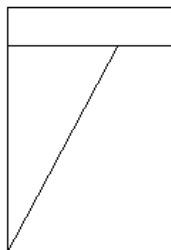
$$\gamma^*(h) = \begin{cases} C_0 + C * \left[ \frac{15h}{8a} - \frac{5}{4} * \left(\frac{h}{a}\right)^3 + \frac{3}{8} * \left(\frac{h}{a}\right)^5 \right] & \dots\dots\dots h \leq a \\ C_0 + C & \dots\dots\dots h > a \end{cases}$$

#### 4.1.1.7 Lineární model s prahem

Lineární model s prahem je jednoduchý a je poměrně často využíván zvláště programy provádějícími interpolaci pomocí krigování na základě automaticky vypočítaného a vyhodnoceného semivariogramu. Při provádění strukturální analýzy se využívá raději jiných přechodových modelů.

$$\gamma^*(h) = \begin{cases} C_0 + p * h & \dots\dots\dots h \leq a \\ C_0 + C & \dots\dots\dots h > a \end{cases}$$

p je směrnici přímky

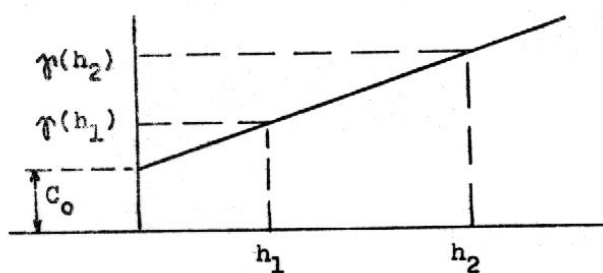


Obr. 4-7 Model lineárního semivariogramu s prahem

#### 4.1.1.8 Lineární model

$$\gamma^*(h) = C_0 + p \cdot h$$

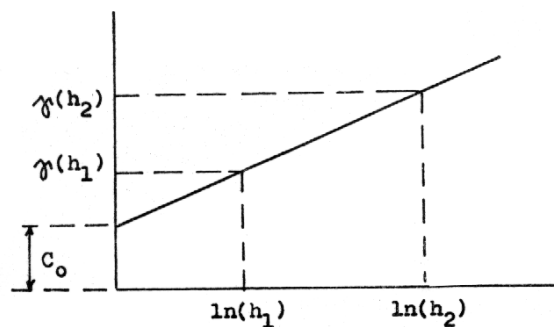
$p$  je směrnici přímky



Obr. 4-8 Lineární model (Schejbal 1985)

#### 4.1.1.9 Logaritmický model

$$\gamma^*(h) = C_0 + p \cdot \ln(h)$$



Obr. 4-9 Logaritmický model (Schejbal 1985)

#### 4.1.1.10 Náhodný model

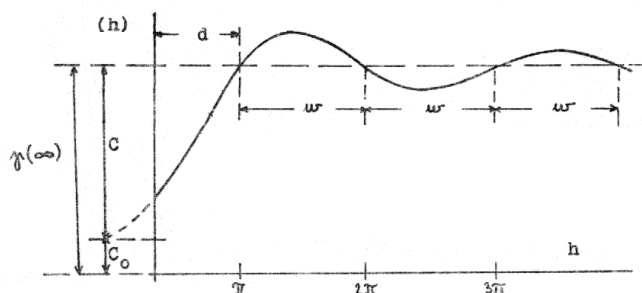
$$\gamma^*(h) = C_0$$

Semivariogram nemá žádnou úvodní rostoucí větev, hodnoty často pouze kolísají kolem prahu. K této situaci dochází, když je studované pole příliš variabilní vzhledem ke zvolenému kroku vzorkování (zjišťování hodnot).

#### 4.1.1.11 Sinový model (model s efektem propadu)

$$\gamma^*(h) = C_0 + C^* \left[ 1 - \frac{\sin(g^*h)}{g^*h} \right]$$

$$\text{kde } g = \frac{\pi}{\omega}$$

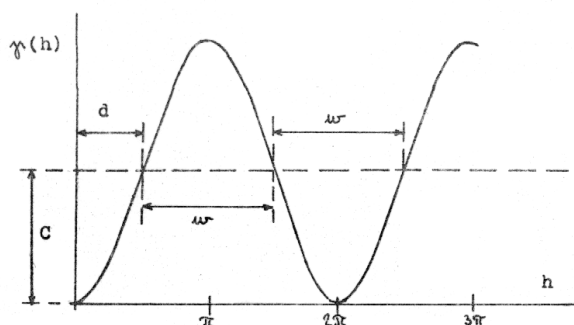


Obr. 4-10 Sinový model (Schejbal 1985)

#### 4.1.1.12 Kosinový model (periodický model)

$$\gamma^*(h) = C_0 + C^* [1 - \cos(g^*h)]$$

$$\text{kde } g = \frac{\pi}{\omega}$$



Obr. 4-11 Kosinový model (Schejbal 1985)

### 4.1.2 Charakteristika zkoumaného pole

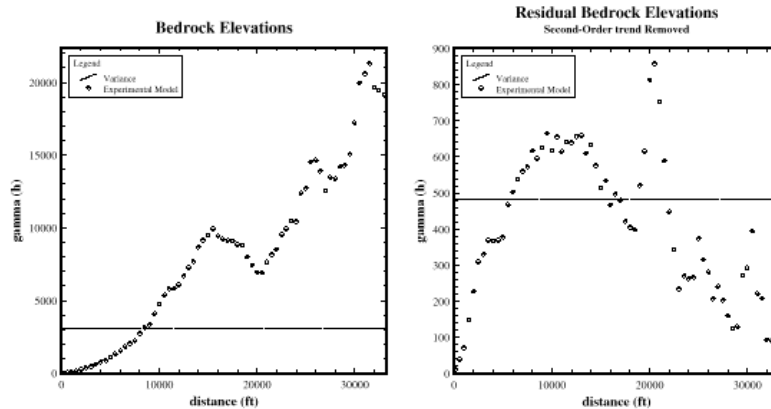
**Kontinuita pole** – popisuje velikost autokorelace. Vyjadřuje se u přechodových modelů dosahem semivariogramu, u oscilačních hodnotou rozpětí a u semivariogramů bez přechodu pomocí směrnice tečny k počátku semivariogramu. Pole s větší kontinuitou má tedy větší "dosah" prostorové korelace mezi hodnotami (tedy vyšší prostorovou autokorelaci).

**Nehomogenita pole** – popisuje existenci výrazných změn ve struktuře pole (výrazné změny ve střední hodnotě a rozptylu). Projevuje se výskytem oscilačního typu semivariogramu. Průměrný rozměr nehomogenity odpovídá polovině periody.

**Nestacionarita pole** – prokázána významná změna průměrné hodnoty proměnné v poli. Projevuje se zpravidla parabolickým nárůstem křivky semivariogramu v jisté vzdálenosti, za kterou již změny průměrné hodnoty v poli nejsou zanedbatelné. Prokazatelná je zvláště v případech, kdy dochází k parabolickému růstu křivky až za hodnotou dosahu, tedy na stabilizované části křivky.



**Anizotropie pole** – vlastnosti pole v různých směrech jsou různé (sledují se zejména rozdíly v průběhu směrových strukturálních funkcí). Anizotropii pole je možné popsat pomocí modelů jednotlivých směrových semivariogramů v poli, resp. elipsou anizotropie, proloženou dosahy směrových semivariogramů.



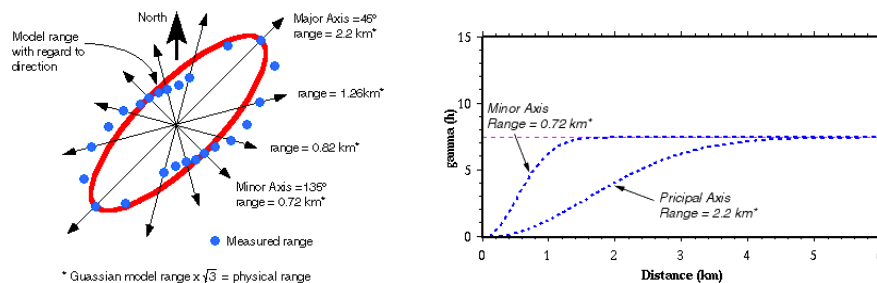
Obr. 4-12 Nestacionarita v poli (vlevo nestacionární semivariogram, vpravo semivariogram vypočtený z reziduí po odstranění trendu). Důvodem nestacionarity je zde klesání stropu sledované vrstvy v celém území k severovýchodu (Wingle, Poeter 1999).

### 4.1.3 Studium anizotropie pole

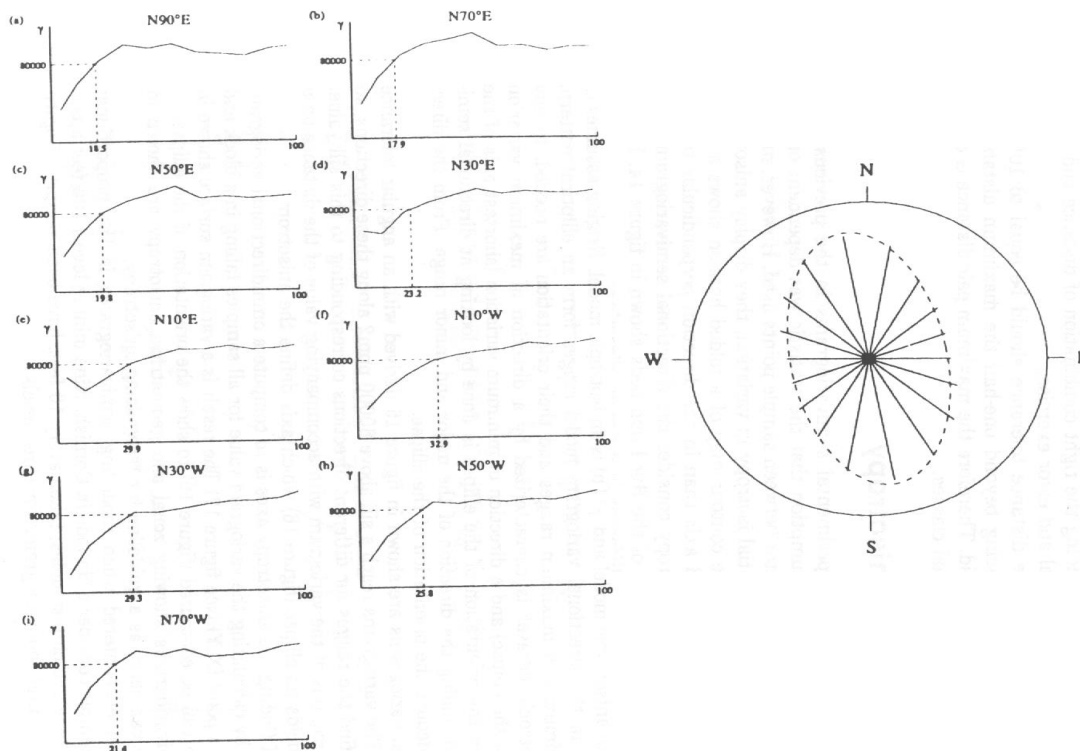
Studium autokorelace v rovině se provádí tak, že se vyhodnotí parametry **směrových semivariogramů** (tedy semivariogramů vypočtených na různých směrech v poli), případně se může vyhodnotit povrch anizotropie (kapitola 6.1.5).

Semivariogramy v jednotlivých směrech se snažíme interpretovat stejným modelem semivariogramu s těžce hodnotou prahu a zbytkového rozptylu, tedy pouze dosah může být rozdílný. Vykreslí se dosahy semivariogramů na liniích odpovídajících výpočtovým směřům. Pokud lze body aproximovat kružnicí, jde o **izotropní pole**. Jestliže body proložíme elipsou, jedná se o anizotropní pole (tzv. **geometrická anizotropie**), kde směr maximální variability pole odpovídá směru s minimálním dosahem (směr vedlejší osy elipsy) a směr minimální variability pole je shodný se směrem maximálního dosahu (tedy směrem hlavní osy elipsy).

Ostatní případy, kdy nelze dodržet stejný model pro všechny směry nebo stejnou hodnotu prahu nebo stejnou hodnotu zbytkového rozptylu, se označují jako **zonální anizotropie**, kterou již nelze jednoduše interpretovat.



Obr. 4-13 Konstrukce elipsy anizotropie z dosahy směrových semivariogramů a její parametry



Obr. 4-14 Interpretace směrových semivariogramů a konstrukce elipsy anizotropie (Meer 1992)

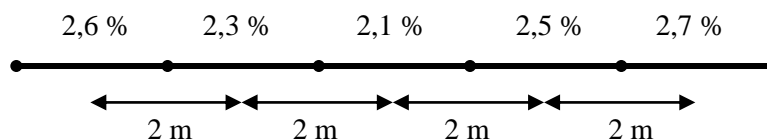
Nestejnou hodnotu prahu resp. zbytkového rozptylu v různých směrech lze eliminovat pomocí složeného modelu semivariogramu. Jednotlivé semivariogramy skládající složený model totiž mohou mít různou orientaci os a velikost poloos u elipsy anizotropie. V případě rozdílné hodnoty prahu lze snadno eliminovat rozdíl pomocí naložení modelu přechodového typu, který bude mít ve směru vysoké hodnoty prahu cílového semivariogramu velice krátký dosah (jeho práh se přičte k základnímu modelu) a ve směru nízké hodnoty prahu naopak dlouhý dosah.

Pokud slouží strukturální analýza v rovině jako základ pro následující krigování, doporučuje se volit jednoduchý a robustní model semivariogramu (rozhodně ne oscilační typy).

#### 4.1.4 Praktický postup analýzy autokorelace na linii

Znamé hodnoty jsou rozmístěny na linii, a to:

a) **pravidelně**, s konstantní vzdáleností mezi body



Obr. 4-15 Pravidelné uspořádání původních hodnot na linii (měření se vztahují ke středům intervalů s délkou 2 m)

Jako první krok  $h_1$  se volí základní vzdálenost mezi sousedními body. Několik chybějících hodnot v řadě není na závadu a není důvodem měnit základní krok.

Příklad:

Primární údaje - vzdálenost od počátku (m) a obsah Cu ve vzorku (%)

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1.20	1.02	0.62	0.20	0.14	0.13	0.24	0.22	0.24	0.22	0.35	0.35	0.34	0.39	0.66

Výpočet semivariogramu:

Pro krok 1 m

$$\gamma_1 = ((1.2-1.02)^2 + (1.02-0.62)^2 + (0.62-0.2)^2 + (0.2-0.14)^2 + (0.14-0.13)^2 + (0.13-0.24)^2 + (0.24-0.22)^2 + (0.22-0.24)^2 + (0.24-0.22)^2 + (0.22-0.35)^2 + (0.35-0.35)^2 + (0.35-0.34)^2 + (0.34-0.39)^2 + (0.39-0.66)^2) / 14 / 2 = 0.017$$

Pro krok 2 m

$$\gamma_2 = ((1.2-0.62)^2 + (1.02-0.2)^2 + (0.62-0.14)^2 + (0.2-0.13)^2 + (0.14-0.24)^2 + (0.13-0.22)^2 + (0.24-0.24)^2 + (0.22-0.22)^2 + (0.24-0.35)^2 + (0.22-0.35)^2 + (0.35-0.34)^2 + (0.35-0.39)^2 + (0.34-0.66)^2) / 13 / 2 = 0.054$$

Pro krok 3 m

$$\gamma_3 = ((1.2-0.2)^2 + (1.02-0.14)^2 + (0.62-0.13)^2 + (0.2-0.24)^2 + (0.14-0.22)^2 + (0.13-0.24)^2 + (0.24-0.22)^2 + (0.22-0.35)^2 + (0.24-0.35)^2 + (0.22-0.34)^2 + (0.35-0.39)^2 + (0.35-0.66)^2) / 12 / 2 = 0.091$$

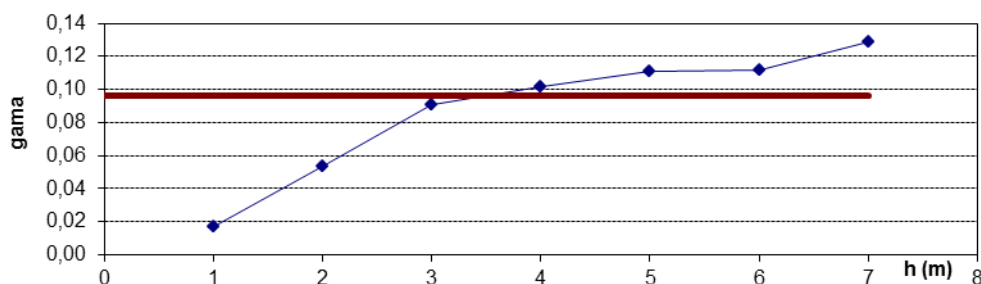
až po krok 7 m

$$\gamma_7 = ((1.2-0.22)^2 + (1.02-0.24)^2 + (0.62-0.22)^2 + (0.2-0.35)^2 + (0.14-0.35)^2 + (0.13-0.34)^2 + (0.24-0.39)^2 + (0.22-0.66)^2) / 8 / 2 = 0.128$$

Výběrový rozptyl je 0.096.

### Průběh semivariogramu

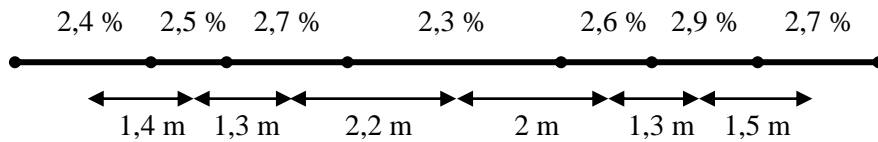
—●— experimentální semivariogram — práh



Obr. 4-16 Výsledný experimentální semivariogram

b) **nepravidelně**, s různými vzdálenostmi mezi sousedními body Pro první krok  $h_1$  se doporučuje průměrná minimální vzdálenost mezi sousedními body. Obecně lze říci, že větší krok dává hladší křivky (výpočet z většího počtu párů a tedy menší vliv nepravidelností), nevýhodou ale je malý počet bodů v počátku křivky, která je pro nás zásadní z hlediska interpretace.

K tomu je ještě nutné udat toleranci délky kroku. Někteří autoři doporučují volit 10% z délky kroku, každopádně horní hranicí je polovina délky kroku. Vzhledem k nepravidelnému rozmístění bodů uvnitř tolerance by měla být v každém kroku vypočítávána skutečná délka kroku jako průměrná vzdálenost všech dvojic bodů.

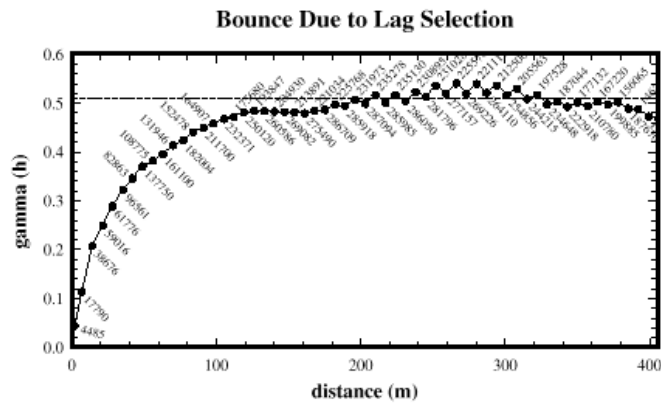


Obr. 4-17 Nepravidelné uspořádání na linii

Základním problémem je volba správné délky prvního kroku, další kroky  $h$  se zpravidla volí jako  $n$ -násobky prvního kroku. Výpočet se doporučuje provádět až do kroku s délkou odpovídající polovině délky celé linie (maximální vzdálenosti bodů), protože za touto vzdáleností již dochází k silnému úbytku párů bodů použitých ve výpočtu a tedy významně klesá spolehlivost vypočtené hodnoty.

Obecně je vhodné evidovat v každém kroku počet párů bodů vstupujících do výpočtu, který vypovídá o věrohodnosti výpočtu semivariogramu v daném kroku, a přihlížet k němu při interpretaci semivariogramu.

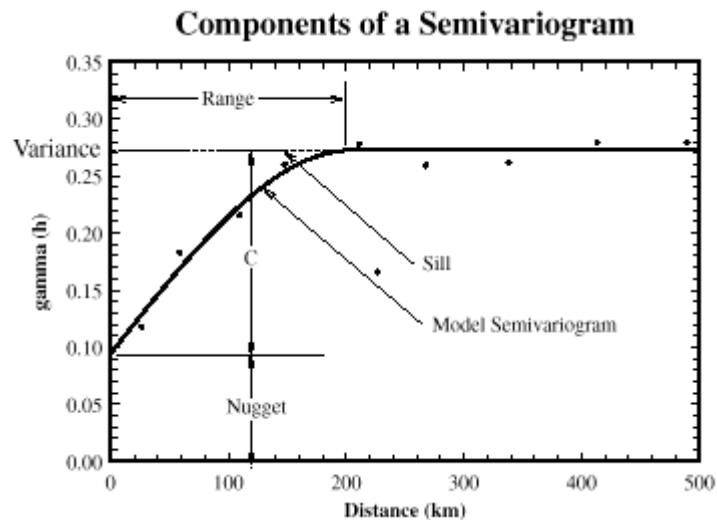
Nevhodně volený krok se projeví např. oscilací počtu párů  $a$ /nebo hodnoty experimentálního semivariogramu v jednotlivých krocích.



Obr. 4-18 Oscilace hodnot semivariogramu v důsledku nevhodně zvoleného kroku

Po provedení výpočtu strukturální funkce se provádí strukturální analýza, v našem případě tedy **interpretace semivariogramu**.

Nakreslí se graf, kde na horizontální ose se vynesou průměrné délky kroku a na vertikální ose příslušné hodnoty semivariogramu. Body se propojí jednoduchou čarou. Počátečními body se dále prokládá přímka, představující tečnu úvodní části křivky semivariogramu (vzhledem k definici semivariogramu nesmí protínat vertikální osu v záporné části). Nakreslí se přímka rovnoběžná s horizontální osou ve vzdálenosti rovné rozptylu původních hodnot (*variance*). Často odpovídá prahu (*sill*). Neodpovídá-li hodnota rozptylu pravděpodobné hodnotě prahu (např. v důsledku uplatnění proporcionalního efektu), je možné zvolit vhodnější práh. Úsek, který vytíná budoucí tečna křivky semivariogramu na prahu, se označuje jako rozpětí  $d$ . Úsek, který uvedená přímka vytíná na vertikální ose, představuje hodnotu zbytkového rozptylu  $C_0$  (*nugget*). Dále se odhaduje model strukturální funkce a hodnota dosahu (*range*).



Obr. 4-19 Parametry semivariogramu

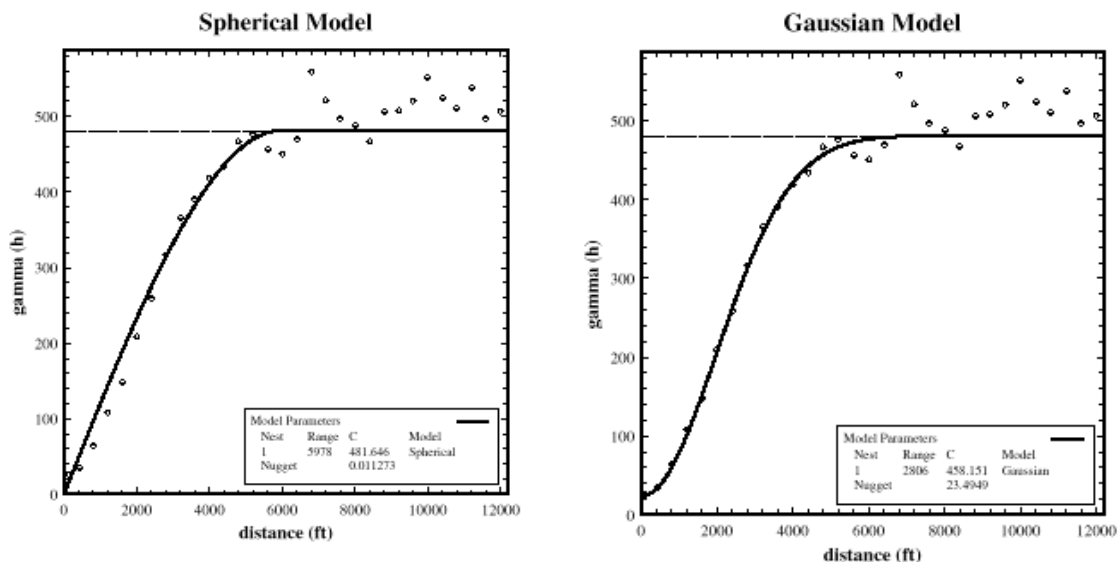
Z odhadnutých parametrů a modelu strukturální funkce se vypočte průběh teoretického semivariogramu a prověřuje se jeho shoda s vypočtenými hodnotami semivariogramu (tj. s experimentálním semivariogramem). Parametry a model se upravují až k dosažení uspokojivého výsledku.

Výsledkem analýzy je identifikovaný typ semivariogramu a jeho parametry.

Strukturální analýza poskytuje cenné poznatky o charakteru prostorové variability sledované proměnné.

Z hlediska hodnocení pole a využití např. při krigování je nejdůležitější správně interpretovat počátek semivariogramu až do dosažení hodnoty prahu.

Pokud již hodnota semivariogramu v 1.kroku je blízká prahu nebo je vyšší, jde zřejmě o čistě náhodný typ semivariogramu. Pokud je to možné, zkraťte krok nebo přidejte další údaje do výpočtu.



Obr. 4-20 Modelování průběhu semivariogramu s využitím sférického a gaussovského modelu

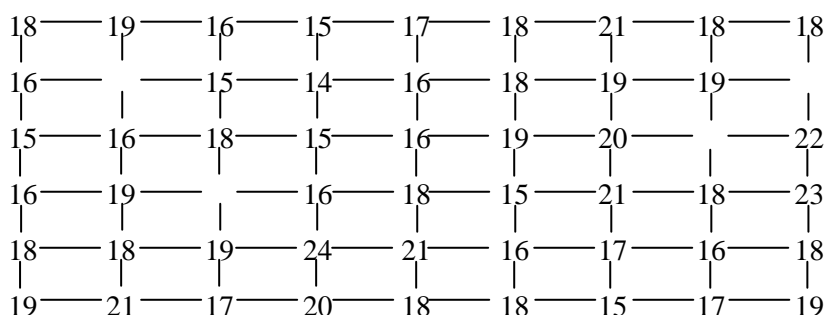
#### 4.1.5 Praktický postup analýzy autokorelace v rovině

Pro výpočet směrových semivariogramů se doporučuje volit 4 - 8 směrů.

Znamé hodnoty jsou rozmístěny v ploše, a to:

a) **pravidelně** např. v čtvercové nebo obdélníkové síti

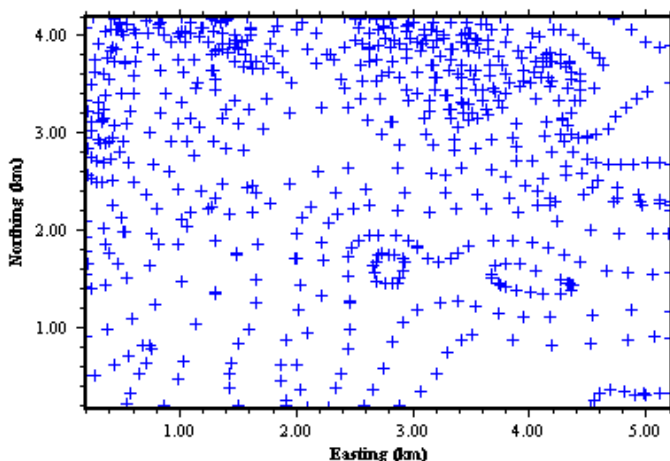
Volí se význačné směry, kde je k dispozici dostatek párů (např. horizontálně, vertikálně, po obou úhlopříčkách).



Obr. 4-21 Pravidelná síť hodnot (chybějící měření nevadí)

b) **nepravidelně**

Zpravidla se nechá proběhnout iterační proces, který zjišťuje vzdálenost a směr libovolného páru bodů a který zjištěné rozdíly hodnot zařadí do tříd podle vzdálenosti a směru.

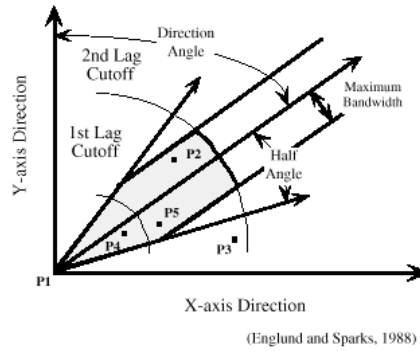


Obr. 4-22 Nepravidelná síť hodnot

K vymezení jednotlivých tříd se stanovuje několik parametrů:

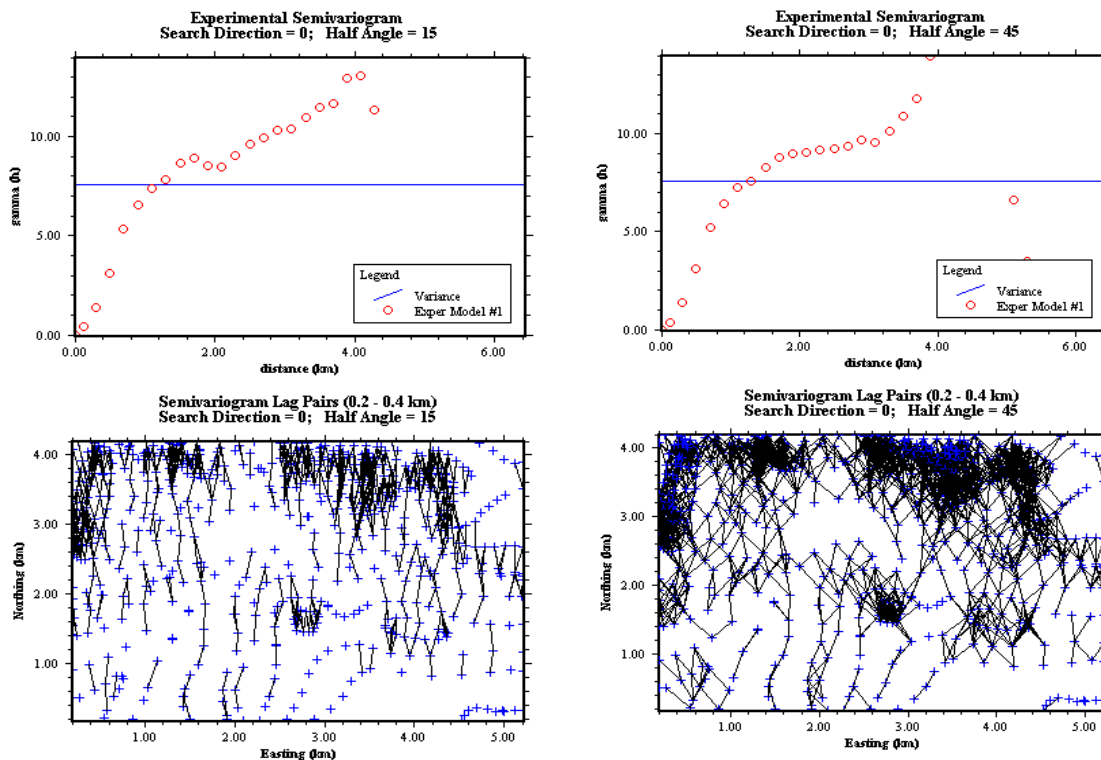
- délková tolerance (která byla vysvětlena v lineárním případě) (vzdálenost od - do)
- úhlová tolerance (*half angle*), kdy se stanovuje úhel, který je možné přičíst a odečíst od každého požadovaného směru a získat minimální a maximální směr pro výpočet směrového semivariogramu
- šířka pásma (*maximum bandwidth*) - používá se k eliminování takových dvojic, které mají větší vzdálenost a současně se značně odchylní z požadovaného směru (využívání maximální úhlové tolerance se připouští v případě malých vzdáleností). Tento parametr se používá tehdy, když chceme zabezpečit, aby výpočet semivariogramu probíhal z dat v téže zóně, např. ve stejné geologické jednotce.

Správné nastavení parametrů je značně závislé na hustotě primárních dat, jejich rozložení, počtu kroků v délce a směru. Vždy se snažíme získat maximální počet párů bodů pro výpočet, aby byly výsledky statisticky věrohodné. Benevolentně stanovené třídy (velká délková či úhlová tolerance, maximální šířka pásma) samozřejmě obsáhnou větší počet párů, výsledná křivka bude vypočtena z více hodnot a bude tedy věrohodnější, na druhou stranu ale může být její průběh rozkolísán tím, že do výsledku jsou zahrnovány příliš heterogenní dvojice, pokud jde o délkové či směrové vymezení. Vhodně zvolená orientace směrových semivariogramů spolu s menší úhlovou tolerancí může poskytnout lepší výsledky.



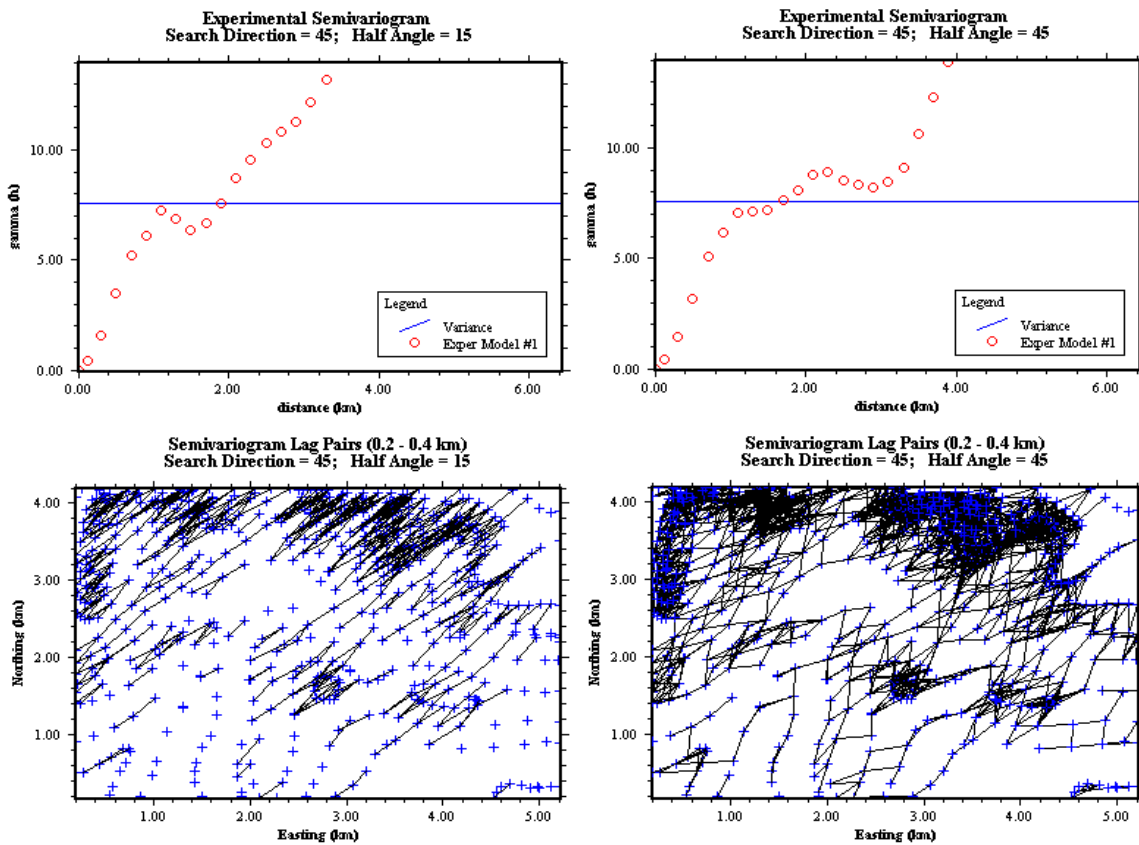
Obr. 4-23 Parametry výpočtu směrových semivariogramů

Dále je demonstrován výpočet směrových semivariogramů pro 2 různé úhlové tolerance (Wingle,

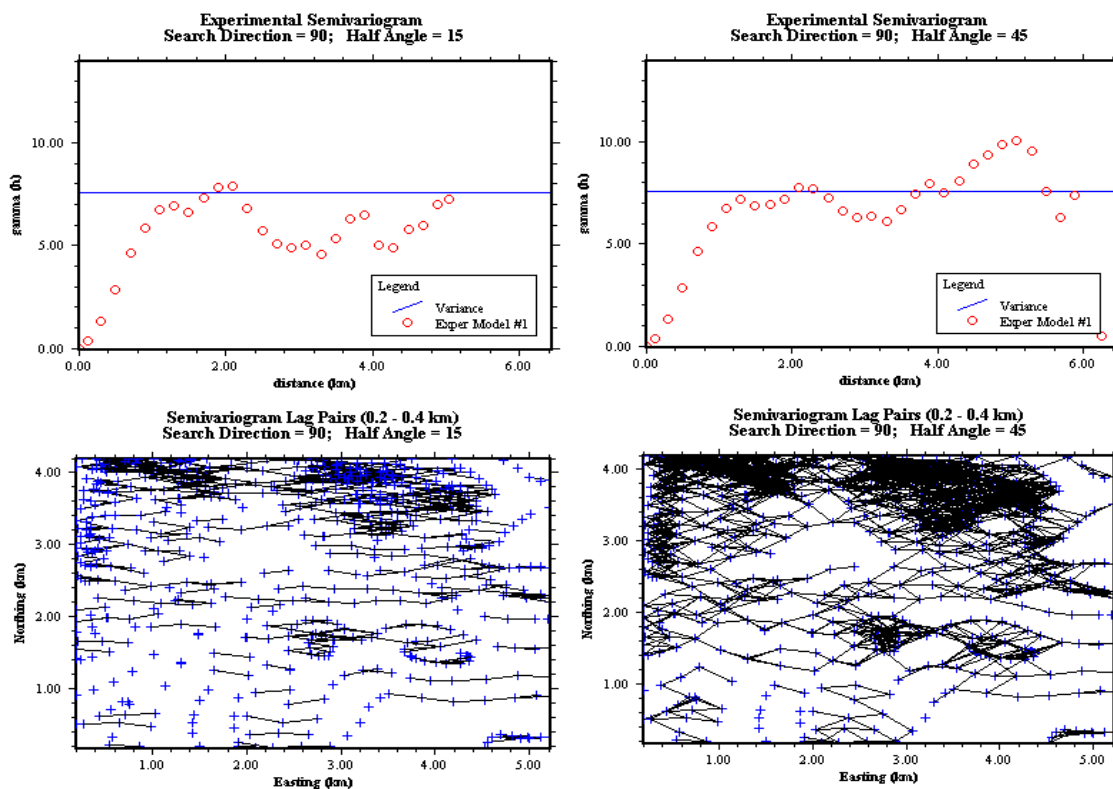


Poeter 1999):

Obr. 4-24 Výpočet semivariogramu pro směr  $0^\circ$  (vlevo tolerance  $15^\circ$ , vpravo  $45^\circ$ ). Spodní obrázky ukazují výběr párů pro výpočet semivariogramu pro krok 0,2-0,4 km.

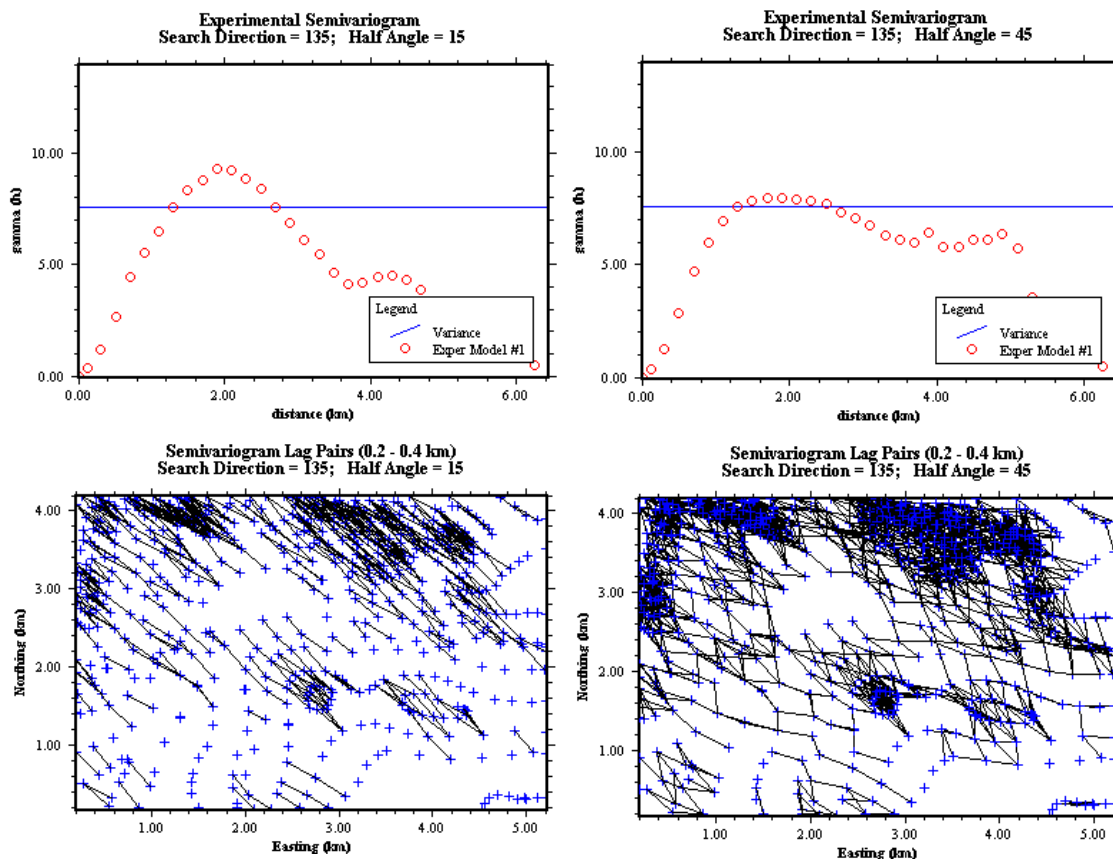


Obr.4-25 Výpočet semivariogramu pro směr  $45^\circ$  (vlevo tolerance  $15^\circ$ , vpravo  $45^\circ$ ). Spodní obrázky ukazují výběr párů pro výpočet semivariogramu pro krok 0,2-0,4 km.



Obr. 4-26 Výpočet semivariogramu pro směr  $90^\circ$  (vlevo tolerance  $15^\circ$ , vpravo  $45^\circ$ ). Spodní obrázky ukazují výběr párů pro výpočet semivariogramu pro krok 0,2-0,4 km.

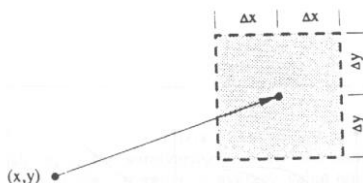




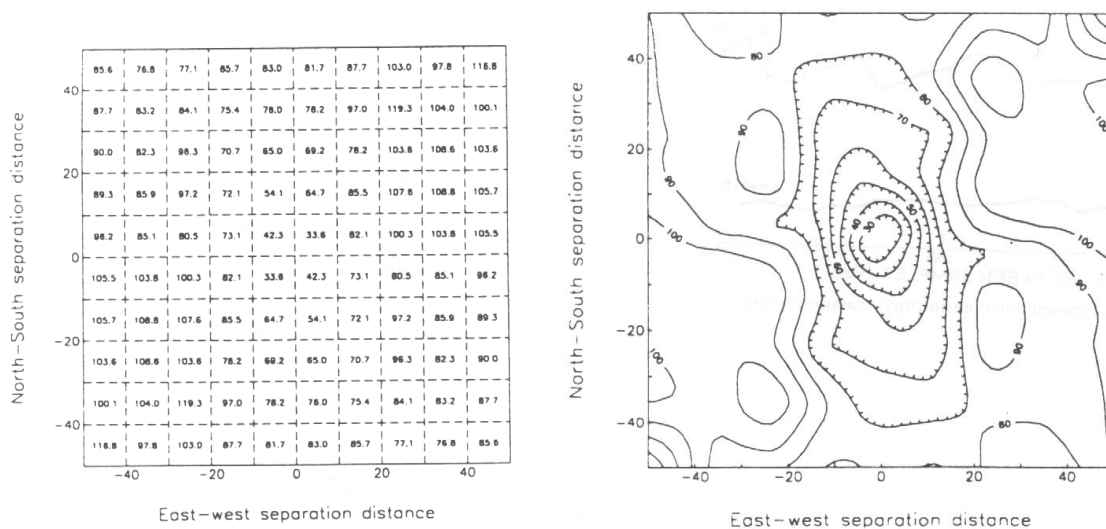
Obr. 4-27 Výpočet semivariogramu pro směr 135° (vlevo tolerance 15°, vpravo 45°). Spodní obrázky ukazují výběr párů pro výpočet semivariogramu pro krok 0,2-0,4 km.

### Studium povrchu anizotropie

Alternativou pro hodnocení anizotropie v rovině je studium povrchu anizotropie. Provádí se tak, že se celá oblast překryje pravouhloú mřížkou a hodnota semivariogramu se vypočítá mezi bodem se souřadnicemi X,Y a všemi body, které se nacházejí v buňkách určité vzdálenosti a směru. Vypočtená hodnota se запиše do této buňky a po skončení výpočtů je možné tento rastr (tvořený buňkami s hodnotou semivariogramu) vykreslit a interpretovat jako povrch anizotropie.



Obr. 4-28 Seskupování všech vzorků spadajících do čtverce určité vzdálenosti a směru od bodu x,y (pro výpočet povrchu anizotropie) (Meer 1992)



Obr. 4-29 Vlevo výpočet hodnoty semivariogramu pro každý čtverec, reprezentující určitou třídu vzdálenosti a směru, vpravo izolinie povrchu anizotropie (Meer 1992)

#### 4.1.6 Koregionalizace

Pokud v poli sledované proměnné nejsou na sobě nezávislé, ale existuje mezi nimi korelace, doporučuje se použít strukturální funkce, které zahrnují např. 2 sledované proměnné. Příkladem může být tzv. **vzájemný semivariogram**:

$$\gamma_{t,u}(h) = \frac{1}{2n_h} \sum_{n_h} [t(x) - t(x+h)] * [u(x) - u(x+h)]$$

$t(x)$  hodnota veličiny  $t$  v bodě  $x$

$u(x)$  hodnota veličiny  $u$  v bodě  $x$

Na rozdíl od "běžného" semivariogramu může vzájemný semivariogram nabývat záporných hodnot v případě, že mezi proměnnými existuje nepřímá korelace. Může stoupat i klesat v závislosti na  $h$  podle toho jaká je korelace mezi veličinami  $t$  a  $u$ .

Dokonce nemusí být mezi veličinami významná korelace v daném místě, ale posunutá o jistý vektor.

Je vyžadována stacionarita 2.řádu (viz předpoklady geostatické analýzy). Vzájemné semivariogramy a semivariogramy jednotlivých proměnných musí být modelované společně. Musí mít stejný model, stejnou anizotropii a mohou se lišit pouze hodnotou  $C$ . Každá struktura (jeden ze skládaných modelů), která se objeví ve vzájemném semivariogramu, se musí objevit i v semivariogramech samotných proměnných.

Je vyžadován lineární model koregionalizace a koregionalizační matice kladně semidefinitní. Aby byla zajištěna kladná hodnota odhadu rozptylu při kokrigingu za všech okolností, kontroluje se

$$|\gamma(h)| \leq \sqrt{\gamma_t(h) * \gamma_u(h)}$$

hodnota vzájemného semivariogramu pomocí Cauchy-Schwartzova vztahu:

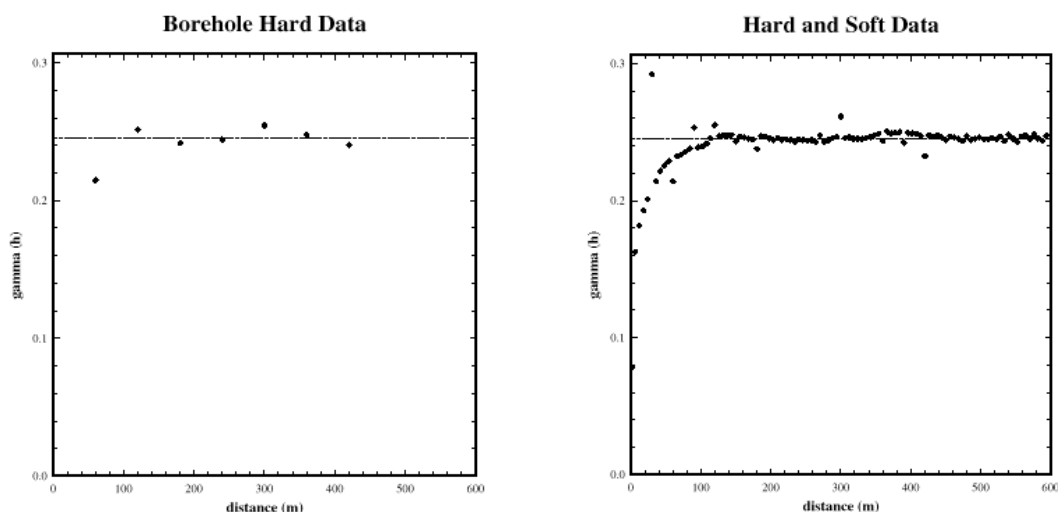
Pro všechny  $h \geq 0$

Pokud se zjišťuje hodnota vzájemného semivariogramu pro pár proměnných, které nejsou měřeny ve stejných místech, označuje se někdy výsledek jako **pseudovzájemný semivariogram** (*pseudo cross-variogram*).

### 4.1.7 Indikátorové a „soft“ semivariogramy

Indikátorové semivariogramy se konstruují a využívají při strukturální analýze nominálních (kvalitativních) dat, typu barva, druh horniny. Primární data se transformují do hodnot 1 a 0 podle splnění indikační podmínky - např. zda je hornina pískovcem. Často slouží jako vstup pro indikátorové krigování (viz kapitola indikátorové krigování).

Soft semivariogramy se využívají v případě nedostatku primárních dat, kdy je možné na základě provedené simulace doplnit další data a usnadnit provedení strukturální analýzy. Interpretace a verifikace je však dosti nesnadná a vyžaduje větší zkušenosti. Soft semivariogramy se často používají při provádění soft krigingu (kapitola 6.2.9).



Obr. 4-30 Využití simulace k doplnění dalších hodnot pro výpočet semivariogramu a k jeho přesnější interpretaci

### 4.1.8 Regularizace a deregularizace

Pozorování a měření sledované veličiny je vždy zjišťováno pro určitého nositele, který má nenulové rozměry. Zjištěná hodnota pak vlastně představuje průměrnou hodnotu veličiny v nositeli. Je evidentní, že s rostoucími rozměry vzorku bude klesat variabilita určení takové veličiny. Proto je i výpočet semivariogramu závislý na velikosti nositele - čím větší nositel, tím vyrovnanější průběh semivariogramu a nižší hodnota prahu.

Pokud je potřebné provést vzájemnou interpretaci několika semivariogramů, které byly vypočteny pro nositele o různé velikosti, je možné hodnoty semivariogramu  $\gamma_v(h)$  transformovat do tzv. bodového semivariogramu  $\gamma(h)$  (tedy pro nositele s nulovou velikostí - pro bod). Tento postup se označuje jako deregularizace. Korekční člen  $\gamma(v,v)$  závisí na tvaru a velikosti nositele (Schejbal 1983).

$$\gamma(h) = \gamma_v(h) + \bar{\gamma}(v,v)$$

Postup transformace bodového semivariogramu na semivariogram s nenulovým nositelem se označuje jako regularizace.

## 4.2 Lokální odhady - krigování

Lokálním odhadem rozumíme výpočet pravděpodobné hodnoty proměnné buď v bodě, kde nebylo provedeno měření (event. zjištění) = **bodový odhad**, anebo v relativně malé ploše = **blokový odhad**.

**Krigování** je základní geostatistickou metodou určování lokálního odhadu. V zahraničí se užívá i akronymu BLUE (Best Linear Unbiased Estimator), který dobře vystihuje výchozí podmínky krigování:

### lineární kombinace vstupních hodnot

$$(1) \quad u^* = \sum \lambda_i \cdot u_i$$

### nestranný odhad (průměrná chyba je rovna 0)

$$(2) \quad \sum(u^* - u_i) = 0$$

### minimalizace rozptylu odhadu

$$(3) \quad \sum(u^* - u_i)^2 = \text{minimum}$$

Podmínku lze vyjádřit pomocí semivariogramu i pomocí kovariační funkce a dalších strukturálních funkcí.

Vstupní podmínky se transformují do soustavy lineárních rovnic, která se řeší a jejím výsledkem je stanovení vah pro jednotlivé body s měřeními a Langrangeova multiplikátoru  $\mu$  (používá se pro zjednodušení rovnic u lineárního programování) pro výpočet rozptylu odhadu.

Krigování se provádí v různých modifikacích, nejdůležitější je základní krigování. Podle cíle odhadu se vyčleňují bodové a blokové odhady.

1. Způsob výpočtu základního krigování u bodového odhadu (kapitola 4.2.1).
2. Způsob výpočtu základního krigování u blokového odhadu (kapitola 4.2.2).
3. Jednodušší výpočet poskytuje jednoduché krigování (kapitola 4.2.3).
4. Vliv trendu je možné zahrnout do výpočtu při univerzálním krigování (kapitola 4.2.4).
5. Lokální odhad na základě 2 či více veličin se provádí pomocí kokrigingu (kapitola 4.2.6).
6. V řadě případů není distribuce hodnot normální, ale odpovídá lognormální distribuci. V takovém případě se používá lognormální krigování (kapitola 4.2.7), lepší variantou je ovšem provést přímo standardizaci proměnné, aby získala normální rozdělení (např. Z-skóre)
7. K neparametrickým geostatistickým metodám patří indikátorové krigování (kapitola 4.2.8), soft kriging (kapitola 4.2.9) a pravděpodobnostní krigování (kapitola 4.2.9).

K ověření kvality odhadu se používá **bumerangový test** (kapitola 4.2.11).

### 4.2.1 Bodový odhad při základním krigování

Podmínky:

- Neznámá, ale konstantní střední hodnota (žádný trend v datech! Pokud existuje, odstraní se jako externí a počítá se s reziduálními hodnotami)
- Intrinziční hypotéza (nejsou rozdíly mezi semivariogramy zjištěnými v různých místech, závisí jen na vektoru  $h$ ).
- suma vah = 1
- Normalita dat

Pro **bodový odhad** při **základním krigování** by soustava rovnic v maticovém tvaru vypadala následovně:

$$\begin{bmatrix} K_{11} & K_{12} & \dots & K_{1n} & 1 \\ K_{21} & K_{22} & \dots & K_{2n} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ K_{n1} & K_{n2} & \dots & K_{nn} & 1 \\ 1 & 1 & 1 & 1 & 0 \end{bmatrix} * \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_n \\ \mu \end{bmatrix} = \begin{bmatrix} K_{10} \\ K_{20} \\ \vdots \\ K_{n0} \\ 1 \end{bmatrix}$$

kde  $K_{ij}$  je kovariance mezi vzorky  $i$  a  $j$ ,  $K_{i0}$  kovariance mezi vzorkem  $i$  a odhadovaným bodem,  $\lambda_i$  je váha vzorku  $i$  a  $\mu$  Lagrangeův multiplikátor.

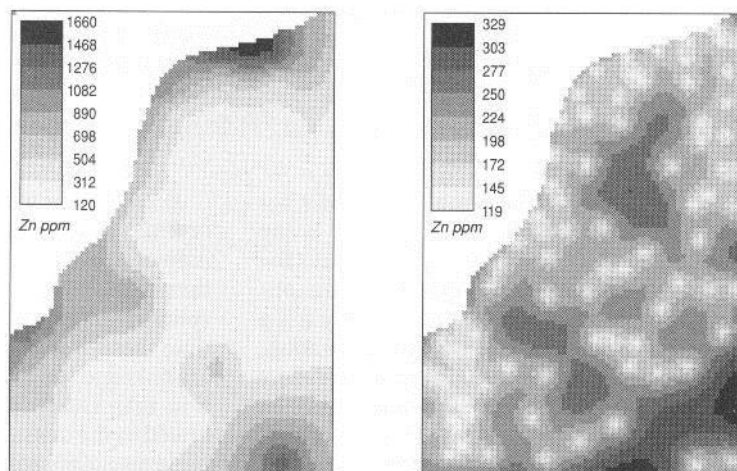
Rozptyl odhadu se vypočte jako:

$$\sigma_E^2 = \sigma_0^2 - \left( \sum_{i=1}^n \lambda_i K_{i0} + \mu \right)$$

kde  $\sigma_0^2$  je kovariance v odhadovaném bodě (=zbytkový rozptyl), význam ostatních veličin viz výše.

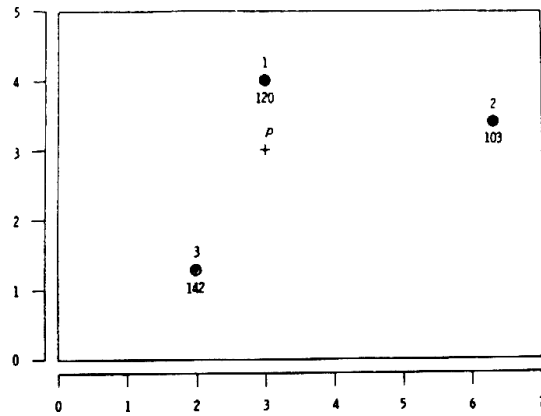
Vypočtené hodnoty rozptylu odhadu ukazují na kvalitu interpolace v ploše. Často se vykresluje místo rozptylu odhadu směrodatná odchylka odhadu, označovaná jako krigovací chyba, protože má stejnou jednotku jako vlastní odhad.

V některých případech (zvláště pokud není kovariance definována) se používá místo kovariance hodnota semivariogramu.



Obr.4-31 Odhad obsah Zn (netransformované hodnoty, vlevo) a krigovací chyba (vpravo) (Burrough, McDonell 1998)

Příklad: Postup odhadu úrovně hladiny podzemní vody v bodě p



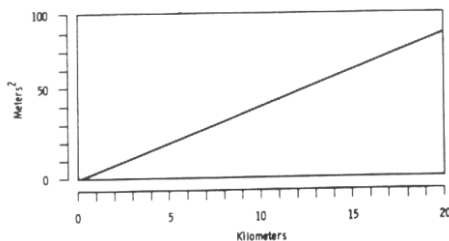
Obr.4-32 Situace - rozmístění známých vrtů a poloha bodu p (Meer 1992)

Tab. Poloha vrtů, měřené hodnoty a poloha bodu p

	X souřadnice	Y souřadnice	Hladina vody (m)
Vrt 1	3.0	4.0	120
Vrt 2	6.3	3.4	103
Vrt 3	2.0	1.3	142
Místo p	3.0	3.0	?

Tab. Vzdálenosti mezi body

	Vrt 1	Vrt 2	Vrt 3	Místo p
Vrt 1	0	3.35	2.88	1.00
Vrt 2		0	4.79	3.32
Vrt 3			0	1.97



Obr.4-33 Příslušný lineární model semivariogramu

Tab. Hodnoty semivariogramu pro příslušné vzdálenosti (předpoklad izotropního pole)

	Vrt 1	Vrt 2	Vrt 3	Místo p
Vrt 1	0	13.42	11.52	4.00
Vrt 2		0	19.14	13.30
Vrt 3			0	7.89

Pro lokální odhad hodnoty v bodě P ze 3 naměřených hodnot by vypadala soustava rovnic následovně:

$$\lambda_1 * \gamma(h_{11}) + \lambda_2 * \gamma(h_{12}) + \lambda_3 * \gamma(h_{13}) + \mu = \gamma(h_{1p})$$

$$\lambda_1 * \gamma(h_{12}) + \lambda_2 * \gamma(h_{22}) + \lambda_3 * \gamma(h_{23}) + \mu = \gamma(h_{2p})$$

$$\lambda_1 * \gamma(h_{13}) + \lambda_2 * \gamma(h_{23}) + \lambda_3 * \gamma(h_{33}) + \mu = \gamma(h_{3p})$$

$$\lambda_1 + \lambda_2 + \lambda_3 = 1$$

kde  $\gamma(h_{12})$  je hodnota semivariogramu pro vzdálenost mezi body 1 a 2 apod.,  $\lambda_i$  je váha pro bod 1.

Tedy v našem případě řešíme soustavu 4 rovnic o 4 neznámých:

$$\lambda_1 * 0 + \lambda_2 * 13.4 + \lambda_3 * 11.5 + \mu = 4$$

$$\lambda_1 * 13.4 + \lambda_2 * 0 + \lambda_3 * 19.1 + \mu = 13.3$$

$$\lambda_1 * 11.5 + \lambda_2 * 19.1 + \lambda_3 * 0 + \mu = 7.9$$

$$\lambda_1 + \lambda_2 + \lambda_3 = 1$$

respektive v maticovém zápisu:

$$\begin{bmatrix} 0 & 13.4 & 11.5 & 1 \\ 13.4 & 0 & 19.1 & 1 \\ 11.5 & 19.1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix} * \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \\ \mu \end{bmatrix} = \begin{bmatrix} 4.0 \\ 13.3 \\ 7.9 \\ 1.0 \end{bmatrix}$$

Inverze matice

$$\begin{bmatrix} -0.0680 & 0.0326 & 0.0354 & 0.1932 \\ 0.0326 & -0.0433 & 0.0106 & 0.4072 \\ 0.0354 & 0.0106 & -0.0461 & 0.3995 \\ 0.1932 & 0.4072 & 0.3995 & -9.5851 \end{bmatrix} * \begin{bmatrix} 4.0 \\ 13.3 \\ 7.9 \\ 1.0 \end{bmatrix} = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \\ \mu \end{bmatrix}$$

Výsledek:

$$\begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \\ \mu \end{bmatrix} = \begin{bmatrix} 0.5954 \\ 0.0975 \\ 0.3071 \\ -0.7298 \end{bmatrix}$$

Výsledkem jsou hodnoty jednotlivých vah a hodnota Langrangeova multiplikátoru.

Váhy se dosadí do vzorce (1) a vypočte se hodnota v bodě P ( $U_p$ ).

$$U_p = 0.5954 * 120 + 0.0975 * 103 + 0.3071 * 142 = 125.1 \text{ m}$$

Rozptyl odhadu veličiny v bodě P se vypočte podle vztahu:

$$s^2 = \lambda_1 * \gamma(h_{1p}) + \lambda_2 * \gamma(h_{2p}) + \lambda_3 * \gamma(h_{3p}) + \mu$$

tedy

$$s_p^2 = 0.5954 * 4 + 0.0975 * 12.1 + 0.3071 * 7.9 - 0.7298 * 1 = 5.25 \text{ m}^2$$

## 4.2.2 Blokový odhad při základním krigování

Pro **blokový odhad** při **základním krigování** by soustava rovnic v maticovém tvaru vypadala podobně jako u bodového odhadu:

$$\begin{bmatrix} K_{11} & K_{12} & \dots & K_{1n} & 1 \\ K_{21} & K_{22} & \dots & K_{2n} & 1 \\ \vdots & & & & \\ K_{n1} & K_{n2} & \dots & K_{nn} & 1 \\ 1 & 1 & 1 & 1 & 0 \end{bmatrix} * \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_n \\ \mu \end{bmatrix} = \begin{bmatrix} K_{1V} \\ K_{2V} \\ \vdots \\ K_{nV} \\ 1 \end{bmatrix}$$

$K_{ij}$  je kovariance mezi vzorky  $i$  a  $j$ ,

$K_{iV}$  kovariance mezi vzorkem  $i$  a odhadovaným blokem  $V$ ,

$\lambda_i$  je váha vzorku  $i$

$\mu$  Langrageův multiplikátor

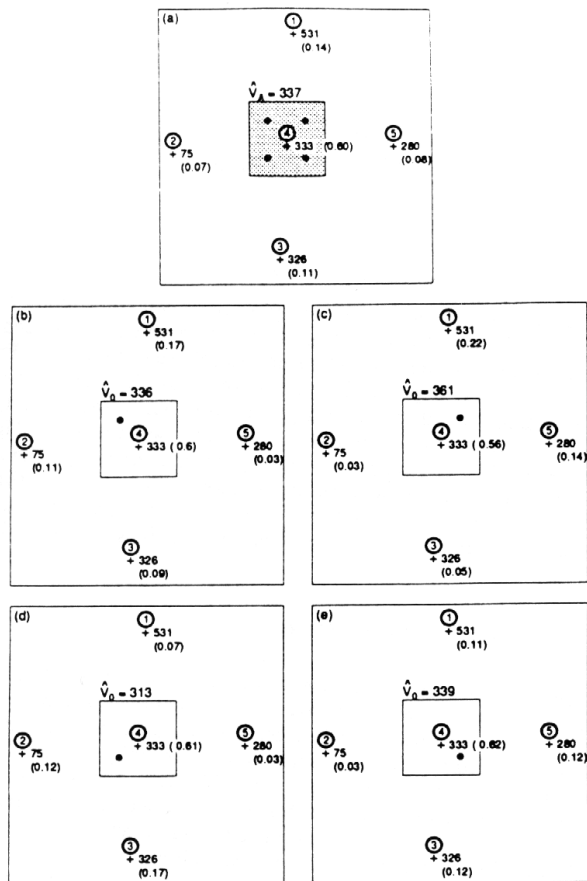
Využívá se zde skutečnosti, že kovariance  $K_{iV}$  mezi vzorkem  $i$  a průměrnou hodnotou veličiny v bloku  $V$  je stejná jako průměr kovariancí mezi vzorkem  $i$  a hodnotou veličiny v náhodně zvoleném bodě uvnitř  $V$ .

Rozptyl odhadu pro blok se vypočte jako:

$$\sigma_{EV}^2 = K_V - \left( \sum_{i=1}^n \lambda_i K_{iV} + \mu \right)$$

$K_V$  je kovariance v odhadovaném bloku, kterou lze vypočítat jako průměrnou hodnotu z kovariancí mezi páry náhodně zvolených bodů uvnitř bloku  $V$ .





Obr. 4-34 Postup výpočtu blokového odhadu (blok vyznačen šedě, aproximace 4 náhodnými body, průměr z těchto 4 bodových odhadů je 337 a odpovídá průměrné hodnotě bloku) (Meer 1992)

### 4.2.3 Jednoduché krigování

Nejjednodušší variantou krigování je tzv. **jednoduché krigování** (simple kriging). K výpočtu je potřebná průměrná hodnota veličiny v poli (m).

Podmínky:

- Vyžadována stacionarita 2.řádu. Konstantní střední hodnota (m) a strukturální fce.

Počítá se s reziduálními hodnotami, tím dojde ke zjednodušení výpočtu (zmizí podmínka suma vah = 1, jednoduchá soustava lineárních rovnic).

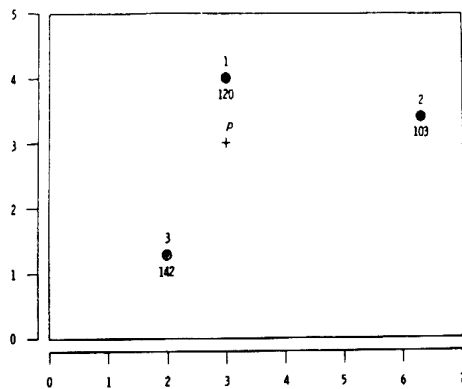
Výsledná hodnota:

$$u^* = m + \sum[\lambda_i^* (u_i - m)]$$

Optimalizuje se soustava:

$$\sum[\lambda_i^* K_{ij}] = K_{j0}$$

Příklad:



$$\lambda_1 * K_{11} + \lambda_2 * K_{12} + \lambda_3 * K_{13} = K_{1p}$$

$$\lambda_1 * K_{21} + \lambda_2 * K_{22} + \lambda_3 * K_{23} = K_{2p}$$

$$\lambda_1 * K_{31} + \lambda_2 * K_{32} + \lambda_3 * K_{33} = K_{3p}$$

Výpočet krigovacího rozptylu:

$$\sigma_{SK}^2 = K_0 - \left( \sum_{i=1}^n \lambda_i K_{i0} \right)$$

Vzhledem k podmínce stacionarity 2.řádu lze adekvátně použít i semivariogramy.

Často se provádí standardizace (NST) a declustering aby nebyly body koncentrovány na některých místech.

#### 4.2.4 Univerzální krigování

V případě nesplnění podmínek stacionarity průměrné hodnoty a strukturální funkce je nutné použít hypotézy **univerzálního krigování**. Prostorová proměnná je pak považována za součet 2 komponent - trendu (driftu), který určuje průměrnou hodnotu v tomto místě, a reziduí. Po výpočtu trendu lze získat hodnot rezidua odečtením hodnoty trendu v daném místě od skutečné hodnoty.

K popisu trendu se používají polynomy zpravidla 1. nebo 2. stupně.

$$M_0 = b_1X + b_2Y$$

nebo

$$M_0 = b_1X + b_2Y + b_3X^2 + b_4XY + b_5Y^2$$

$M_0$  je hodnota trendu v sledovaném bodě.

Strukturální funkce se pak vypočítají z reziduí hodnot. Krigování je prováděno pro rezidua hodnot, k výsledné hodnotě je nutno ještě připočítat hodnotu trendu v daném místě.

V praxi pak stačí rozšířit soustavu rovnic o členy popisující trend (koeficienty pro jednotlivé členy). Soustava rovnic pro polynom 1.stupně:

$$\begin{bmatrix} K_{11} & K_{12} & \dots & K_{1n} & 1 & X_1 & Y_1 \\ K_{21} & K_{22} & & K_{2n} & 1 & X_2 & Y_2 \\ \vdots & & & & & & \\ K_{n1} & K_{n2} & & K_{nn} & 1 & X_n & Y_n \\ 1 & 1 & & 1 & 0 & 0 & 0 \\ X_1 & X_2 & & X_n & 0 & 0 & 0 \\ Y_1 & Y_2 & & Y_n & 0 & 0 & 0 \end{bmatrix} * \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_n \\ \mu \\ b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} K_{10} \\ K_{20} \\ \vdots \\ K_{n0} \\ 1 \\ X_0 \\ Y_0 \end{bmatrix}$$

$X_i$  a  $Y_i$  jsou souřadnice bodu  $i$ ,  
 $b_1$  a  $b_2$  koeficienty v rovnici polynomu.

#### 4.2.5 Disjunktivní krigování (DK)

Jde o nelineární zobecnění krigování.

Namísto  $u_i = m + e_i$  se používá  $f(u_i)$ , kde  $f$  je libovolná funkce.

Výsledná hodnota (prediktor) je

$$g^*(u) = \Sigma f(u_i)$$

Podmínky:

- vyžaduje alespoň dvourozměrnou normalitu (bivariate normality) a aproximace funkcí  $f$ . Předpoklady je obtížné ověřit a výpočet je komplikovaný.

Doporučuje se provádět nejdříve standardizaci (NST).

#### 4.2.6 Kokriging

V případě vzájemné závislosti více zkoumaných veličin je možné provádět tzv. **kokriging** (cokriging), který provádí odhad proměnné na základě hodnot korelovaných veličin a jejich vztahů popsaných pomocí vzájemného semivariogramu. Důvodem použití může být situace, kdy vedle přímých měření zkoumané veličiny máme k dispozici měření či stanovení jiných veličin, která jsou mnohem levnější a máme jich k dispozici mnohem více. Např. se vedle přímého měření veličiny použijí výsledky nepřímých metod stanovení. Lze použít dokonce i korelované veličiny, které mají charakter faktorů jako např. vzdálenost od pravděpodobného zdroje.

Kokriging se provádí jak pro bodový tak pro blokový odhad, jednoduché nebo základní krigování. K výpočtu se používají vzájemné semivariogramy (cross-variogram, viz strukturální funkce).

Předpokládejme, že máme k dispozici měření  $Z$  sledované veličiny a také  $k$  měření jiných veličin  $V_k$  s indexem 1, 2, 3 až  $V$  na místech  $\mathbf{n}_V$ . Potom výpočet proměnné  $\mathbf{z}$  v místě  $\mathbf{x}_0$ :

$$z'(x_0) = \sum_{k=1}^V \sum_{i=1}^{n_V} \lambda_{ik} * z_{ik}$$

Nestrannost odhadu u základního krigování zajistíme splněním následujících rovnic:

$$\sum_{i=1}^{n_V} \lambda_{ik} = 1 \dots Z = V$$

$$\sum_{i=1}^{n_V} \lambda_{ik} = 0 \dots Z \neq V$$

První podmínka stanovuje, že musí být nejméně jedno pozorování Z k dispozici pro kokriging. Interpolační váhy se určí na základě minimalizace rozptylu:

$$\sigma_z^2(x_0) = \frac{1}{n} \sum (z(x_0) - z'(x_0))^2$$

Pro každou kombinaci místa a atributu je jedna rovnice, takže celý systém rovnic zahrnuje k proměnných na g místech:

$$\sum_{j=1}^V \sum_{i=1}^{n_v} \lambda_{ij} \gamma_{ij}(x_{ij}, x_{gk}) + \varphi_k = \gamma_{zV}(x_0, x_{gk})$$

pro všechny g=1 až n<sub>v</sub> a k=1 až V

kde  $\varphi_k$  je Lagrangeův multiplikátor.

Pokud jsou všechny další veličiny měřeny ve stejných místech jako Z, nezískáme výsledky odlišné od základního krigování (tedy má smysl provádět jen pro různá místa).

#### 4.2.7 Lognormální krigování

Lognormální krigování představuje krigování upravené pro data, která vykazují lognormální distribuci. Data jsou nejdříve transformována do přirozených či dekadických logaritmů. Výpočet a modelování semivariogramů (strukturální analýza) se provádí s logaritmovanými hodnotami stejně jako následující krigování. Výsledek se musí nejenom zpětně transformovat, ale navíc ještě korigovat o vznikající systematickou chybu:

$$u^* = \exp\left(u_{\ln}^* + \frac{s_{\ln}^2(u)}{2} - \mu\right)$$

kde  $u_{\ln}^*$  je interpolovaná hodnota

$s_{\ln}^2$  je rozptyl odhadu

$\mu$  je Lagrangeův multiplikátor

Deutsch a Journel (1992) upozorňují na vysokou citlivost na odlehlé hodnoty při zpětné transformaci a doporučují multigaussovské krigování nebo indikátorové krigování (kapitola 4.2.8).

#### 4.2.8 Indikátorové krigování (Indicator Kriging)

Zatím uváděné metody krigování prováděly lokální odhad na základě přímo naměřených (zjištěných) hodnot. Co však v případě, že nás nezajímá, zda je v daném místě nebo na dané ploše nejpravděpodobnější průměrná koncentrace 0.016 nebo 0.015, ale odhad pravděpodobnosti, s jakou je překročena limitní hodnota např. 0.012?

Řešením je právě indikátorové krigování ze skupiny neparametrických geostatistických metod

Indikátor nabývá pouze 2 hodnot - ano/ne resp. 1/0. Původní zjištěné hodnoty se nahradí hodnotou indikátoru. Indikátor má hodnotu 1, pokud původní údaj v místě splňuje stanovenou podmínku (např. barva=červená nebo hodnota přesahuje zvolený limit), a hodnotu 0, pokud tomu tak není. Tímto způsobem je možné krigovat i kvalitativní údaje i vykreslovat např. mapy rizika špatné klasifikace.

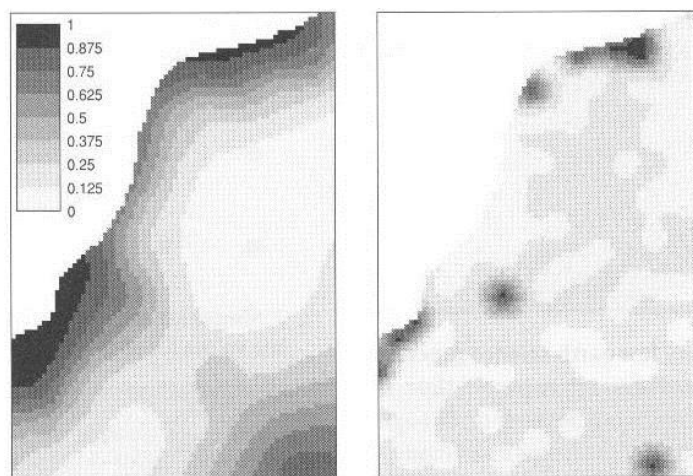
Transformované pole hodnot slouží jako vstup pro krigování, kde se používá standardních postupů. Výsledkem lokálního odhadu je pak pravděpodobnost, se kterou je v daném místě splněna testovaná podmínka.

Ve výše uvedeném příkladě tedy pravděpodobnost překročení hodnoty 0.012 na daném místě.

Často se volí několik limitů, pro které se hodnoty kategorizují (nahrazují 1 nebo 0). Z výsledků lze pak usoudit na % zastoupení jednotlivých "stupňů" znečištění, kovnatosti apod. Uplatňuje se pak představa, že se na sledovaném místě vyskytují v určitém zastoupení všechny kategorie - např. 70% nejnižší kategorie, 20% střední kategorie a 10% nejvyšší kategorie.

Následně lze také odhadnout průměrnou hodnotu v daném místě. Proč vypočítávat průměrnou hodnotu takto složitě a ne přímo z naměřených hodnot? Protože průměrná hodnota vypočtená z indikátorového krigování není při správné volbě limitů vychýlena extrémními hodnotami, je k nim mnohem méně citlivá.

Při indikátorovém krigování se větší měrou vyskytují problémy s interpretací strukturálních funkcí (vznikají velké oblasti totožných hodnot - nul nebo jedniček). V tom případě se doporučuje volit limity podle mediánů (1. limit = medián celého souboru, 2. vyšší limit = medián horní poloviny tj. horní kvartil celého souboru, atd.), které poskytují rozumné výsledky (Issaks and Shrivastava 1989).



Obr. 4-35 Indikátorové krigování pro obsah Zn v půdě - pravděpodobnost překročení hodnoty 500ppm (vlevo), pravděpodobnost překročení hodnoty 1000ppm (vpravo). (Burrough, McDonell 1998)

#### Příklad 1:

K dispozici je 5 vzorků s obsahem Au, cílem je stanovit zastoupení jednotlivých typů rudy ve vymezeném bloku a celkovou střední hodnotu obsahu Au v bloku.

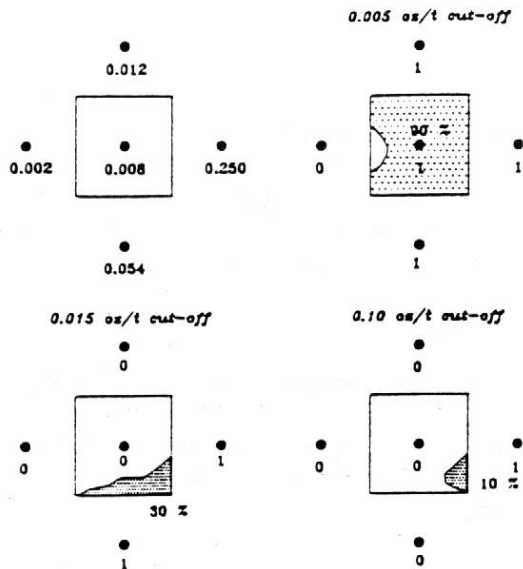
Limity: 0.005 oz/t, 0.015 oz/t, 0.1 oz/t.

Výsledky:

Blok obsahuje 10% horniny s nulovým obsahem Au (pod 0.005 oz/t), 60% nízkokovnaté rudy (0.005 - 0.015 oz/t) se střední hodnotou 0.009 oz/t, 20% středněkovnaté rudy (0.015 - 0.1 oz/t) se střední hodnotou 0.039 oz/t a 10% vysokokovnaté rudy (nad 0.1 oz/t) se střední hodnotou 0.332 oz/t.

Střední hodnota:

$$0.1 \times 0 + 0.6 \times 0.009 + 0.2 \times 0.039 + 0.1 \times 0.332 = \underline{0.046 \text{ oz/t}}$$



Obr. 4-36 Výpočet indikátorů pro stanovené limitní kovnatosti bloku rudy (příklad 1) (Meer 1992)

Běžně vypočítávané indikátory mohou mít problémy s nekonzistencí výsledků, protože se indikátory pro jednotlivé kategorie počítají a vyhodnocují nezávisle, což nemusí odpovídat realitě. Typicky v případě stanovování postupných limitů může dojít k tomu, že obsah vyšší kategorie v daném místě je vyšší než součet nižší a vyšší kategorie. Tedy např. pro limit  $Z > 5$  vypočítáme množství 12%, ale pro limit  $Z > 10$  vypočítáme množství 15%, což by znamenalo záporné množství v kategorii  $Z$  mezi 5 a 10.

Řešením je použití relativních indikátorů. Indikátory jsou stanovovány pouze pro ty primární údaje, které již překročily předchozí limit. Všechny primární údaje, které již předchozí limity nesplnily, dostanou automaticky hodnotu 0 bez ohledu na jejich skutečnou původní hodnotu.

Výpočet se někdy označuje jako **hnízděné indikátorové krigování**.

#### Příklad 2:

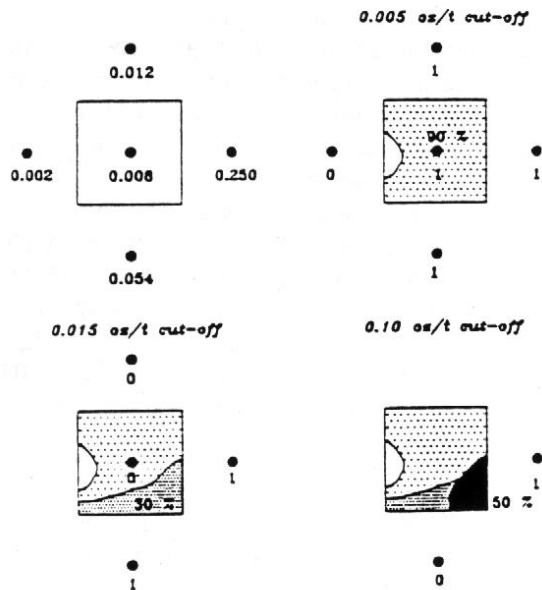
Zadání je stejné jako v 1. příkladu.

Výsledky:

Blok obsahuje 10% horniny s nulovým obsahem Au (pod 0.005 oz/t), 63% nízkokovnaté rudy (0.005 - 0.015 oz/t) se střední hodnotou 0.009 oz/t, 13.5% středněkovnaté rudy (0.015 - 0.1 oz/t) se střední hodnotou 0.039 oz/t a 13.5% vysokokovnaté rudy (nad 0.1 oz/t) se střední hodnotou 0.332 oz/t.

Střední hodnota:

$$0.1 \times 0 + 0.63 \times 0.009 + 0.135 \times 0.039 + 0.135 \times 0.332 = \underline{0.054 \text{ oz/t}}$$



Obr.4-37 Výpočet relativních indikátorů pro stanovené limitní kovnatosti bloku rudy (příklad 2)

Výtkou vůči indikátorovému krigování může být, že nevyužívá celou informaci (tedy naměřenou hodnotu). Z kombinace indikátorové hodnoty s původní naměřenou hodnotou vychází tzv. pravděpodobnostní krigování.

#### 4.2.9 Pravděpodobnostní krigování (*Probability Kriging*)

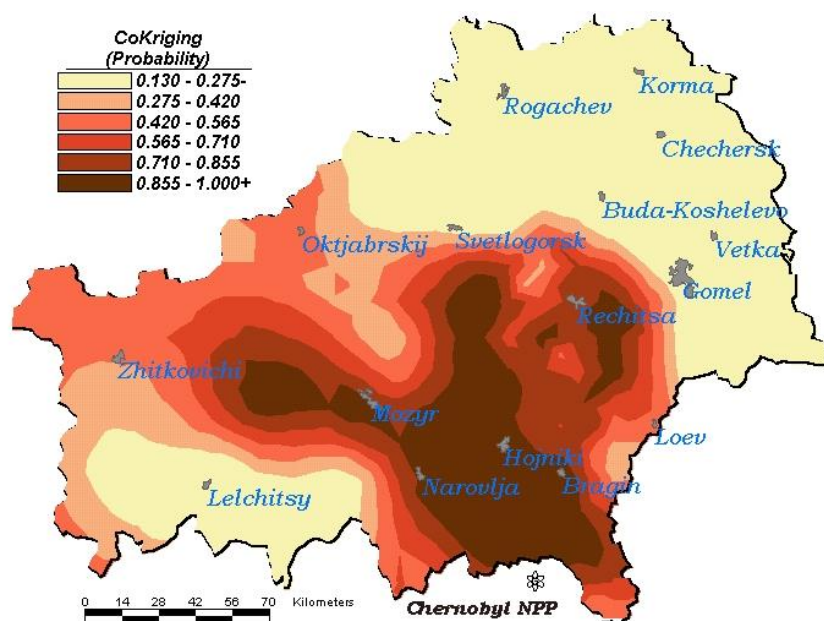
Lokální odhad se vypočte jako

$$u^* = \sum a_j \cdot i_{z,u} + \sum b_j \cdot u_j$$

kde  $a_j$  a  $b_j$  představuje váhy,

$i_{z,u}$  hodnotu indikátoru pro hodnotu  $u_j$  při limitu  $z$ ,

$u_j$  je původní naměřená hodnota.



Obr. 4-38 Pravděpodobnostní krigování - pravděpodobnost překročení limitní hodnoty 100 Bq/m<sup>2</sup> u <sup>241</sup>Am v půdě v oblasti severně od Černobylu v roce 1992. (Krivoruchko 1999)

Při výpočtu se využívá jak indikátorové funkce tak i původních hodnot. Původní hodnoty jsou často z důvodu nutnosti používání stejné jednotky nahrazeny pořadím původního údaje ve variační řadě.

#### 4.2.10 Soft Kriging

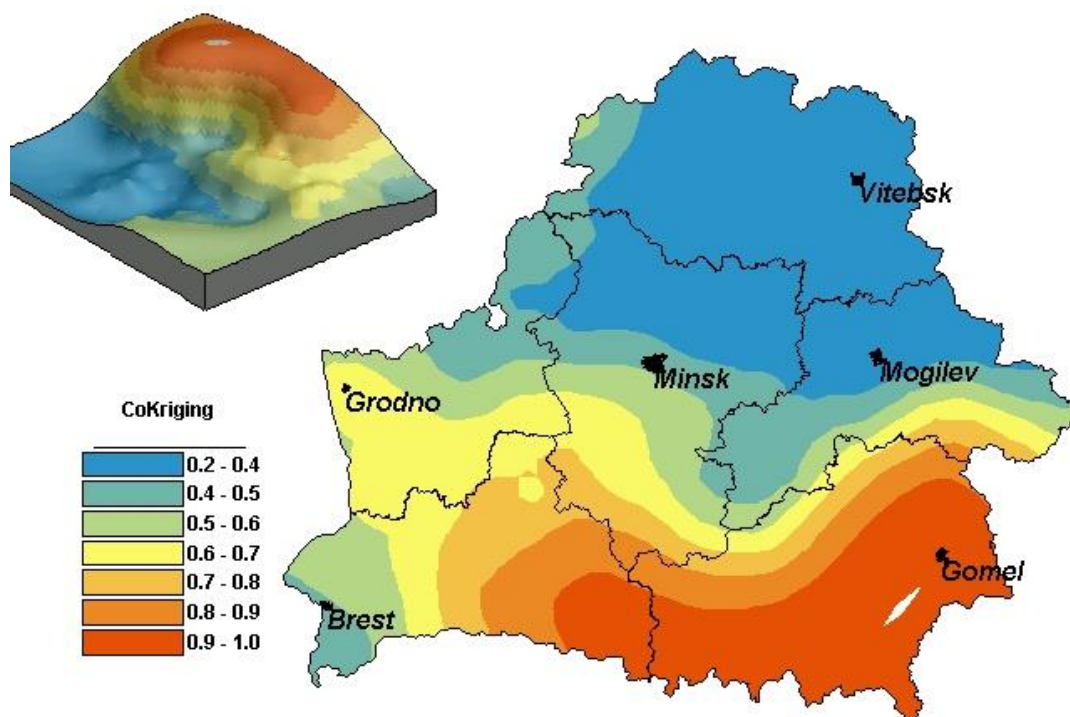
Dosavadní techniky považovali naměřené (nebo jinak zjištěné) hodnoty za skutečné hodnoty v daném místě. To ale často nebývá pravda. Naměřené údaje jsou zatíženy chybou, jejíž velikost závisí na řadě faktorů. Často jsme schopni stanovit měřenou hodnotu v daném místě pomocí intervalu:

$$a(x_j) \leq z(x_j) \leq b(x_j)$$

nebo pomocí distribuce pravděpodobnosti.

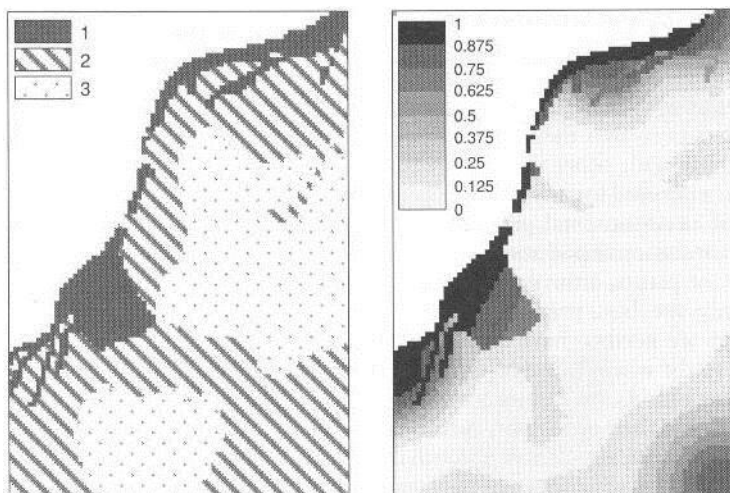
Naměřené údaje dokonce můžeme zkombinovat s dalšími znalostmi o daném fenoménu - např. i bez jakéhokoliv měření víme, že reálný rozsah hodnot v daném místě je od a do b. Všechny tyto informace lze využít při výpočtu lokálních odhadů pomocí soft krigingu. Řada autorů uvádí podstatné zlepšení kvality odhadu při použití doplňkové „soft“ informace. Výpočet se zpravidla provádí pomocí kokrigingu (kapitola 4.2.6).

Krivoruchko uvádí použití soft krigovací techniky v případě zjišťování pravděpodobnosti výskytu rakoviny štítné žlázy u dětí, kde se modifikovala pravděpodobnost zjištěná z výskytu rakoviny vahou velikosti osídlení (předpokládaly se větší váhy údajů zjištěných pro větší osídlení).



Obr. 4-39 Soft krigování - podmíněná pravděpodobnost, že výskyt rakoviny štítné žlázy u dětí je vyšší než 1 z 10000 v Bělorusku pro období 1986 - 1995. (Krivoruchko 1999)

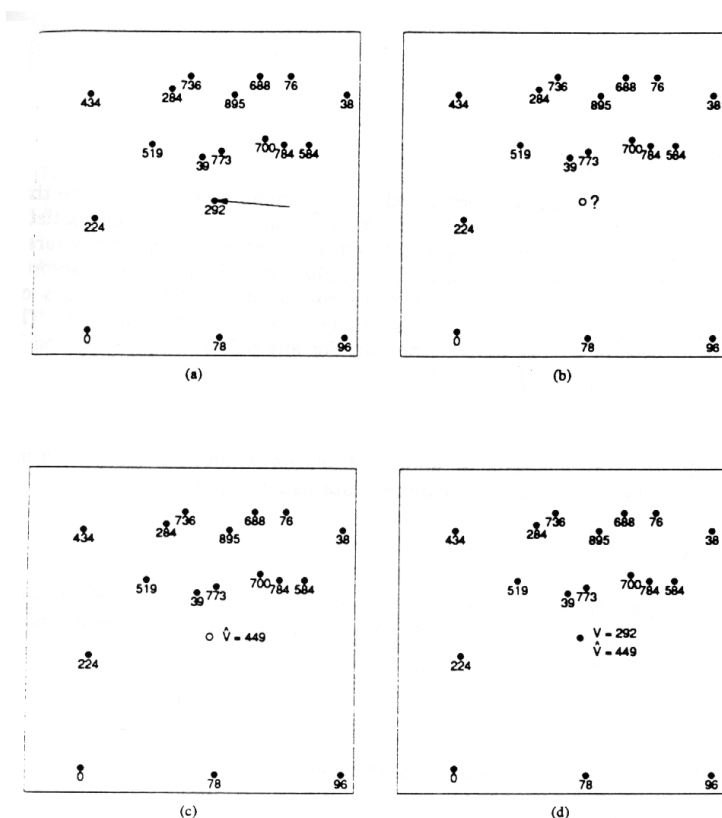




Obr. 4-40 Soft krigování pro obsah Zn v půdě - vlevo třídy četnosti záplav, vpravo výsledek, tj. pravděpodobnost překročení hodnoty 500ppm s využitím pravděpodobné hodnoty obsahu Zn v jednotlivých třídách četnosti záplav (Burrough, McDonell 1998)

#### 4.2.11 Bumerangový test

Ověřování kvality odhadu je možné provádět pomocí **bumerangového testu**, kdy se pro bod se známou hodnotou provede výpočet lokálního odhadu z ostatních hodnot. Výsledkem je vypočtená hodnota v místě, kde známe skutečnou hodnotu. Můžeme tedy stanovit chybu odhadu v tomto místě. Testování opakujeme pro náhodně vybrané body nebo pro všechny body.



Obr.4-41 Demonstrace bumerangového testu (provádí se výpočet pro místo označené šipkou, původní hodnota 292 slouží pro ověření velikosti chyby odhadu) (Meer 1992)

#### Chyby z bumerangové metody

Statistické hodnocení za celý soubor:

- Průměrná chyba (mean error)
- Střední chyba (RMSE)
- Průměrná krigovací chyba (Average standard error)
- Průměrná chyba standardizovaná (Mean standardized)
- Střední chyba standardizovaná (RMSE standardized )

Mapa chyb:

- Krigovací chyba
- Chyba indikátorů

#### 4.2.12 Omezení krigování

Krigování má také svá **omezení**. Krigování bylo definováno jako nejlepší nestranný lineární odhad. V této definici jsou skryta další omezení - proč by měl být odhad vypočítáván pouze z lineární kombinace hodnot, proč by mělo být jediným optimalizačním kritériem minimum rozptylu? Objevují se nové techniky, které využívají nelineárních odhadů či používají jiné optimalizační kritérium (běžně se definuje "ztrátová funkce", která se minimalizuje). Jinou možností provádění lokálních odhadů představuje simulace (následující kapitola).

Krigování se soustřeďuje se na co nejvěrnější odhad střední hodnoty v daném místě, avšak **degraduje rozptyl** (nové pole má již minimální rozptyl proti původnímu, je značně vyhlazené). I když se odvozuje např. horní a dolní odhad (rozsah odhadu) s pomocí interpolované hodnoty a směrodatné odchylky, jsou i tyto povrchy značně vyhlazené, protože jsou vypočteny pouze z několika hodnot.

V řadě případů však nepožadujeme přesný odhad střední hodnoty v daném místě, ale spíše realistický odhad situace v daném místě, tj. pravděpodobné hodnoty, variability. Např. pro výběr vhodných dobývacích či surovinu upravujících metod je potřebné dobře odhadnout i rozptyl hodnot, tedy přirozenou variabilitu suroviny.

Jiným důvodem může být situace, kdy potřebujeme určit **spojenou pravděpodobnost pro jistou sadu míst**. Např. nás zajímá pravděpodobnost toho, že sada míst  $x_i$  ( $i=1$  až  $N$ ) má hodnotu  $Z$  vyšší než je určitá limitní hodnota. Taková sada míst může představovat např. puklinový systém v hornině, který zabezpečuje mnohem vyšší propustnost pro tekutiny než je okolní, relativně homogenní prostředí horniny.

$$K(z_0; x_j, j = 1, \dots, N | (n)) = \prod_{j=1}^n (1 - F(z_0; x_j | (n)))$$

Takový požadavek lze vyjádřit následovně:

Spojená pravděpodobnost je tedy součinem jednotlivých pravděpodobností. Z toho plyne, že výsledná pravděpodobnost se významně snižuje u delších řetězců (větší  $n$ ). Ve skutečnosti však je pravděpodobnost o něco vyšší díky prostorové závislosti.

Mapy lokálních odhadů pomocí krigování nemusí odhalit existenci takových cest či bariér, protože jsou založeny na jiném optimalizačním kritériu a nezajistí tudíž vhodnou reprodukci textury. V těchto případech je však mnohem významnější správně odhadnout existenci „cest“ než získat lokálně přesné odhady.

Řešením je v obou případech použití **simulace**.

### 4.3 Simulace pro modelování regionalizované proměnné

Nejjednodušší simulační postupy využívají metody Monte Carlo, která simuluje hodnoty v závislosti na zvolené distribuci pravděpodobnosti se známou střední hodnotou a variabilitou. Výsledkem je stacionární šum.

Základní rozšíření představuje využití informace o vývoji prostorové variability v poli při znalosti semivariogramu či jiné strukturální funkce. Používají se **metody nejbližšího souseda** a **metoda rotujících pásem** (*turning bands*), která simuluje dobře vývoj v poli (snaží se modelovat pole tak, aby zůstala zachována nejenom průměrná hodnota, ale i prostorová kovariance a tedy i rozptyl v poli). Tato metoda zachovává prostorovou kontinuitu, ale nerespektuje naměřené hodnoty.

#### Metoda rotujících pásem

Redukuje 3D (zřejmě častěji 2D) simulace do několika nezávislých 1D simulací podél linií, které jsou pak rotovány ve 3D. Uvažujme linii  $D_1$  ve 3D, na které je definována kovariance  $C_{(hD_1)}$  a střední hodnota je nulová. Pak  $x_{D_1}$  bude projekcí nějakého bodu  $x$  ve 3D na linii  $D_1$  a zavede se 3D-náhodná funkce taková, že  $Z_1(x) = Y(x_{D_1})$ . 3D kovariance je pak  $E\{Z_1(x) * Z_1(x+h)\} = E\{Y(x_{D_1}) * Y(x_{D_1}+h_{D_1})\} = C_{(hD_1)}$ .

Předpokládejme, že sledujeme kovnatost jistého ložiska  $z_o(x)$ . Pomocí simulace získáme numerický model ložiska s hodnotami  $z_s(x)$ . Pro něj platí, že rozptyl  $z_o(x)$  a  $z_s(x)$  je stejný. Zásadní výhodou takového numerického modelu je, že díky sadě simulací máme k dispozici v každém místě sadu např. 500 simulovaných hodnot, takže jsme schopni v každém místě určit průměr, rozptyl a další statistické charakteristiky (při nejmenším do řádu 2).

**Podmíněná stochastická simulace** při modelování pole využívá i sady naměřených dat a zajišťuje, že v místech zjištěných hodnot bude mít výsledné, simulované pole stejnou hodnotu. Výsledkem je (nespojité podle Burrougha) pole hodnot  $z_{sc}(x)$ . Podmíněná simulace může být vylepšena využitím dalších druhů kvalitativních informací (např. existence zlomů, zaplavovaná území).

V každém bodě  $x$  však není  $z_{sc}(x)$  nejlepším odhadem  $z_o(x)$ . Rozptyl  $z_{sc}(x)$  je 2x větší než krigovací rozptyl (potvrzení podcenění a značného vyhlazení krigovacího rozptylu).

Pro simulační výpočet odhadu se používá technika jednoduchého krigování (kapitola 4.2.3) (*simple kriging*) s konstantním průměrem hodnoty.

#### Princip podmíněné simulace:

Regionalizovaná proměnná je realizací stacionárního náhodného pole  $z_o(x)$  s očekávanou hodnotou  $m$  a kovariancí  $c(h)$  nebo semivariogramu  $\gamma(h)$ . Charakteristickou vlastností krigování je, že krigovací chyba  $[z_o(x) - z^*_{ok}(x)]$  je ortogonální ke krigovaným hodnotám tj.

$$E\{z^*_{ok}(y) * [z_o(x) - z^*_{ok}(x)]\} = 0$$

Vyjdeme z rovnice

$$z_o(x) = z^*_{ok}(x) + [z_o(x) - z^*_{ok}(x)]$$

Výraz  $[z_o(x) - z^*_{ok}(x)]$  neznáme, ale lze ho nahradit v případě podmíněné simulace chybou  $[z_s(x) - z^*_{sk}(x)]$ .

Pak

$$z_{sc}(x) = z^*_{ok}(x) + [z_s(x) - z^*_{sk}(x)]$$

kde  $z^*_{sk}(x)$  je odhad hodnoty  $z_o(x)$  pomocí jednoduchého krigování

Rozptyl podmíněné simulace se vypočte jako

$$E\{[z_o(x) - z_{sc}(x)]^2\} = E\{[z_o(x) - z^*_{ok}(x)]^2\} + E\{[z_s(x) - z^*_{sk}(x)]^2\} = 2E\{[z_o(x) - z^*_{ok}(x)]^2\}$$

Jde tedy o dvojnásobek krigovacího rozptylu.

Prakticky:

Uvažujme N linií  $D_1, D_2, \dots, D_N$  odpovídající směrům jednotlivých vektorů  $k_1, k_2, \dots, k_N$  rovnoměrně rozložených v kouli. Na každé linii  $D_i$  je realizace  $y(x_{D_i})$  generována z náhodného pole  $Y(x_{D_i})$ .

Konečná hodnota v bodě  $x$  je dána sumací N příspěvků z N linií:

$$z_s(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N z_i(x)$$

N linií může být odvozeno jedno z druhé pomocí rotace, proto se metoda označuje jako metoda rotujících pásem. Výsledek  $z_s(x)$  je realizací 3D náhodného pole  $z_s(x) = z_s(h, v, m)$ , který je stacionární druhého řádu, s nulovou střední hodnotou a kovariancí rovné  $C(h)$ . Pokud se N blíží k nekonečnu,  $C(h)$  směřuje k izotropické kovarianci  $C(r)$ .

$$C(r) = \frac{1}{2\pi} \int C(h, k) dk$$

Integrace probíhá přes 1/2 sférické jednotky .

Kde  $(h, k)$  je projekce vektoru  $h$  do osy  $k$ .

$$r = |h| = \sqrt{h_h^2 + h_v^2 + h_w^2}$$

Izotropickou kovarianci  $C(r)$  lze převést přes sférické souřadnice na tvar

$$C(r) = \frac{1}{r} \int_0^r C(s) ds$$

V praxi se používá  $C(r)$  konstantní ve 3D a 1D kovariance je simulována v každé linii derivací

$$C(s) = \frac{\partial}{\partial s} sC(s)$$

*Poznámka 1:*

Obecně nemusí být střední hodnota nulová. Proto přidáváme konstantu  $m$  ke každé simulované hodnotě  $z_s(x)$ . Tato aditivní procedura může být využita pro simulaci nestacionárního náhodného pole, interpretovaného jako součet driftu a stacionárního zbytku  $Z(x) = m(x) + R(x)$ . Stacionární zbytek  $R(x)$  s nulovou střední hodnotou a známou kovariancí je simulován opět metodou rotujících pásem. V každém bodě je připojena simulace driftu  $m(x)$ . Simulace  $m(x)$  může být založena na rozložení  $m(x)$  do řady známých funkcí (zpravidla polynomické funkce).

*Poznámka 2:*

Dosud uvažovaná 3D kovariance byla izotropická. Každou anizotropní kovarianci však můžeme složit ze sady izotropních modelů menší dimenze (3D + 2D + 1D). Pak se provádí simulace jednotlivých komponent nezávisle a následně se sčítají.

Např. anizotropická kovariance typu  $C(h_u, h_v, h_w) = K_0 C_0(r) + K_1 C_1(h_w) + K_2 C_2(h_u^2 + h_v^2)^{1/2} + K_3 C_3(r)$

$K_0, K_1, K_2, K_3$  kladné konstanty

$C_0, C_1, C_2, C_3$  izotropické kovariance definované v 3D, 1D, 2D, 3D.

Tyto kovariance mají např. sférický model s dosahy  $a_0, a_1, a_2, a_3$  a stejným prahem. Pak

$$z(s) = \sum_{i=0}^3 T_i(x)$$

$T_0(x)$  charakterizuje izotropický nugetový efekt (mikrostruktury s  $C_0(r)$  a dosahem  $a_0$ , který je velmi malý vzhledem k experimentální vzdálenosti a simulační vzdálenosti).

Všechny  $T_0(x)$  jsou prostorově nekorelovány a mají nulovou střední hodnotu a rozptyl  $E(T_0^2(x))=K_0$ .  $T_0(x)$  může být snadno simulováno výpočtem náhodných čísel z distribuce s rozptylem  $K_0$ , nemusí se tedy použít metoda rotujících pásem.

$T_1(x)$  - charakterizuje 1D variabilitu s kovariancí  $K_1C_1(h_w)$ . Tato struktura může často odpovídat vertikální změně (např. změna obsahu ve vertikálním směru, rozdíly dané sedimentací), přitom jsou v horizontální rovině všechny  $T_1(x)$  stejné. Mohou být simulovány podél vertikální linie nebo je linie odvozena z dřívější 3D simulace s izotropickou kovariancí. Dosah  $a_1$ .

$T_2(x)$  - 2D izotropická variabilita s kovariancí  $K_2C_2(h_u^2+h_v^2)^{1/2}$ . Často jde o změny v horizontální rovině. Může být získán simulací na horizontální rovině nebo odvozena ze souboru 3D simulace s izotropickou kovariancí  $K_2C_2(r)$ . Dosah  $a_2$ , u sedimentárních struktur výrazně větší než  $a_1$  nebo  $a_3$ .

$T_3(x)$  - 3D struktura s izotropickou variabilitou s kovariancí  $K_3C_3(r)$ . Může odpovídat variabilitě malého dosahu  $a_3$ . Mění se ve všech směrech a vyžaduje metodu rotujících pásem.

*Poznámka 3: Ikosaedrická aproximace*

Metoda rotujících pásem předpokládá nekonečný počet linií  $D_i$ . V praxi se používá cca 100 linií. Vhodnou aproximací, která je méně nákladná, je použití 15 linií spojujících střední body protějších hran pravidelného ikosaedru (dvacetistěn). Výsledkem je

$$z(s) = \frac{\sum_{i=1}^{15} z_i}{\sqrt{15}}$$

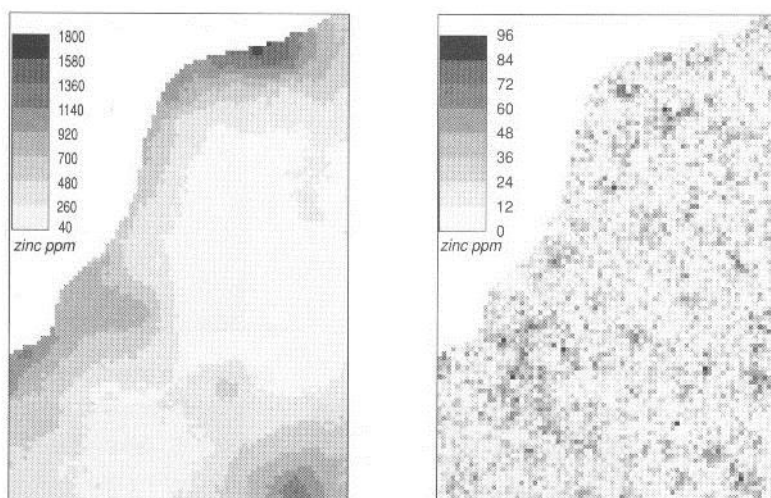
### 4.3.1 Doporučený postup pro podmíněnou stochastickou simulaci

**Možný obecný postup provádění simulace (podle Burrough, McDonell, 1998):**

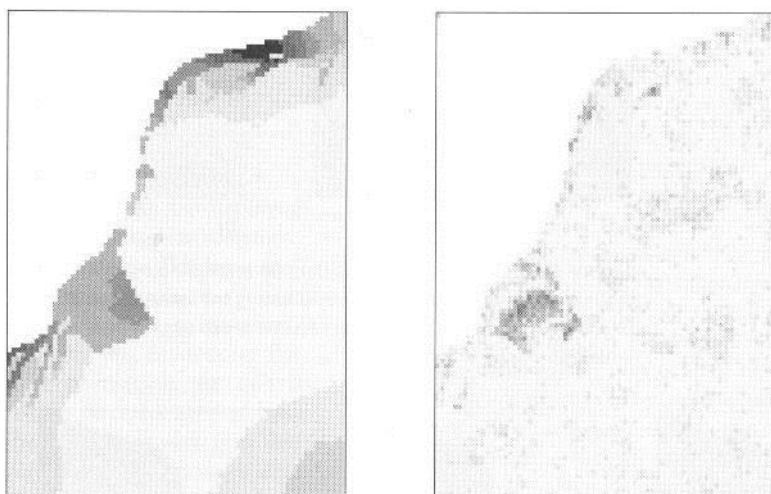
- 1) Vyberte náhodně dosud nezpracovaný bod (neznámé hodnoty)
- 2) Vypočítejte pro něj pomocí krigování odhad hodnoty a rozptyl odhadu s využitím známých (naměřených nebo již vypočítaných) hodnot
- 3) Stanovte náhodnou hodnotu z distribuce pravděpodobnosti určené zjištěným průměrem (zjištěný odhad hodnoty) a rozptylem (zjištěný rozptyl odhadu). Výsledek představuje simulovanou hodnotu pro dané místo. Bod se stává místem se známou hodnotou a vstupuje do dalších výpočtů.
- 4) Opakujte první 3 kroky, až je pokryto celé území.
- 5) Opakujte první 4 kroky tolikrát, kolik simulací je potřebné provést pro dostatečně věrohodný odhad modelu (např. 100 krát).
- 6) Ze simulovaných hodnot vypočítejte průměrnou hodnotu a rozptyl (resp. směrodatnou odchylku) pro každé místo.

Podmíněná stochastická simulace je značně výpočetně náročná, ale může poskytnout cenné a realistické výsledky. Doporučuje se provádět více než 500 simulací, kdy se získávají výsledky podobné jako při základním krigování. Při méně než 500 simulacích se objevují odchylky především při stanovování směrodatné odchylky.

Podmíněná simulace může také využívat jiných informací ze zkoumaného území, zvláště tzv. soft informací (tak označujeme zpravidla neměřená, ale odvozená, klasifikovaná data apod.).



Obr. 4-42 Podmíněná simulace obsahu Zn v půdě po 100 iteracích - vlevo očekávaná hodnota Zn, vpravo chyba odhadu. (Burrough, McDonell 1998)



Obr. 4-43 Podmíněná simulace obsahu Zn v půdě po 100 iteracích s využitím informace o pravděpodobné hodnotě obsahu Zn v jednotlivých třídách četnosti záplav - vlevo očekávaná hodnota Zn, vpravo chyba odhadu. (Burrough, McDonell 1998)

Pokud bychom měli srovnat výsledky odhadů (zvláště krigování) a simulace, nejsou kompatibilní. Kdybychom vykreslili profil přes sledované území, zjistíme, že krigovací odhad  $z^*(x)$  je bližší průběhu reálné křivky z hodnot  $z_o(x)$ , ale profil ze simulovaných hodnot  $z_{sc}(x)$  je lepší reprodukcí fluktuaace reálné křivky. Tedy krigování poskytuje vhodné odhady pro lokální a globální odhady hodnot, simulace jsou lepší pro odhad rozptylu.

## 5 Linie

Linie reprezentují na mapách 2 příbuzné fenomény:

- a) reprezentují a lokalizují skutečně lineární geografické fenomény (řeky, silnice, potrubí)
- b) rozdělují plochy a povrchy (hraniční linie, lomové linie)

Z hlediska analýzy liniového vzorku je první typ zajímavější. Statistické metody pro zkoumání liniového vzorku jsou ale méně často využívány ve srovnání s analýzami bodů nebo areálů.

Vybrané metody se používají ke statistickému popisu liniového vzorku (kapitola 5.7) a k určení jeho náhodnosti.

Pro popis liniového vzorku je nutno si uvědomit subjektivní charakter záznamu "přírodní" linie.

Linie jsou popisovány páry souřadnic, mezi 2 páry souřadnic je linie rovná. Většina linií na mapě (i při velkém měřítku) vykazuje určité zakřivení. Např. silniční linie mezi 2 městy není rovná, ale zakřivuje se, často podle jiných geoprvků. Proto je zavádějící měřit přímou liniovou vzdálenost na meandrující řece. Citlivější přístup využívá měření podél středu kanálu řeky paralelně s břehy. S tím pochopitelně narůstá počet vertexů popisujících linii. Součet délky dílčích linií pak udává vzdálenost mezi 2 uzly (počátečním a koncovým). Někdy se místo měření vzdálenosti v délkových jednotkách používá cestovní čas a dopravní náklady.

Pro analýzu liniového vzorku je významný také směr a spojení. Existence spojení mezi soustavou bodů, které tvoří linii, znamená, že lokace (body) na sobě nejsou nezávislé, ale jsou spojené v určitém směru. Body spojené v určitém pořadí musí zachovávat tuto posloupnost.

**Interakční metody** umožňují analýzu interakčních dat, tj. dat, popisujících tok mezi zdroji a cíli. Na rozdíl od jiných typů dat jsou tato data typicky vztahována k páru míst (tj. 1 zdroj a 1 cíl). Tok vyjadřuje přenos materiálu, zboží, lidí, zvířat, ale také energie nebo myšlenek mezi minimálně 2 místy. Primárním cílem analýz je porozumění prostorovým tokům a jejich modelování.

Vzhledem k rozdílnému charakteru přenášené substance se také liší prostředky, kterými se tok realizuje. Velmi často se realizuje prostřednictvím dopravních sítí. Pro popis těchto sítí a různé typy optimalizačních úloh v síti se využívá především teorie grafů (kapitola 5.2). K nejběžnějším úlohám patří hledání optimální cesty (kapitola 5.3).

V analýze interakčních dat (kapitola 5.1) se uplatňují 3 hlavní sady faktorů:

- a. Prostorové vazby mezi zdroji a cíly toků
- b. Charakteristiky určující velikost toků ze zdrojů
- c. Charakteristiky určující velikost toků do cílů

K praktickým úlohám patří hodnocení dopravní dostupnosti (kapitola 5.4), lokalizační a alokační úlohy (kapitola 5.5) a aplikace gravitační teorie (kapitola 5.6).

Z hlediska geoinformatiky je na prvním místě studium **prostorových vazeb** mezi zdroji a cíli. Vazby mezi zdroji a cíli mohou být vyjadřovány pomocí struktur liniových či vektorových, zpravidla bez důrazu na topografickou přesnost. Takové „nelokalizované“ linie zaznamenávají schématické spojení ze zdroje do cíle, u kterého nesledujeme trasu toku, ale pouze skutečnost existence (a velikosti, případně struktury) toku.

Vizualizace interakčních dat (kapitola 5.7) a tedy nejjednodušší interakční metody jsou spojeny s **liniovým vzorkem**, jeho reprezentací a popisem. Uplatňuje se především ve formě pohybové linie.

### 5.1 Analýza interakčních dat

V analýze interakčních dat se uplatňují 3 hlavní sady faktorů:

1. Prostorové vazby mezi zdroji a cíly toků (kapitola 5.1.1)

2. Charakteristiky určující velikost toků ze zdrojů (kapitola 5.1.2)
3. Charakteristiky určující velikost toků do cílů (kapitola 5.1.3)

### 5.1.1 Prostorové vazby mezi zdroji a cíly toků

Velikost separace mezi zdroji a cíli se popisuje pomocí obecně chápané vzdálenosti. Vzdálenost může být vyjadřována prostou euklidovskou vzdáleností (přímá vzdálenost mezi objekty), sférickou vzdáleností (smysl pouze u velkých vzdáleností), stále častěji se ale vyjadřuje pomocí cestní vzdálenosti v dopravní síti nebo pomocí času či nákladů potřebných k překonání vzdálenosti, eventuelně se používají jiné míry vyjádření vzdálenosti.

Existuje tedy řada variant popisu obecné vzdálenosti. Ve všech případech se ale očekává, že s rostoucí vzdáleností klesá interakce mezi zdrojem a cílem a klesá tak i velikost toku. Způsob popisu závislosti velikosti toku na vzdálenosti je opět často předmětem modelování, které hledá pro konkrétní případ odpovídající závislost.

### 5.1.2 Charakteristiky určující velikost toků ze zdrojů

Charakteristiky určující velikost toků ze zdrojů popisují schopnost zdroje generovat tok. V jednoduchých modelech může být vyjádřena i 1 proměnnou, kterou podle charakteru úlohy může být např. velikost ohrožené populace ve zdravotnické aplikaci nebo analýza objemu známých nebo očekávaných výdajů na nákupy. V základním gravitačním modelu je tato vlastnost popisována parametrem  $\lambda$ . Ve složitějších modelech (např. s omezením cílů) je možné tento parametr nahradit funkcí, která nám dovoluje namísto konstantních parametrů sady zdrojů měnit jejich produktivitu v závislosti na vybavenosti, demografickém růstu.

Předpokládá se, že mezi faktorem popisujícím zdroj a velikostí generovaného toku je přímá úměrnost.

### 5.1.3 Charakteristiky určující velikost toků do cílů

Charakteristiky určující velikost toků ze zdrojů do cílů popisují atraktivnost cílů, tedy jeho schopnost přitahovat toky. Opět mohou být vyjádřeny jedním parametrem („atraktivností“), jehož význam může být např. počet lůžek ve zdravotnických aplikacích nebo velikost prodejní plochy v marketingových aplikacích. V úvahu přichází také funkce, která popisuje závislost atraktivnosti na řadě jiných proměnných.

I v tomto případě musí platit přímá úměrnost mezi hodnotou atraktivnosti a velikostí přitahovaných toků.

## 5.2 Teorie grafů a její aplikace pro analýzu interakčních dat

Teorie grafů představuje dnes již samostatně rozvinutou matematickou disciplínu, jejíž aplikace nacházíme v řadě oblastí. Široce se uplatňuje např. v operačním výzkumu, významné místo ale zaujímá v prostorové analýze, a to především interakčních dat. Postupně se však uplatňuje i pro topologický popis a analýzu areálů.

Popis teorie grafů lze nalézt v Dudorkin (1997) nebo Veverka (1989), základní pojmy a možnosti vyjádření grafů také např. v Schejbal (1993), kde je použita mírně odlišná terminologie. Protože terminologie teorie grafů není všeobecně známá ani jednotná, je potřebné uvést a vysvětlit některé pojmy.

Na prvním místě je potřebné vymežit samotný pojem **graf**. Za obecnou můžeme považovat definici Veverky (1989), který vymezuje graf jako matematickou strukturu modelující skutečnost, že v nějaké množině prvků existují vazby. Pro prvky se používá označení uzly, pro vazby označení hrany.



Grafy jsou tedy chápány jako jistý matematický útvar, jehož základními konstruktory jsou **uzly** a **hrany**. Hrany spojují uzly a skutečnost, že hrana vstupuje nebo vystupuje z uzlu, se označuje jako **incidence**.

Tab. 5-1 Základní typologie hran a posloupnosti hran a uzlů podle Dudorkina (1997)

Pojem	Vysvětlení, příklad	
Hrany	<b>Neorientovaná hrana</b>	hrana bez vyznačení směru (bez orientace)
	<b>Orientovaná hrana</b>	hrana s vyznačením směru
	<b>Smyčka</b>	hrana, která inciduje jen s 1 uzlem
	<b>Rovnoběžné hrany</b>	hrany, které incidují se stejnými uzly
	<b>Násobné hrany</b>	Rovnoběžné hrany, které jsou všechny stejně orientované nebo ani jedna z nich není orientovaná
Sledy	<b>Sled</b> mezi uzly $u_0$ a $u_n$	uspořádaná posloupnost uzlů a hran mezi uzly $u_0$ a $u_n$
	<b>Uzavřený sled</b>	sled, kde $u_0 \equiv u_n$
	<b>Tah</b>	sled, ve kterém se každá hrana vyskytuje nejvýše 1x
	<b>Cesta</b>	tah, ve kterém se každý uzel vyskytuje nejvýše 1x
	<b>Kružnice</b>	uzavřená cesta
	<b>Orientované spojení</b>	uspořádaná posloupnost uzlů a orientovaných hran mezi uzly $u_0$ a $u_n$ , kde hrany jsou orientovány ve směru odpovídajícímu pořadí v posloupnosti
	<b>Cyklus</b>	uzavřená orientovaná cesta

Na základě typologie hran a sady hran lze vymezit základní typy grafů (tab. 5-2).

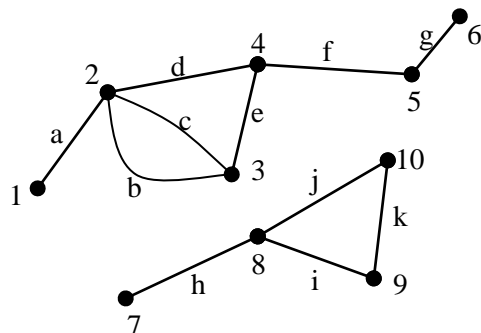
Tab. 5-2 Základní typy grafů podle Dudorkina (1997)

Označení	Vysvětlení, příklad
<b>Orientovaný graf</b>	obsahuje pouze orientované hrany
<b>Neorientovaný graf</b>	obsahuje pouze neorientované hrany
<b>Smíšený graf</b>	obsahuje orientované i neorientované hrany
<b>Prostý graf</b>	neobsahuje násobné hrany
<b>Jednoduchý graf</b>	prostý graf bez smyček
<b>Konečný graf</b>	jeho množina uzlů a hran je konečná. Většina praktických úloh řeší právě konečné grafy.
<b>Plošný (planární) graf</b>	je zobrazitelný v rovině bez protínání hran. Hrany nemají žádné společné body kromě uzlů. Neplošný graf nám dovoluje realizovat např. mimoúrovňová křížení v dopravní síti.
<b>Úplný graf</b>	Mezi každými 2 uzly existuje právě 1 hrana. Úplný neorientovaný graf o $n$ uzlech tak obsahuje $n*(n-1)/2$ neorientovaných hran.
<b>Izomorfní grafy</b>	Grafy, které se liší pouze označením uzlů a hran a způsobem zakreslení (diagramem). Kardinalita vztahu mezi uzly je 1:1 a stejně tak pro hrany.
<b>Faktor grafu</b>	část grafu, která má tutéž množinu uzlů
<b>Souvislý graf</b>	mezi libovolnými 2 uzly existuje sled
<b>Silně souvislý graf</b>	mezi libovolnými 2 uzly existuje orientovaná cesta tam i zpět
<b>Strom</b>	Souvislý graf bez kružnic
<b>Kostra grafu</b>	Faktor grafu, který je jeho stromem
<b>Cyklický graf</b>	Orientovaný graf obsahující alespoň jeden cyklus. Příkladem může být dálniční síť nebo letecká dopravní síť
<b>Acyklické grafy</b>	Orientovaný graf bez cyklu. Příkladem může být kanalizační síť.
<b>Ohodnocený graf</b>	Hrany nebo uzly jsou ohodnoceny. Ohodnocením se zaznamenávají kvalitativní a kvantitativní charakteristiky hran a uzlů. Používají se zpravidla kladná, reálná čísla. Prvním krokem ohodnocení prvků grafu je jejich indexace.

Pro popis grafů a optimalizaci v grafech se používají maticové nebo relační struktury. Především jsou:

- matice incidence
- matice sousednosti
- matice dosažitelnosti

Do **matice incidence** se zapisuje existence spojení mezi uzly a hranami (incidence). Hanuš (1992) používá pojem uzlohranová matice.



Obr. 5-1 Příklad neorientovaného grafu

**Matice incidence** pro graf na obr.5-1

		HRANY										
		a	b	c	d	e	f	g	h	i	j	k
U Z L Y	1	1	0	0	0	0	0	0	0	0	0	0
	2	1	1	1	1	0	0	0	0	0	0	0
	3	0	1	1	0	1	0	0	0	0	0	0
	4	0	0	0	1	1	1	0	0	0	0	0
	5	0	0	0	0	0	1	1	0	0	0	0
	6	0	0	0	0	0	0	1	0	0	0	0
	7	0	0	0	0	0	0	0	1	0	0	0
	8	0	0	0	0	0	0	0	1	1	1	0
	9	0	0	0	0	0	0	0	0	1	0	1
	10	0	0	0	0	0	0	0	0	0	1	1

Pro orientované grafy je možné u matice incidence používat i hodnotu -1. Tedy:

$a_{ij} = -1$  když j-tá hrana vstupuje do uzlu  $U_i$

$a_{ij} = +1$  když j-tá hrana vystupuje z uzlu  $U_i$

**Matice sousednosti** (někdy i matice spojitosti) zapisuje počet hran mezi 2 sousedními uzly.

**Matice sousednosti** pro graf na obr.5-1

		UZLY									
		1	2	3	4	5	6	7	8	9	10
U Z L Y	1	0	1	0	0	0	0	0	0	0	0
	2	1	0	2	1	0	0	0	0	0	0
	3	0	2	0	1	0	0	0	0	0	0
	4	0	1	1	0	1	0	0	0	0	0
	5	0	0	0	1	0	1	0	0	0	0
	6	0	0	0	0	1	0	0	0	0	0
	7	0	0	0	0	0	0	0	1	0	0
	8	0	0	0	0	0	0	0	1	0	1
	9	0	0	0	0	0	0	0	0	1	0
	10	0	0	0	0	0	0	0	0	1	1

**Maticе dosažitelnosti** obsahuje prvky  $b_{ij}$ , které nabývají hodnoty 1, pokud existuje sled mezi uzly  $U_i$  a  $U_j$ .

**Maticе dosažitelnosti** pro graf na obr.5-1

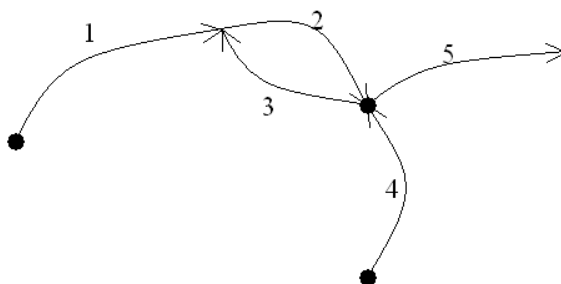
		UZLY									
		1	2	3	4	5	6	7	8	9	10
U Z L Y	1										
	2	1									
	3	1	1								
	4	1	1	1							
	5	1	1	1	1						
	6	1	1	1	1	1					
	7	0	0	0	0	0	0				
	8	0	0	0	0	0	0	1			
	9	0	0	0	0	0	0	1	1		
	10	0	0	0	0	0	0	1	1	1	

Maticе dosažitelnosti lze odvodit z matice sousednosti pomocí Booleovského součtu a součinu.

Ohodnocení hran je možno zapisovat pomocí **matice cen**. Vzdálenostní poměry (především délku nejkratší cesty mezi  $i$  a  $j$ ) lze zapisovat do **distanční matice**. Distanční matice je vyjádřením nepodobnosti objektů, na rozdíl od matice sousednosti, kterou můžeme chápat jako míru podobnosti. Distanční matici je možné transformovat do matice sousednosti (v obecném pojetí) pomocí jistých klasifikačních pravidel (podrobně je popisuje Tiefelsdorf 2000).

Pro praktickou implementaci teorie grafů pro síťové analýzy (např. pro optimalizaci cesty) je potřebné popsat odpor proti pohybu podél jednotlivých liniových segmentů. Ohodnocení hran, které popisuje velikost tření podle jednotlivých liniových segmentů, se zapisuje do tzv. **tabulky impedance**.

U grafů s mnoha cykly a rovnoběžnými hranami (např. uliční síť) je nedostatečné znát incidenci hran s uzly a impedanci segmentů (hran). Stejně důležité je vědět, zda se mohou volně pohybovat mezi jednotlivými hranami nebo zda existují nějaké bariery pro přechod (např. umístění semaforů, zákaz odbočení apod.). Musíme proto popisovat atributy párů hran nebo páry hrana-uzel. K tomu se využívá **odbočovací matice** (*turn matrix*) nebo odbočovací relace, kde se zapisuje možnost spojení (tedy odbočení) z jednoho segmentu do druhého.



Obr. 5-2 Příklad orientovaného grafu

**Odbočovací matice** pro graf na obr.5-2

	Do hrany 1	Do hrany 2	Do hrany 3	Do hrany 4
Z hrany 1	Ne	Ano	Ne	-
Z hrany 2	Ne	Ne	Ano	Ne
Z hrany 3	Ne	Ano	Ne	Ne
Z hrany 4	-	Ne	Ano	Ne

Tímto způsobem můžeme popisovat neplanární grafy, konkrétně např. dopravní síť s podjezdy, tunely, kde je křížení nekorektně sémanticky zaznamenáno jako průsečík.

**Odbočovací relace pro graf na obr.5-2**

Z hrany	Do hrany	Typ kontroly	Čas čekání (s)
1	2	Semafor	20
2	5	Dovoleno odbočit vždy	10
4	3	Semafor	30

Jak uvádí Laurini, Thompson (1994), alternativně může být atribut zaznamenán k určité kombinaci uzel-hrana.

*Praktické úlohy řešitelné v prostředí GIS pomocí teorie grafů jsou např.:*

- *Jaká je relativní přístupnost jednotlivých měst leteckými linkami?*
- *Jak bude ovlivněna doprava postavením nového mostu?*
- *Kolik alternativních cest je k dispozici při cestě z domu do práce?*

K nejčastějším úlohám optimalizace v grafu patří hledání optimálních cest (kapitola 5.3).

### 5.3 Hledání optimální cesty

Při hledání optimálních cest rozeznáváme úlohy:

- nejkratší cesty (tj. cesty s nejmenším počtem hran, s nejmenším počtem přesezení),
- nejlevnější cesty (tj. cesty s nejmenším součtem ohodnocení hran).

Někdy se nehledá pouze 1 cesta, ale více cest nejlépe splňujících dané kritérium.

Každá z těchto úloh může být dále modifikována podle toho, zda se jedná o nalezení cest:

- mezi dvojicí zadaných uzlů,
- ze zadaného uzlu do všech ostatních,
- ze všech ostatních do zadaného koncového uzlu,
- pro všechny (uspořádané) dvojice uzlů.

Pro řešení těchto úloh je k dispozici řada algoritmů. Jako příklad je možné uvést:

- Dantzigův algoritmus, který umožňuje nalézt všechny nejlevnější cesty ze zadaného uzlu do všech ostatních uzlů.
- Floydův algoritmus, který se používá pro nalezení všech nejkratších orientovaných cest mezi všemi dvojicemi uzlů (pro jednoduchý orientovaný graf).
- Dijkstrův algoritmus, který dovoluje najít nejkratší cestu a je vhodný i pro ohodnocené grafy.
- Hladový algoritmus, který hledá nejlevnější kostru grafu.
- Algoritmus pro hledání Eulerovy cesty.
- Algoritmus logického rozhodovacího stromu pro řešení úlohy obchodního cestujícího.
- Algoritmus maximálních úspor (*saving algoritmus*) pro řešení úlohy obchodního cestujícího.

Vysvětlení a praktickou ukázkou algoritmů 1, 2, 4 uvádí např. Dudorkin (1997), algoritmus 3 popisuje Raper (1993, in Tuček 1998), algoritmy 5 a 6 Veverka (1989), algoritmus 7 dokumentuje Píšek, Hanuš (1992).

#### 5.3.1 Dantzigův algoritmus

Dantzigův algoritmus umožňuje nalézt všechny nejlevnější cesty ze zadaného úhlu s do všech ostatních úhlů. Označme  $d(z)$  cenu nejlevnější cesty ze zadaného výchozího uzlu s do uzlu z. Uzly, které mají stanovenou cenu nejlevnější cesty z uzlu s, považujeme za prozkoumané, ostatní za neprozkoumané (Dudorkin (1997)).

- krok: Pro zadaný výchozí uzel položíme  $d(s)=0$  a tím jej považujeme za prozkoumaný. Vyškrtneme všechny hrany končící v s.

2. krok: Pro každý prozkoumaný uzel  $x$ , který je bezprostředním předchůdcem alespoň jednoho neprozkoumaného uzlu  $y$ , vypočteme součet jeho známé minimální ceny  $d(x)$  a ceny  $c(x,y)$  hrany  $(x,y)$  a nalezneme minimální z těchto součtů, tj.

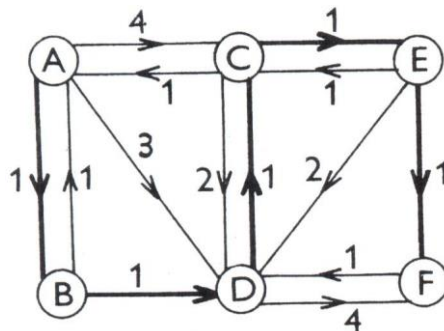
$$d(z) = \min_{x \in U_p} \left[ d(x) + c(x,y) \right]_{y \in U_n}$$

kde  $U_p$  značí množinu prozkoumaných uzlů  $U_n$  množinu neprozkoumaných uzlů.

Uzel  $z$  je novým prozkoumaným uzlem. Jestliže  $d(z) = \min[d(x)+c(x,y)]=d(r) = \dots=d(q)$ , budou novými prozkoumanými uzly  $y, r, \dots, q$  se stejnou minimální cenou nejlacinější cesty od uzlu  $s$ .

3. krok: Zaznamenáme vstupní hranu  $(x,y)$  a z dalších úvah vyškrtíme všechny ostatní vstupní hrany do uzlu  $z$ . Jestliže byly prozkoumány všechny uzly (nebo alespoň uzly, které nás zajímají z hlediska nejlevnější cesty), je výpočet u konce, jinak následuje krok 2.

Příklad (Dudorkin 1997)



Obr. 5-3 Nejlevnější cesty v grafu Dudorkin (1997)

V orientovaném grafu zadaném výkresem na obr. 5-3 je třeba nalézt nejlevnější cestu z uzlu A do ostatních uzlů. Výchozí údaje a výsledky výpočtu přehledně zapíšeme do tabulky xxx, v jejímž sloupci jsou k příslušnému uzlu  $u$  uvedeny všechny hrany z tohoto uzlu začínající v rostoucím pořadí jejich cen. Dle popsaného algoritmu položíme  $s(A)=0$  a vyškrtíme hrany končící v uzlu A (BA,CA). Vypočteme  $d(B)=\min(0+1, 0+3, 0+4)=1$ . Zarámujeme hranu AB, kde toto minimum nastalo a vyškrtíme nezarámované hrany končící v B (zde takové nejsou). Vypočteme pro zkoumané uzly A, B minimum  $\min(0+3, 0+4, 1+1)=2=d(D)$ , zarámujeme hranu BD a vyškrtíme všechny ostatní hrany končící v uzlu D (AD, CD, ED, FD). Opakováním tohoto postupu nalézáme nejlacinější cesty z uzlu A do všech ostatních uzlů. Jejich ceny  $D(z)$  a průběh jsou snadno zjistitelné z tab. 3.1. Nejlevnější cesty jsou vyznačeny v obr. 3.27 silně.

Tab.: Výpočet nejlevnějších cest

Uzel z	A	B	C	D	E	F
$d(z)$	0	1	3	2	4	5
$c(z,y)$	<b>AB(1)</b> <del>AD(3)</del> <del>AC(4)</del>	<del>BA(1)</del> <b>BD(1)</b>	<del>CA(1)</del> <del>CD(2)</del> <b>CE(1)</b>	<b>DC(1)</b> <del>DE(4)</del>	<del>EC(1)</del> <b>EF(1)</b> <del>ED(2)</del>	<del>FD(1)</del>

### 5.3.2 Floydův algoritmus

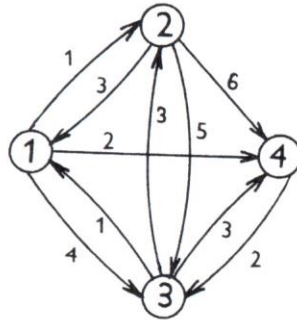
Nalezení všech nejkratších orientovaných cest mezi všemi dvojicemi uzlů jednoduchého orientovaného grafu umožňuje Floydův algoritmus (Dudorkin 1997). Je použitelný i pro grafy se záporným ohodnocením hran, avšak bez cyklů záporné ceny. Ve Floydově algoritmu postupně vyčíslujeme matice  $D_0=C, D_1, \dots, D_n$ . Prvek  $d_{ij}^{(k)}$  matice  $D_k$  ( $i=1, \dots, n; j=1, \dots, n$ ) vypočteme z prvků matice  $D_{k-1}=[d_{ij}^{(k-1)}]$  dle vztahu  $d_{ij}^{(k)}=\min(d_{ij}^{(k-1)}, d_{ik}^{(k-1)}+d_{kj}^{(k-1)})$   $i=1, \dots, n; j=1, \dots, n; k=1, \dots, n$ .

Distanční matice  $D_n=[d_{ij}^{(n)}]=[d(u_i,u_j)]=D$ , udává ceny  $d(u_i,u_j)$ . Jestliže se v některé matici  $D_k$  ( $k=1, \dots, n$ ) vyskytne záporný prvek na diagonále  $d_{ii}^{(k)}<0$ , potom v grafu existuje cyklus záporné délky

procházející uzlem  $u_i$  a výpočet je třeba přerušit. Ve Floydově algoritmu jsou postupně počítány ceny  $d_{ij}^{(k)}$  nejlevnějších cest z  $i$ -tého uzlu do  $j$ -tého uzlu, které procházejí přes uzly  $1, 2, \dots, k$  (kromě uzlů  $i$  a  $j$ ). Při vlastním výpočtu zřejmě předcházejí  $k$ -tý sloupec a  $k$ -tý řádek matice  $D_{k-1}$  beze změny do matice  $D_k$ . Floydův algoritmus nalezne pouze ceny nejlevnějších cest. Pro zjištění jejich průběhu rozšíříme algoritmus o výpočet matic

$P_k = [p_{ij}^{(k)}]_n^n$ ,  $k=1, \dots, n$ ,  $P_0 = [p_{ij}^{(0)}=1]_n^n$ . Prvek  $p_{ij}^{(k)} = p_{ij}^{(k-1)}$ , platí-li  $d_{ij}^{(k)} = d_{ij}^{(k-1)}$  a  $p_{ij}^{(k)} = p_{kj}^{(k-1)}$ , platí-li  $d_{ij}^{(k)} = d_{ik}^{(k-1)} + d_{kj}^{(k-1)}$ . Nastanou-li oba tyto případy současně, zvolíme libovolný z nich. Prvek  $p_{ij}^{(n)}$  matice  $P_n$  udává číslo uzlu bezprostředně předcházejícího uzlu  $j$  na nejlevnější cestě z  $i$  do  $j$ . Z matice  $P_n$  je tedy možno snadno sestavit průběh nejlevnější cesty mezi dvěma uzly. Existuje-li mezi dvěma uzly více nejlevnějších cest, je v matici  $P_n$  zaznamenána pouze jedna z nich.

Příklad (Dudorkin 1997):



Obr. 5-4: Orientovaný graf (Dudorkin 1997)

Mějme orientovaný graf zadaný diagramem na obr. 5- s maticí cen  $C = \begin{bmatrix} 0 & 1 & 4 & 2 \\ 3 & 0 & 5 & 6 \\ 1 & 3 & 0 & 3 \\ M & M & 2 & 0 \end{bmatrix}$

Z matice  $C=D_0$  postupně vypočteme matice  $D_1, \dots, D_4$  a  $P_0, \dots, P_4$  dle uvedeného postupu.

$$D_0 = \begin{bmatrix} 0 & 1 & 4 & 2 \\ 3 & 0 & 5 & 6 \\ 1 & 3 & 0 & 3 \\ M & M & 2 & 0 \end{bmatrix}$$

$$P_0 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 \\ 4 & 4 & 4 & 4 \end{bmatrix} \quad k=0$$

$$D_1 = \begin{bmatrix} 0 & 1 & 4 & 2 \\ 3 & 0 & 5 & 5 \\ 1 & 2 & 0 & 3 \\ M & M & 2 & 0 \end{bmatrix}$$

$$P_1 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 1 \\ 3 & 1 & 3 & 3 \\ 4 & 4 & 4 & 4 \end{bmatrix} \quad k=1$$

$$D_2 = \begin{bmatrix} 0 & 1 & 4 & 2 \\ 3 & 0 & 5 & 5 \\ 1 & 2 & 0 & 3 \\ M & M & 2 & 0 \end{bmatrix}$$

$$P_2 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 1 \\ 3 & 1 & 3 & 3 \\ 4 & 4 & 4 & 4 \end{bmatrix} \quad k=2$$

$$D_3 = \begin{bmatrix} 0 & 1 & 4 & 2 \\ 3 & 0 & 5 & 5 \\ 1 & 2 & 0 & 3 \\ 3 & 4 & 2 & 0 \end{bmatrix}$$

$$P_3 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 1 \\ 3 & 1 & 3 & 3 \\ 3 & 3 & 4 & 4 \end{bmatrix} \quad k=3$$

$$D=D_4 = \begin{bmatrix} 0 & 1 & 4 & 2 \\ 3 & 0 & 5 & 5 \\ 1 & 2 & 0 & 3 \\ 3 & 4 & 2 & 0 \end{bmatrix}$$

$$P=P_4 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 1 \\ 3 & 1 & 3 & 3 \\ 3 & 3 & 4 & 4 \end{bmatrix} \quad k=4$$

V maticích jsou vyznačeny řádky a sloupce, které přecházejí do následující matice beze změny a z nichž se počítají v případě matice D součty  $d_{ik}+d_{kj}$ . Matice Dosahuje ceny nejlevnějších cest. Všimněte si např. výpočtu prvku  $d_{24}^{(1)} = \min(d_{24}^{(0)}; d_{21}^{(0)} + d_{14}^{(0)}) = \min(6; 3+2) = 5$ . Protože minimum nastalo pro součet  $3+2$ , bude příslušný prvek matice  $P_1$  roven  $p_{24}^{(1)} = p_{14}^{(0)} = 1$ . Z matice D zjišťujeme např. cenu nejlevnější cesty z uzlu 1 do uzlu 4  $d_{14} = 3$ . Z matice P zjišťujeme, že uzlu 4 na nejlevnější cestě 1-4 předchází uzel 3, neboť  $p_{14} = 3$ . Uzlu 3 na nejlevnější cestě 1-3 předchází uzel 1, neboť  $p_{13} = 1$ . Tedy nejlevnější cesta z 1 do 4 vede přes uzel 3.

### 5.3.3 Dijkstrův algoritmus

Dijkstrův algoritmus je konečný, protože v každém průchodu jeho cyklem se do množiny navštívených uzlů přidá právě jeden uzel. Průchodů cyklem je nejvýše tolik, kolik má graf vrcholů. V grafu se nesmí vyskytnout záporné hodnoty.

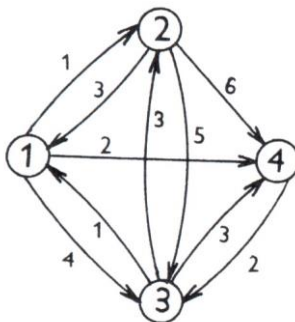
Mějme graf G, v němž hledáme nejkratší cestu. Řekněme, že V je množina všech vrcholů grafu G a množina E obsahuje všechny hrany grafu G. Algoritmus pracuje tak, že si pro každý vrchol v z V pamatuje délku nejkratší cesty, kterou se k němu dá dostat. Označme tuto hodnotu jako  $d[v]$ . Na začátku mají všechny vrcholy v hodnotu  $d[v] = \infty$ , kromě počátečního vrcholu s, který má  $d[s] = 0$ . Nekonečno symbolizuje, že neznáme cestu k vrcholu.

Dále si algoritmus udržuje množiny Z a N, kde Z obsahuje už navštívené vrcholy a N dosud nenavštívené. Algoritmus pracuje v cyklu tak dlouho, dokud N není prázdná. V každém průchodu cyklu se přidá jeden vrchol vmin z N do Z, a to takový, který má nejmenší hodnotu  $d[v]$  ze všech vrcholů v z N.

Pro každý vrchol u, do kterého vede hrana (označme její délku jako  $l(vmin, u)$ ) z vmin, se provede následující operace: pokud  $(d[vmin] + l(vmin, u)) < d[u]$ , pak do  $d[u]$  přiřad' hodnotu  $d[vmin] + l(vmin, u)$ , jinak neprováděj nic.

Až algoritmus skončí, potom pro každý vrchol v z V je délka jeho nejkratší cesty od počátečního vrcholu s uložena v  $d[v]$ . ([www.wikipedia.cz](http://www.wikipedia.cz))

Příklad řešení:



Hledejme cestu z uzlu 3.

uzel u	hodnota cesty d(u) z uzlu č.3
1	$\infty$
2	$\infty$
3	0
4	$\infty$

$v_{min}=3$ , testujeme možnost přímého cestování z uzlu číslo 3

Množina Z krok 1

uzel u	hodnota cesty d(u) z uzlu č.3
3	0

Množina N krok 1

uzel u	hodnota cesty d(u) z uzlu č.3
1	$\infty$ 1
2	$\infty$ 3
4	$\infty$ 3

$v_{min}=1$ , testujeme možnost přesedání v uzlu číslo 1

Množina Z krok 2

uzel u	hodnota cesty d(u) z uzlu č.3
3	0
1	1

Množina N krok 2

uzel u	hodnota cesty d(u) z uzlu č.3
2	<del>5</del> 2
4	3

$v_{min}=2$ , testujeme možnost přesedání v uzlu číslo 2 (žádná změna,  $1+2=3$ , která tam už je

Množina Z krok 3

uzel u	hodnota cesty d(u) z uzlu č.3
3	0
1	1
2	2

Množina N krok 3

uzel u	hodnota cesty d(u) z uzlu č.3
4	3

$v_{min}=4$ , testujeme možnost přesedání v uzlu číslo 4, ale již není žádný neznámý uzel

Množina Z krok 4

uzel u	hodnota cesty d(u) z uzlu č.3
3	0
1	1
2	2
4	3

Pro řídké grafy je vhodné použít Dijkstrův algoritmus jehož časová složitost se blíží  $O(V * E)$ , Floydův algoritmus (časová složitost výpočtu Floydovým algoritmem je  $O(V^3)$ ) je naopak vhodný pro husté grafy i díky jeho jednoduché implementaci.

Podstatný rozdíl mezi algoritmy spočívá ve schopnosti Floydova algoritmu provést výpočet i nad grafem se zápornými váhami hran (ale bez záporných cyklů), kdežto u Dijkstrova algoritmu v některých případech nikoliv



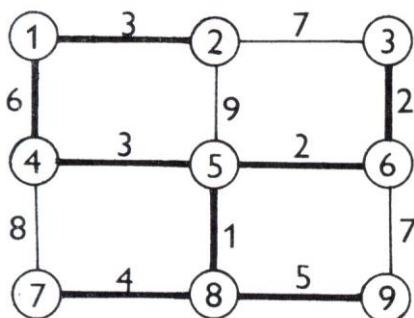
### 5.3.4 Základní algoritmus pro nalezení nejlevnější kostry

Mějme dán úplný neorientovaný graf  $G=(U,H,R)$  a matici cen  $C \geq 0$ . Máme nalézt takovou jeho kostru  $G_1=(U,H_1,R_1)$ , která je nejlevnější, tj. jejíž součet ohodnocení všech jejích hran je minimální ze všech možných koster. Metoda k nalezení této nejlevnější kostry je poměrně jednoduchá a popíšeme ji slovně (Dudorkin 1997):

1. krok: Zvolíme libovolný uzel grafu  $G$  a spojíme jej s cenově „nejbližším“ uzlem.
2. krok: Nalezneme dosud „nepřipojený“ uzel, který je cenově „nejblíže“ k některému z již „připojených“ uzlů a spojíme jej s ním. 2.krok opakujeme tak dlouho, dokud nebudou propojeny všechny uzly.

### 5.3.5 Hladový algoritmus

Nejlevnější kostru je možno také nalézt tzv. hladovým algoritmem (Dudorkin 1997). Seřadíme hrany grafu vzestupně dle jejich cen. V tomto pořadí je postupně vybíráme, přičemž pro každou vybíranou hranu kontrolujeme, zda netvoří s již vybranými hranami kružnici. Tvoří-li kružnici, zamítneme ji; jinak ji zařadíme do seznamu vybraných hran. Postup opakujeme tak dlouho, dokud nezískáme kostru grafu. Lze ukázat, že je nejlevnější. V příkladu na obr. 5-5 zařazujeme hrany v pořadí 5-8, 5-6, 3-6, 4-5, 1-2, 7-8-, 8-9-, 1-4. Přidání dalších hran by již vedlo ke vzniku kružnice.



Obr. 5-5 Nejlevnější kostra grafu (Dudorkin 1997)

## 5.4 Způsoby hodnocení dopravní dostupnosti

Jedním z nástrojů zkoumání prostorové interakce mezi zdroji a cíli a popisování prostorových vazeb je i průzkum a popis dopravní dostupnosti.

Pojem dostupnosti geografických objektů byl rozpracován již v 50. a 60. letech. K popisu dostupnosti se používají různé míry dostupnosti, významně se uplatňuje teorie grafů (kapitola 5.2).

**Dostupnost** je chápána jako určitý ukazatel, který na základě přístupnosti nebo dosažitelnosti daného objektu k ostatním objektům určuje jeho postavení v rámci dané prostorové struktury (Kusendová 1996). Někteří autoři používají pojem akcesibilita (Maryáš et al. 1995).

Dostupnost je chápána jako geografický pojem, geografická charakteristika objektu. Stanovuje se na základě vzdálenostních charakteristik v rámci bodové nebo liniové struktury.

Míry dostupnosti dovolují popisovat dostupnost geografických objektů a uplatňují se především v socioekonomické geografii.

Míry dostupnosti můžeme podle použité metriky (zjednodušeně podle použitých jednotek) dělit na:

- 1) Metrické (kapitola 5.4.1)
- 2) Časové (kapitola 5.4.2)
- 3) Topologické (kapitola 5.4.3)

- 4) Cenové (nákladové) (kapitola 5.4.4)
- 5) Ostatní

Příkladem posledně jmenované míry je např. fyziologický index únavnosti, který zahrnuje informaci o délce pěších cest k dopravnímu prostředku a hodnotí i jakost přepravy - její únavnost (Hůrský 1969).

Vedle základních měr dostupnosti se uplatňují i vážené míry dostupnosti (kapitola 5.4.5).

Dostupnost lze dělit i podle jiných hledisek, např. podle dopravního prostředku, pro který je zjišťován. Pak bychom dělili dostupnost na určenou podle provozně organizačního hlediska pro hromadnou a individuální dopravu, podle provozně technického hlediska na veřejnou a neveřejnou dopravu. Ze všech kombinací má smysl sledovat především neveřejnou individuální dopravu a veřejnou hromadnou dopravu.

K termínu „dostupnost“ má relativně blízko pojem dopravní **obslužnost**, která je vztahována k veřejné hromadné dopravě.

### 5.4.1 Metrické míry

Metrické míry používají vzdálenosti měřené jako přímé (vzdušné, euklidovské) a nebo cestní (po komunikaci).

#### Míra přímé dostupnosti euklidovské

U této míry dostupnosti není potřebná konstrukce grafu, využívají se pouze euklidovské (vzdušné) vzdálenosti, takže ji lze snadno vypočítat ze souřadnic zkoumaných míst. Nejlepší dostupnost má místo s nejmenší hodnotou přímé euklidovské vzdálenosti, což odpovídá těžišti cílových objektů.

$$D_i^P = \sum_j d_{ij}^v$$

$D_i^P$	míra přímé dostupnosti v místě i
$d_{ij}^v$	euklidovská vzdálenost mezi místy i a j
j	index cíle

#### Míra cestní dostupnosti

Míra cestní dostupnosti používá výpočet vzdálenosti po trase přesunu, tedy vlastně délky cest v grafu.

Cestní vzdálenost se stanovuje zpravidla na základě určitého modelu dopravní sítě, jehož přesnost je závislá na měřítku a úrovni generalizace. Často se odvozuje v prostředí GIS pomocí síťových funkcí typu nejkratší cesta. Proto je výpočet nejkratší vzdálenosti mezi jednotlivými místy nutně zatížen jistou chybou. Proto je vhodnější používání označení „cestní vzdálenost“ na místo „skutečná vzdálenost“.

Nejlepší cestní dostupnost má místo s nejmenší hodnotou ukazatele:

$$D_i^C = \sum_j d_{ij}^c$$

$D_i^C$	míra cestní dostupnosti v místě i
$d_{ij}^c$	délka nejkratší cesty z místa i do j
j	index cíle

Při řešení nejkratší cesty se používají především techniky lineárního programování (viz hledání optimální cesty (kapitola 5.3).

## Jiné

Vedle základních údajů typu přímá vzdálenost a cestní vzdálenost se používají i další ukazatele typu rozvoj čáry (poměr mezi cestní a přímou vzdáleností) či koeficient okliky (o kolik % je cestní vzdálenost větší než přímá vzdálenost)

Teoreticky by bylo možné vyčlenit metody přímé sférické založené na použití sférické geometrie, avšak hodnocení dostupnosti velkých regionů, kde se již uplatní zakřivení Země, se prakticky neprovádí.

### 5.4.2 Časové míry dostupnosti

Do skupiny **časových měr** řadíme především **časovou dostupnost**. Vyjadřuje celkovou dobu cestování ze zkoumaného zdroje do všech cílů hvězdicovým způsobem. Nejlepší časovou dostupnost má potom uzel (místo) s nejmenší hodnotou časové dostupnosti.

$$D_i^t = \sum_j t_{ij}$$

$D_i^t$	míra časové dostupnosti v místě $i$
$t_{ij}$	doba nejkratšího přesunu z místa $i$ do $j$
$j$	index cíle

Zde se analogicky k cestní dostupnosti sčítají cestovní časy mezi uzly  $i$  a  $j$ . Cestovní časy můžeme chápat jako časovou vzdálenost (uvažují se potom pouze doby samostatného nejrychlejšího přesunu) nebo jako časovou ztrátu. Ta má smysl v případě využívání veřejné hromadné dopravy, kdy do doby cestování zahrnujeme i dobu čekání na odjezd dopravního prostředku.

### 5.4.3 Topologické míry dostupnosti

Topologické míry dostupnosti využívají teorie grafů (kapitola 5.2).

#### Přímá topologická dostupnost

Vyjadřuje celkový počet sousedních uzlů v grafu. Místo (uzel) s nejvyšším počtem sousedů má nejlepší přímou topologickou dostupnost.

$$D_i^U = \sum_j I_{ij}$$

$D_i^U$	míra přímé topologické dostupnosti v místě $i$
$I_{ij}$	indikátor sousedství uzlu $j$ vzhledem k uzlu $i$ (nabývá hodnoty 1 v případě existence sousedství, jinak 0, lze získat z matice sousednosti)
$j$	index cíle

#### Nepřímá topologická dostupnost

Vzdálenosti mezi uzly jsou vyjadřovány počtem hran na nejkratší cestě mezi nimi. Nejlepší nepřímou topologickou dostupnost bude mít uzel s nejmenší hodnotou ukazatele, podle teorie grafů se jedná o střed grafu, tedy o uzel s minimální excentricitou.

$$D_i^H = \sum_j d_{ij}^h$$

$D_i^H$	míra nepřímé topologické dostupnosti v místě $i$
---------	--

$d_{ij}^h$  počet hran na nejkratší cestě mezi místy i a j  
j index cíle

Pokud každý z uzlů představuje konečnou stanici dopravního prostředku, můžeme porovnávat dostupnost uzlů z hlediska minimálního počtu přesezení potřebného k cestování do ostatních uzlů sítě.

#### 5.4.4 Cenové míry dostupnosti

Cenové míry dostupnosti jsou založeny na ceně dopravy, v případě individuální dopravy na nákladech dopravy.

U veřejné hromadné dopravy se sleduje cena placená za přepravu mezi jednotlivými místy (zpravidla základní jízdné bez různých slev). V některých případech se vybírá dopravní prostředek, v jiných měřeních se povoluje přestupovat mezi prostředky. Více variant nabízí sledování individuální dopravy, kde vedle výběru dopravního prostředku se může sledovat jen spotřeba pohonné látky (přepočtená na cenu) nebo se může zahrnout i amortizace vozidla.

$$D_i^F = \sum_j c_{ij}$$

$D_i^F$  míra cenové dostupnosti v místě i  
 $c_{ij}$  cena nejlevnější přepravy z místa i do j  
j index cíle

#### 5.4.5 Vážené míry dostupnosti

Jednoduché míry dostupnosti považují všechny geografické objekty, které představují zdroje (resp. cíle) toků za rovnocenné a přidělují jim stejnou váhu. Proto prvním rozšířením uvedených základních modelů je zahrnutí atraktivity center, tedy cíle cestování.

Příkladem může být **vážená časová dostupnost** vyjádřená

$$D_i^{tv} = \frac{\sum_j t_{ij} * F_j}{\sum_j F_j}$$

$D_i^{tv}$  míra časové dostupnosti v místě i  
 $F_j$  atraktivita cíle j

Např. Jánošíková, Kubáni (2000) demonstují použití takového ukazatele dostupnosti při analýze dopravní dostupnosti obcí Žilinského kraje. Atraktivita cíle je u dojížděky do zaměstnání vyjadřována počtem dojíždějícího ekonomicky aktivního obyvatelstva a pro případ dojížděky do školy počtem dojíždějících studentů.

Podobně i další míry dostupnosti (např. cestní dostupnost) by bylo možné vážit.

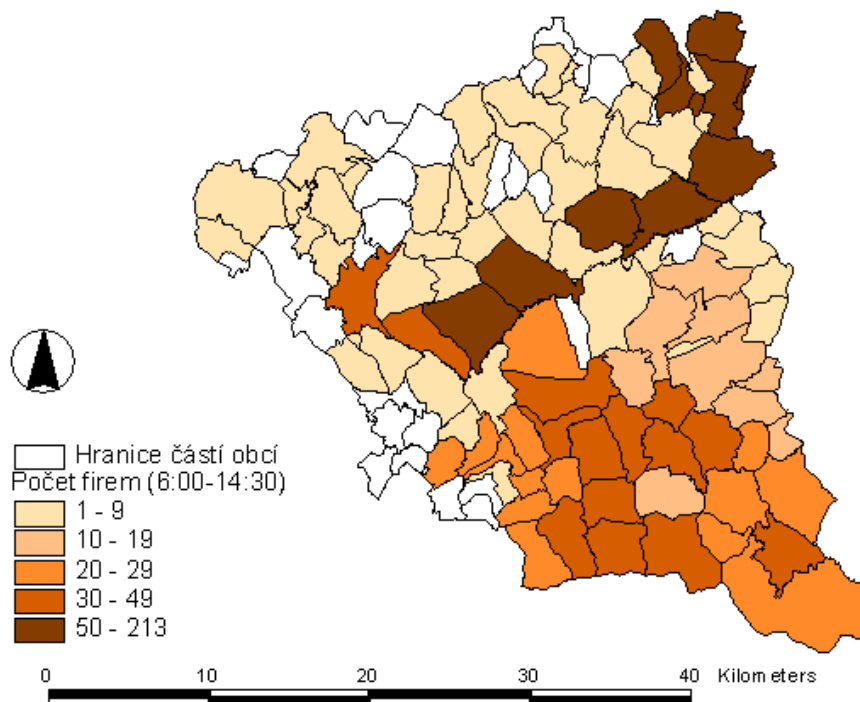
Zahrnutí atraktivity představuje první krok na přechodu k využívání zobecněných gravitačních modelů (kapitola 5.6) (resp. modelů maximalizujících entropii).

#### 5.4.6 Hodnocení veřejné dopravy

Rozšíření modelů dostupnosti může zahrnout komplexní vyhodnocení dopravy, zvláště veřejné hromadné dopravy. Zde často nestačí v praktických socioekonomických úlohách sledovat např. nejkratší dobu  $t_{ij}^{\min}$  cestování mezi 2 geografickými objekty (např. mezi 2 obcemi), ale spíše:

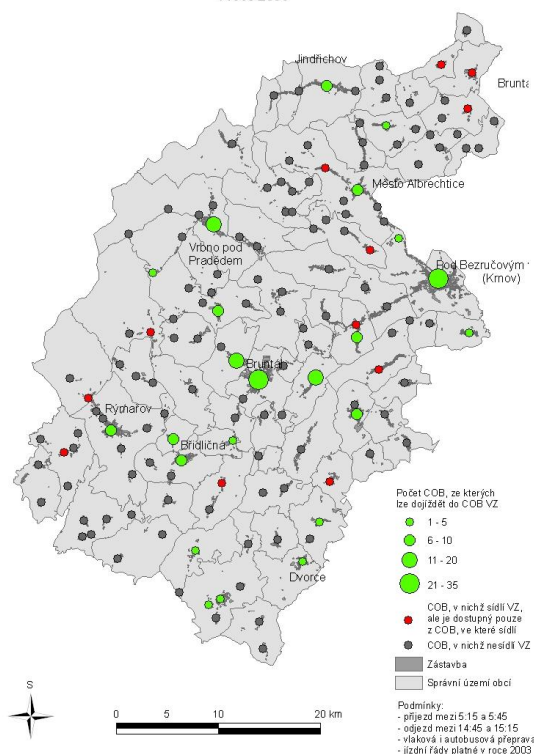
- dobu  $t_{ij}^T$  v čase  $T$  požadované přepravy (cestování do zaměstnání, do školy, na úřady, tedy všude, kde je stanoven požadovaný čas nástupu),
- počet spojení v určitý časový interval (frekvenční akcesibilita),
- počet přesezení (a s tím spojené zvýšení nepohodlí, riziko přerušení či zdržení spojení),
- komfort cestování (obsazení vozidel veřejné hromadné dopravy, jejich vybavenost) a použitelnost veřejné hromadné dopravy (bezbariérový přístup pro imobilní občany apod.).

Významnou roli mají také geografické faktory kvality dopravy - především přístupnost zastávek veřejné hromadné dopravy.

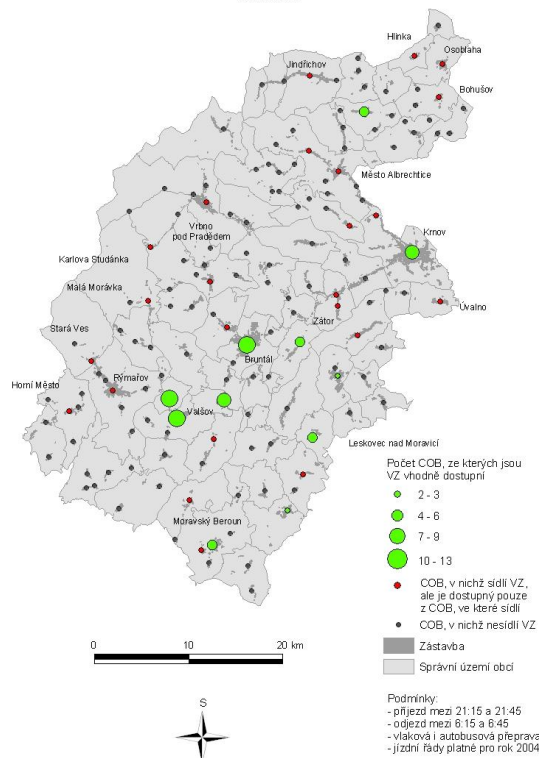


Obr. 5-6 Počet velkých firem dostupných z částí obcí v okrese Nový Jičín v roce 2000 (práce 6:00-14:30)

DOPRAVNÍ DOSTUPNOST COB VZ VHODNOU DOJÍŽDKOU NA 6:00 A ZPĚT  
v roce 2003



DOPRAVNÍ DOSTUPNOST COB VZ VHODNOU DOJÍŽDKOU NA 22:00 A ZPĚT  
v roce 2004



Obr. 5-7 Dopravní dostupnost částí obcí okresu Bruntál do zaměstnání v 6:00 a 22:00 v letech 2003 a 2004

## 5.5 Lokalizační a alokační úlohy

Se studiem interakčních dat a jejich modelováním jsou spojovány často používané lokalizační a alokační úlohy.

**Lokalizační úlohy** řeší problém optimálního umístění zařízení či objektu. Umístění je optimalizováno z hlediska optimalizačního kritéria, které závisí na charakteru úlohy a funkci zařízení. Kritérium může zahrnovat náklady na umístění zařízení (cena pozemku, výstavby a vybavení), tak především maximalizaci poskytovaných funkcí.

**Alokační úlohy** se zaměřují na problém optimálního zásobování. Existující zařízení je třeba optimálně „vybavit“, aby dobře plnila své funkce.

Je zřejmé, že oba typy úloh spolu úzce souvisí. Optimálně lze ovšem přiřazovat i úkoly jednotlivým pracovištím apod. - takové úlohy někteří autoři (především z oblasti operačního výzkumu) označují jako přiřazovací úlohy.

Lokalizační úlohy jsou často využívány socioekonomickými geografi, dříve se však zabývaly především rozmístěním výrobních kapacit ve vazbě na trh a dopravní náklady. Klasickým modelem v tomto směru je Thünenův model rozmístění jednotlivých druhů zemědělské produkce od centrálně situovaného trhu. V současnosti se lokalizační a alokační úlohy využívají při optimalizaci rozmístování nových administrativních, komerčních či obslužných objektů a optimalizaci jejich vybavení.

*Příklady lokalizačních úloh - umístění optimálního centra pro síť zákazníků, úložiště jaderných odpadů z několika elektráren, středisko údržby pro několik závodů, umístění strojů ve výrobní hale. Příkladem alokačního problému je požadavek rozmístit zboží v adekvátním množství do příslušných obchodů vzhledem k známému rozmístění obchodů a zákazníků, často s využitím dopravní sítě.*

Lokalizační úloha se zabývá volbou umístění zařízení, zpravidla ve vztahu k rozmístění zákazníků.

Pojmem **zařízení** budeme označovat objekty libovolného druhu, která poskytují jisté služby a jejichž umístění či vybavení závisí na lidské aktivitě. V typických aplikacích se jedná o obchody, zdravotnická zařízení, hasičské stanice, školy, úřady, firmy, informační centra, skládky, ale i např. krmná zařízení či útočiště pro zvířata, přírodní rezervace, zóny klidu apod.

Jako **zákazníky** označíme spotřebitele nebo uživatele služeb poskytovaných zařízeními, jejichž umístění či vybavení je optimalizováno. Mohou to být nakupující dojíždějící do obchodů, děti dojíždějící do škol, občané k úřadům, zaměstnanci do zaměstnání, pacienti k lékařům, zvířata cestující za potravou apod. Mohou to být také organizace odebírající zboží z velkoskladu, dopravní prostředky využívající čerpací stanice, výrobní stroje odebírající elektřinu z transformátoru, odpady atd.

V operačním výzkumu se často řeší úlohy, kde optimalizujeme umístění dalšího zařízení pouze v závislosti na existujících zařízeních.

Pro řešení lokalizačních úloh je nutno zvážit **podmínky** prováděné analýzy. Především je důležité, zda:

- můžeme umístit 1 nebo více nových zařízení,
- se může nebo nemůže měnit počet a velikost zařízení (dynamické versus statické řešení),
- zařízení jsou reprezentována bodem nebo areálem,
- umístění nových zařízení závisí nebo nezávisí na existujících,
- množina přípustných umístění je spojitá (nové zařízení lze umístit kamkoliv) nebo diskrétní (např. výběr z nabídky pozemků). Druhá varianta se lépe řeší s využitím teorie grafů a heuristických postupů.

Při řešení lokalizačních a alokačních úloh se uplatňuje matematické programování (lineární i nelineární), optimalizace na grafech a heuristické postupy. Některé metody neobsahují podmínku optimalizace vzdálenosti, protože postačuje pouze splnění vzdálenostního limitu (např. pro optimalizaci rozmístění zařízení naléhavé potřeby na území města při splnění limitu časové dostupnosti).

Základní krokem je volba optimalizačního kritéria. Obecně jím může být například:

- maximalizace zisku (tj. rozdílu mezi příjmy plynoucími z poskytování služeb zákazníkům a náklady na výstavbu zařízení a jejich provoz),
- maximalizace objemu toků (co nejvyšší objem poskytovaných služeb),
- minimalizace délky toků (časové, metrické či nákladové) při dosažení limitního (resp. nadlimitního) objemu toků či splnění jiného limitu (např. doba cestování nesmí překročit jistý limit),
- maximalizace efektu pro zákazníka,
- minimalizace nákladů pro zákazníka.

Klasické metody jsou deterministické a předpokládají, že toky směřují k nejbližšímu cíli. Tímto způsobem vznikají jasně vymezené a nepřekrývající se spádové oblasti kolem jednotlivých zařízení. To je vhodné pro lokalizaci administrativních zařízení nebo v optimalizaci umístění zařízení naléhavé potřeby.

Ve skutečnosti je chování toků (inicializované zákazníky) málokdy deterministické, ale spíše neurčité a je vhodné ho modelovat pomocí pravděpodobnostních a behaviorálních modelů. U nich pak vzdálenost vystupuje sice jako důležitý faktor, ale ne jako striktní omezení. Připouštějí takové rozložení toků, kde největší tok směřuje k nejbližšímu cíli, ale i vzdálenější zařízení jsou cílem jistého menšího toku z daného místa.

V některých případech může být výhodné použití etapového postupu při určování cílů, který odpovídá praktické zkušenosti, kdy v případě malých rozdílů mezi vzdálenostmi/náklady již tyto nehrají při výběru žádnou roli. Thill a Horowitz (1997) popisují dvouetapový model s tím, že v první

etapě se vybírá jistá množina kandidátů podle cestovní vzdálenosti a ceny dopravy, a teprve ve druhé etapě se určují vhodné alternativy z této skupiny kandidátů.

Optimalizačních kritérií, a tedy i prostorových interakčních modelů, může být celá řada. Uvedeme 3 základní kritéria:

#### minimalizace délky toků

Jednodušší model používá jako kritérium minimalizaci procestovaných vzdáleností v systému za předpokladu cestování zákazníků k nejbližšímu zařízení. Umísťuje se  $p$  zařízení do  $n$  uzlů (diskrétní úloha).

Kritéria:

$$\min \left\{ \sum_{i=1}^n \sum_{j=1}^n a_i * I_{ij} * d_{ij} \right\} \quad \sum_{j=1}^n I_{ij} = 1 \quad \sum_{i=1}^n I_{ij} \geq 1$$

$d_{ij}$  cena cestování z  $i$  do  $j$  (např. vzdálenost)

$a_i$  velikost požadavku na službu od zákazníka  $i$  (zjednodušeně objem požadovaných služeb v místě  $i$ ) - např. počet obyvatel, velikost určité části populace, cena majetku

$I_{ij}$  indikátor; nabývá hodnoty 1 pokud je v místě  $j$  poskytována služba pro zákazníka  $i$ , nebo hodnoty 0, pokud tomu tak není.

Druhá podmínka (suma  $I_{ij}=1$ ) musí být splněna pro všechna  $i = 1 \dots n$ , protože každý požadavek je vyřizován právě jedním zařízením.

Není omezení kapacity cíle.

Jednoduchý model předpokládá lineární závislost celkových nákladů na vzdálenosti nového zařízení od zákazníků (délky toků). Lineární závislost je samozřejmě nejsnáze řešitelná (např. metodou nejmenších čtverců). Obecněji lze uvažovat o náhradě  $d_{ij}$  za funkci  $g(d_{ij})$ , která může zahrnovat jak lineární tak nelineární modely (musí však být k dispozici vhodné řešení matematické úlohy).

Pro  $p=1$  se úloha označuje jako Weberův problém, pro obecný počet zařízení  $p$  jako  $p$ -mediánové umístění. Praktická řešení pro umístění jednoho i  $p$  zařízení uvádí např. Dudorkin (1997), který popisuje algoritmy pro euklidovskou vzdálenost, manhatonskou vzdálenost (rektilineární) i pro kvadrát euklidovské vzdálenosti, doporučovaný např. pro optimalizaci umístění hasičských stanic. Pro jednoduché přiřazovací úlohy se doporučuje tzv. maďarská metoda.

Řešení těchto úloh je k dispozici jak pro kontinuální tak pro diskrétní možnost lokalizace, pro diskrétní je ovšem algoritmy obtížná vzhledem k rozsahu reálných úloh a tak se často používají různé heuristické postupy.

Rozšíření základního modelu dovoluje zohlednit náklady umístění zařízení v daném místě, kapacitní omezení v daném místě (např. kapacita nemocnice či počet pracovních míst), omezení na maximální délku toku (maximální vzdálenost, kterou je zákazník ochoten překonat, či např. cena za dopravu, kterou je zaměstnanec ochoten při dojíždění do zaměstnání zaplatit). S rozšiřováním základního modelu dochází často k automatickému uvolnění původního požadavku na cestování do nejbližšího místa.

#### minimaxové řešení

Kritériem je minimalizace maximální vzdálenosti toků. Praktické řešení úlohy pro 1 zařízení poskytuje např. Elzingův algoritmus (Dudorkin 1997).



### maximalizace pokrytí

Kritériem je maximalizace počtu zákazníků při splnění limitu, kladeného na tok (čas, vzdálenost nebo náklady jsou nižší než stanovený limit)

$$\min \left\{ \sum_{i=1}^n a_i * I_i(h) \right\} \quad \sum_{j=1}^n J_{ij} \geq 1$$

- $d_{ij}$  cena cestování z  $i$  do  $j$   
 $a_i$  velikost požadavku na službu od zákazníka  $i$   
 $I_i(h)$  indikátor; nabývá hodnoty 1 pokud místo  $i$  je za limitní vzdáleností  $h$  od každého zařízení (tedy pokud není zákazník „pokryt“), nebo hodnoty 0 pokud tomu tak není  
 $J_{ij}$  indikátor; nabývá hodnoty 1 pokud není zařízením  $j$  poskytována služba pro zákazníka  $i$ , nebo hodnoty 0, pokud zařízení  $j$  poskytuje službu pro zákazníka  $i$ ; platí pro všechna  $d_{ij} \leq h$

První podmínka minimalizuje množství nepokrytých zákazníků.

Poslední uvedený model má zajímavé vlastnosti. Model je obecnější - např. pro velké  $h$  přechází v minimaxové řešení. Výhodou je rovněž možnost sledování chování v případě změny  $p$  a zjišťování např. kolika zařízení ( $p$ ) je potřeba pro dosažení určité úrovně pokrytí.

Vzhledem k výpočetní náročnosti se opět často uplatňují heuristické postupy.

### **Rozšíření základních modelů**

Tyto základní modely je možné upravit - např. vyloučení obslužených zákazníků a využití vzdálenostního hendikepu pro eliminaci vlivu zanedbatelných rozdílů ve vzdálenosti.

Další rozšíření modelu zahrnují např. **hierarchické lokalizačně-alokační modely**, které dovolují umístit více zařízení, které však neposkytují služby stejné úrovně. Příkladem mohou být zdravotnická zařízení, z nichž některé jsou jen se základním vybavením a jiná mají i specializovaná pracoviště.

Některé modely zohledňují i různou **atraktivitu** zařízení, vycházející např. z jejich velikosti (analogie gravitačních modelů).

Konečně se uplatňují i **multikriteriální modely**, kde je současně pro optimalizaci použito více kritérií, problémy jsou pak především v realizaci řešení.

V případě **dynamického řešení** se optimalizace rozšiřuje a používá se sumace vah za zkoumané časové období.

Z konstrukce modelů vyplývají minimální požadavky na vstupní data a specifikace vztahů.

V minimální variantě musí být známo:

- optimalizační kritérium,
- umístění zákazníků,
- umístění stávajících zařízení (alespoň pro statické řešení),
- způsob realizace toků (např. dopravní spojení, elektrické vedení),
- vyjadřování vzdálenosti (např. měrné náklady na realizaci spojení v závislosti na vzdálenosti) .

Diskrétní úlohy jsou zpravidla simulovány pomocí matematických grafů. V prostředí GIS se však setkáváme i s jinými typy úloh. Některé lokalizační problémy nejsou orientované na grafy.

*Např. lidské nebo přírodní toky se mohou vyskytovat mezi určitými body nebo mezi bodovými a liniovými prvky (např. plošný tok znečištění přes plochu k řece, kde končí). Takový problém může*

*vyžadovat speciální hledání cesty přes povrch s různými podmínkami, která jsou reprezentována polygony nebo buňkami mřížky nebo který může zahrnovat tvorbu virtuálních spojení mezi body. Předpověď objemů přepravovaných vzduchem (znečištění vzduchu) uvnitř města byla dříve často modelována gravitačním modelem, vyžadujícím výpočet velikosti uzlů a jejich vzdálenost v mřížce.*

Na závěr je potřebné upozornit na některé specifické problémy při praktických aplikacích.

Především je nezbytné zohlednit tzv. hraniční problém. Modely mohou být prakticky omezené hranicemi sledovaného území, které se nemusí krýt s funkčním vymezením pomocí optimalizačního kritéria (tedy se spádovou oblastí). Fyzické omezení modelu může narušovat proces optimalizace i vlastní využívání modelu.

Řada modelů využívá pro informaci o zákaznících informace o obyvatelstvu s trvalým bydlištěm v dané zóně (a zjišťují model závislosti určitého podílu zákazníků z celkového množství obyvatelstva) v závislosti na vzdálenosti od služby). Reálná situace je ovšem komplikována mobilitou zákazníků (dojíždění za prací apod.) a rovněž i změnami v demografické situaci (časový vývoj v populaci).

Na druhou stranu je třeba upozornit, že nemá smysl vytvářet zbytečně komplexní modely. Do praktického rozhodování stejně vstupují další hlediska a požadavky, které nelze nebo není efektivní zahrnout do prostorové optimalizace. Málokdy jsou modelová řešení použita ke skutečné optimalizaci, zpravidla jsou použita pro podporu rozhodování při srovnávání variant či průzkumu situace. Pro jejich správnou interpretaci je pak nezbytná jasná koncepce a konzistence modelu.

## 5.6 Gravitační teorie a její zobecnění

Požadavkem některých analýz interakčních dat je vytvoření vhodného modelu popisujícího velikost interakcí.

Obecný model pro sledování interakčních dat předpokládá, že každá zjištěná hodnota toku  $Y_{ij}$  mezi zdrojem  $i$  a cílem  $j$  se skládá z pravidelné a náhodné složky. Pravidelnou složku označíme jako  $\mu_{ij}$ , náhodnou složku (chyby) jako  $\varepsilon_{ij}$ .

Modelování se zabývá především správným popisem pravidelné složky.

Základním modelem pro popis pravidelné složky interakčních dat je gravitační model. Patří do obecné rodiny modelů maximalizujících entropii systému a dnes již existuje v celé řadě variant. Jeho podstatou je závislost velikosti interakce (tedy objemu toku) na velikosti zdroje, velikosti cíle a nepřímo úměrně na jisté míře vzdálenosti zdroje a cíle (původně nepřímo úměrně na čtverci euklidovské vzdálenosti).

### Základní model

Gravitační model popisuje např. Pavlík, Kühnl (1981) pomocí vztahu:

$$P_{ij} = \frac{M_i * M_j}{d_{ij}^b}$$

$P_{ij}$	síla vzájemného působení hmot
$M_i, M_j$	„hmoty“ v místě $i$ a $j$ , tj. velikost zdroje resp. cíle
$d_{ij}$	vzdálenost
$b$	koeficient vlivu vzdálenosti

V geografických aplikacích vystupuje v roli „hmoty“ (tj. velikosti) např. počet obyvatel, počet ekonomickou aktivních obyvatel i složitější faktory typu „počet obyvatel \* průměrný příjem“. Vzdálenost může být vyjádřena jako metrická, často se používá časová, vzdálenostní nebo cenová. Vliv vzdálenosti ( $b$ ) se mění podle typu dopravního prostředku.

Jako 2 hlavní aplikace se uvádí výpočet hraničního bodu a výpočet přitažlivosti obchodního střediska.

Základní model gravitačního zákona však nezabezpečí soulad se zjištěnými hodnotami toků. Proto se prakticky uplatňují oboustranně omezený gravitační model (kapitola 5.6.1) nebo jiné gravitační modely (kapitola 5.6.2).

Všechny tyto modely se ale zaměřují pouze na modelování efektu 1.řádu (tedy středních hodnot) a naopak často při řešení vyžadují nezávislost chyb a tedy nulový efekt 2.řádu (řešení např. metodou maximální věrohodnosti).

Problém s modelováním efektu 2.řádu je specifický pro interakční data. Projevuje se zde totiž šíření chyby, tedy zjevná tranzitivní závislost velikosti efektu 2.řádu v jednotlivých místech.

Gravitační modely se využívají např. pro modelování dojíždění. Někdy zahrnují dojížděku i migraci obyvatel. Gravitační modely (nebo techniky maximalizace entropie) dovolují řešit chybějící interakční matice (popis toků pracujících). K tomu je potřebné znát počet dojíždějících, vybrat vhodnou míru vzdálenosti (čas nebo cena) a vhodnou míru „hmotnosti“ (např. počet pracovních míst v regionu).

V demograficky orientovaných studiích se používají často Markovovy modely, které stanovují pravděpodobnost transitu pro jisté % regionální pracovní síly, která pak dojíždí do jiných regionů (Schubert et al. 1987).

### 5.6.1 Oboustranně omezený gravitační model

Základní podoba gravitačního zákona má jednu zásadní vadu a sice že připouští nekonzistenci s pozorovanými toky.

Chceme-li dodržet konzistenci modelu s pozorovanou situací, musíme zabezpečit 3 podmínky:

1. Suma toků ze zdroje **i** musí odpovídat zjištěné hodnotě

$$\sum_j \mu_{ij} = a_i$$

2. Suma toků do cíle **j** musí odpovídat zjištěné hodnotě

$$\sum_i \mu_{ij} = b_j$$

3. Celková cena cestování v systému je konstantní.

$$\sum_i \sum_j d_{ij} * \mu_{ij} = c$$

Na základě maximalizace funkce entropie lze získat obecný gravitační či prostorový interakční model:

$$\mu_{ij} = \alpha_i * \beta_j * e^{\gamma * d_{ij}}$$

$\alpha_i$  sada parametrů popisujících vlastnost zdroje **i** generovat toky  
 $\beta_j$  sada parametrů popisujících vlastnost cíle **j** přitahovat toky

$\gamma$  popisuje vzdálenostní efekt  
 $d_{ij}$  vzdálenost mezi zdrojem **i** a cílem **j**

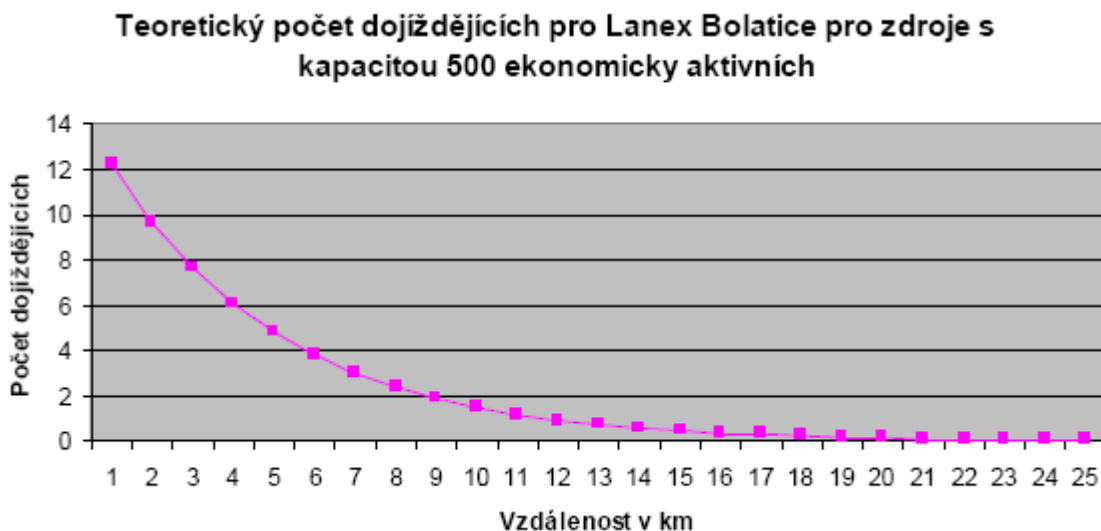
Nakonec tedy můžeme vyjádřit modelový vztah pro pozorované toky jako:

$$Y_{ij} = \alpha_i * \beta_j * e^{\gamma * d_{ij}} + \varepsilon_{ij}$$

$Y_{ij}$  zjištěná hodnota toku mezi zdrojem **i** a cílem **j**  
 $\varepsilon_{ij}$  náhodná složka pozorování (chyby)

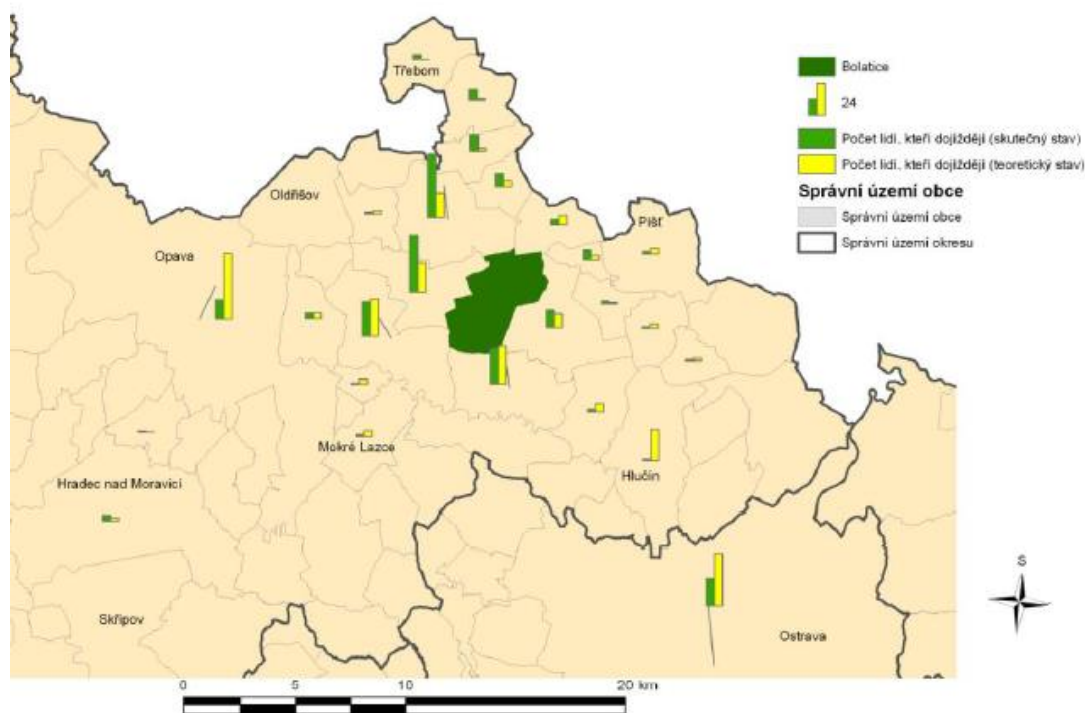
Uvedený vztah je možno linearizovat logaritmováním a následně použít pro řešení např. metodu nejmenších čtverců.

Pokud uplatníme výše uvedené okrajové podmínky, hovoříme o dvojnásobně omezeném modelu (velikost zdrojů zůstává konstantní a rovněž atraktivnost cílů zůstává konstantní).



Obr. 5-8 Výpočet teoretických toků z gravitačního modelu pro dojíždějící do Lanexu Bolatice (Vojta 2009)

## POROVNÁNÍ SKUTEČNÉ A TEORETICKÉ DOJÍŽDKY - LANEX BOLATICE



Obr. 5-9 Porovnání skuteční a teoretické dojížděky z gravitačního modelu pro Lanex Bolatice (Vojta 2009)

### 5.6.2 Jiné gravitační modely a aplikace

V některých případech využíváme jednostranně omezené modely.

První variantou je **model omezený ve zdrojích**. Předpokládá se, že velikost zdrojů je fixní a nemění se. Oproti tomu uvolnění této podmínky pro cíle dovoluje náhradu parametrů atraktivnosti  $\beta_j$  za funkci. Výsledkem je model, který dovoluje sledovat chování toků na změny atraktivnosti cílů.

Z hlediska modelování trhu práce by bylo možné takový model aplikovat pro sledování velikosti dojížděky do jednotlivých center v závislosti na atraktivnosti zaměstnání, profesní skladbě, výši mezd i atraktivnosti centra jako takového (kumulování důvodů dojíždění).

Druhou variantou je **model omezený v cílech**. Atraktivnost cílů se nemění a velikost zdrojů je nahrazena funkcí. Výsledný model dovoluje sledovat např. předpovídat rozmístění zaměstnanců v jednotlivých sídlech podle jejich vzdálenosti, vybavenosti apod. za předpokladu konstantní situace v poptávce po práci.

Další možnosti nabízí jsou:

- model bez okrajových omezení, který dovoluje modelovat např. migraci obyvatel,
- modely používající různé vzdálenostní parametry pro různé cíle - např. nemocnice nemá konstantní spádovou oblast, ale ta se liší podle specialistů v nemocnici - některé služby mají velký dosah (neurologie, radioizotopie), jiné velmi krátký (geriatrie).
- modely s různým vlivem vzdálenosti v různých místech (prostorové rozšíření modelů) - rozdílný vliv vzdálenosti v různých částech sledovaného území
- modely zohledňující soutěžení cílů a jiné formy konkurenčních výhod - řešení pomocí relativní přístupnosti destinace

Zajímavou možností nabízí využívání generalizované ceny dopravy (zohledňuje čas, vzdálenost i další poplatky spojené s dopravou).

## 5.7 Popisná statistika

Kruhový liniový vzorek (modelovaný jako cyklický graf) zahrnuje smyčky (zpravidla cykly, tedy orientované uzavřené cesty) – vhodným příkladem jsou dopravní sítě. Popis tohoto lineárního vzorku se zaměřuje na topologické spojení pomocí matice dosažitelnosti. Tato matice zaznamenává přítomnost nebo nepřítomnost spojení mezi všemi páry bodů v síti. Umožňuje porovnat různé sítě (např. srovnání počtu hran v každém uzlu sítě - nodalita - a suma těchto spojení vůči maximálnímu možnému počtu) nebo posoudit, zda je pohyb v síti mezi různými uzly volný nebo omezený.

### 5.7.1 Popis liniového vzorku

Pro popis liniového vzorku se používá řada ukazatelů:

- **Gama index** hodnotí poměr zjištěného počtu hran (linií) k maximálnímu možnému počtu

$$\gamma_{index} = \frac{l}{l_{max}} = \frac{l}{3(n-2)}$$

kde  $l$  je pozorovaný počet linií (hran)

$l_{max}$  je maximální možný počet linií.

U plošného grafu lze dokázat, že je maximální počet vždy  $3*(n-2)$ , kde  $n$  je počet uzlů; u neplošného (např. letecká spojení) je maximální počet roven  $n*(n-1)/2$ .

- **Alfa index** hodnotí poměr zjištěného počtu smyček k maximálnímu možnému počtu

$$\alpha_{index} = \frac{c}{c_{max}} = \frac{c}{2n-5}$$

kde  $c$  je počet kružnic

$$c_{max} = 2*n-5$$

Dostupnost pro celou síť může být oceněna pomocí **součtu nebo průměru nodality uzlů** (viz topologická dostupnost (kapitola 5.4.3)).

Další ukazatelé využívají charakteristiky hran:

- minimální nebo maximální počet hran pro daný počet uzlů, který je potřebný pro vznik souvislého grafu;
- poměr počtu hran ku počtu uzlů.

Je potřebné poukázat na skutečnost, že plocha uvnitř smyčky grafu není zpravidla dle teorie grafů předmětem zájmu (vyjma výpočtu Eulerovy rovnice). Předmětem zájmu jsou uzly a hrany, nikoliv plocha mezi nimi. U grafu jsou totiž atributy připojovány k hranám, aby vznikl hodnocený graf - např. počet letadel a čas potřebný k překonání určité hrany. Uzly mohou mít připojeny např. údaj typu počet celkem přepravovaných cestujících (tedy vytížení letiště).

## 5.7.2 Náhodnost distribuce pro liniový vzorek

Vedle momentového popisu liniového vzorku je možné se zabývat otázkou náhodnosti distribuce linií - tedy zda je distribuce linií náhodná, tj. zda lze pozorované uspořádání vysvětlit náhodou? Přitom se sledují rozdíly mezi pozorovanou délkou a směrem, hustotou linií, jejich zakřiveností a četností jejich délek.

Naneštěstí binomická, Poissonova i normální distribuce nejsou vhodné pro posuzování distribuce linií ze dvou důvodů. Pro ocenění délek cest, které mohou být oceněny libovolnou hodnotou a ne omezeny na celočíselné výsledky jako v případě četnosti bodů, je vhodná kontinuální hustotní funkce pravděpodobnosti.

Na distribuci pravděpodobnosti má vliv rovněž velikost studované oblasti a její tvar. Např. cesta křížící obdélníkový tvar má pravděpodobně přísně bimodální distribuci. Cesta přecházející čtverec je pravděpodobně jednodimální.

Většina připravených analytických procedur počítá jenom přímé cesty přes oblast. Relativně malá pozornost je věnována zakřiveným liniím, jejichž délka není omezena maximální cestou přes oblast, nebo liniím, které začínají nebo končí uvnitř hranic oblasti.

Statistické analýzy stromového vzorku (příklad kanalizační sítě) se koncentruje na uspořádání (stupeň) takových sítí. Základem procedur je testování, zda pozorovaný liniový vzorek – např. 4 toky prvního stupně, 2 druhého stupně a 1 tok prvního stupně – mohl být vytvořen náhodným procesem.

Obecně se vypočítává počet alternativ, kterými může být náhodným spojovacím postupem dosaženo profilu kanálů pozorované sítě (např. kombinace 4,2,1). Počet kombinací se dělí celkovým počtem možných sítí dávajících tyto proporce. Výsledek je interpretován jako pravděpodobnost dosažení pozorovaného uspořádání kanálů náhodou.

Statistická analýza kruhového vzorku (cyklického grafu) se provádí podobně. Východiskem pro kruhový vzorek je počet uzlů a hran v pozorované síti, který má jasně kritický vliv na to, kolik topologicky rozdílných kruhů může být definováno. Statistické testy zaměřují svou pozornost na nodalitu nebo vnitřní spojitost v síti; jinými slovy jak často je každý uzel zapojen do sítě.

Pro daný počet uzlů a hran může být nodalita vyjádřena numericky rozsahem od 0 (uzel je nespojen) do  $n-1$  ( $n$  je počet uzlů), které znamená, že je každý uzel (nod) přímo spojen se všemi ostatními. Pravděpodobnost pozorované nodality je pak oceněna proti náhodnému spojení uzlů a hran.

Prvním problémem je určit náhodnou distribuci pravděpodobnosti. Na otázku, zda je uzel přímo spojen s jiným, jsou jen 2 odpovědi (nominální data). Protože je celkový počet linií v síti také konstantní, pravděpodobnost spojení uzlu s jiným uzlem se změní, když jsou již linie (hrany) využity (jde tedy o náhodný výběr bez opakování). Vhodným modelem může být hypergeometrická distribuce pravděpodobnosti. Pravděpodobnost jevu je vypočítána jako:

$$P(m, n, q) = \frac{\binom{q}{m} \binom{N-q}{n-1-m}}{\binom{N}{n-1}}$$

$n, q$  .. pozorované počty uzlů (nodů) a hran (linií)

$m$  .... žádaná hodnota nodality

$N$ .... celkový možný počet hran (linií)

Pozorovaná nodalita  $n$  je určena ze sumace 0 a 1 v řádku matice sousednosti (kde počítáme jen s 1 hranou mezi uzly, tj. žádné násobné hrany) při vynechání jedniček na diagonále. Hodnoty pozorované a očekávané nodality mohou být vyjádřeny v absolutních i kumulativních distribucích četnosti a

maximální rozdíl mezi kumulativními součty poskytuje D test, který může být vyzkoušen pomocí Kolmogorov-Smirnovova testu významnosti.

## 5.8 Vizualizace interakčních dat

Vizualizace interakčních dat a tedy nejjednodušší interakční metody jsou spojeny s **liniovým vzorkem**, jeho reprezentací a popisem. Uplatňuje se především ve formě pohybové linie.

K vizualizaci toku podél linie lze použít liniově lokalizovaných **kartodiagramů**, které jsou označovány i jako liniové, stuhové, pásové či proužkové kartodiagramy (Voženílek 2001). Dopravní trasy přitom bývají vyjadřovány schématicky, protože cílem je především znázornění toku ze zdroje do cíle a konkrétní průběh toku geografickým prostorem je až druhořadý. Pomocí těchto kartodiagramů lze zobrazit směr, přepravovanou kvantitu, její strukturu (složené toky), případně další charakteristiky.

Podle toho se vymezuje jednoduchý, složený, součtový, strukturní, srovnávací a izochronní liniový kartodiagram, případně i dynamický typ. V zásadě lze pro vizualizaci interakčních dat velmi dobře použít všechny uvedené typy.

Kaňok (1999) rozlišuje vektorové a stuhové kartodiagramy. U vektorového kartodiagramu je vektor popisován svým počátečním bodem, směrem a délkou. Vymezuje se kartodiagram vektorový dosahový a vektorový proudový, kde pro zobrazení interakčních dat lze použít především první z nich, neboť u druhého jde sice také o interakční data, která však nemají jednoznačně lokalizovány zdroje a cíle toků.

Charakteristiku interakčních dat lze vztáhnout i k jinému nositeli než je linie. Především lze popisovat vlastnosti cílů z hlediska interakčních vazeb, často ve formě vyjádření dostupnosti těchto cílů. Uplatňují se pak metody kartodiagramů (především bodově lokalizovaných), méně často plošně lokalizovaných kartodiagramů, metody izolinií (např. ekvidistanty či izochrony), metody kartogramů, metody anamorfózy.

V práci Hůrského (1969) jsou podrobně popsány kartografické metody znázornění dojížděky do zaměstnání. Zabývá se přitom zobrazením dat ze sčítání lidu, domu a bytů, resp. pro malé areály rovněž daty ze speciálních šetření. Používá řadu hledisek pro klasifikaci připravených map - např. měřítko, věcného obsahu, praktického účelu a forem kartografického zobrazení. Podle věcného obsahu vymezuje:

- mapy intenzity dojížděky (s důrazem na přesnost topografickou, tématického obsahu nebo strukturní diferenciaci tématického obsahu),
- mapy dojížděkových proudů (členěny dále podle důrazu podobného jako u první skupiny),
- mapy regionalizace dojížděky (s důrazem na topografickou přesnost a přesnost intenzity jevu, na topografickou přesnost a relativní poměr mezi centry dojížděky, na geografickou názornost a přesnost intenzity jevu, na geografickou názornost a relativní poměr mezi centry dojížděky).

S jistým zjednodušením lze říci, že mapy intenzity dojížděky využívají především kartografické metody znakové (používání symbolů) resp. bodově lokalizovaných kartodiagramů, mapy dojížděkových proudů využívají především kartografické metody liniově lokalizovaných kartodiagramů a mapy regionalizace dojížděky především metody kartogramů nebo plošně lokalizovaných kartodiagramů.



## 6 Polygony

Objekty, u kterých nemůžeme s ohledem na odpovídající reprezentaci (z důvodu vizualizace a hlavně analýzy) zanedbat jeho minimálně 2 rozměry, používáme ve 2D zobrazení k reprezentaci polygony. V některých případech objekty netvoří souvislou plochu a pak pro reprezentaci objektu použijeme skupinu polygonů se stejnou hodnotou (se stejným identifikátorem). Tato skupina polygonů se označuje jako areál.

Popisná charakteristika (kapitola 6.1) v případě polygonů nebývá problémem. Zjišťování náhodnosti distribuce velikosti ploch či jiných charakteristik se nevyužívá, spíše jsou zajímavé techniky sledující náhodnost a autokorelaci hodnot přidělených polygonům (kap. 6.5).

Z celé řady metod, které jsou aplikovány na tento typ reprezentace prostorových dat, vybereme jen některé specifické:

- Problém plošné interpolace (kapitola 6.2)
- Problém regionalizace (kapitola 6.3)
- Vyhlazování areálových dat (kapitola 6.4)
- Sledování autokorelace (kapitola 6.5)
- Multivariační techniky (kapitola 6.6)
- Regresní modelování (kapitola 6.7)

### 6.1 Popisná statistika pro polygony

Charakterizovat soustavu polygonů lze snadno s využitím jednoduchých statistických ukazatelů typu: průměrná plocha polygonu, průměrná plocha areálu (multipolygonu), průměrný obvod polygonu, variabilita plochy polygonu (rozptyl, směrodatná odchylka), statistika zastoupení polygonů v areálu (průměrný počet polygonů v 1 areálu, max. a min. počet, variabilita atd.).

Lze využít i tvarové ukazatele typu koeficient zakulacenosti apod.

K důležitým charakteristikám, které hodnotí zastoupení polygonů, tříd, hranic ve sledovaném území, patří ukazatele krajinné metriky.

Rovněž lze hodnotit míry konektivity mezi polygonu (např. pomocí definice sousedství viz kap. 6.5 a vyjadřování příslušných statistických ukazatelů).

### 6.2 Problém plošné interpolace

Častým problémem je situace, kdy máme analyzovat 2 faktory, které jsou vázány každý k jiným polygonům. Např. cenová mapa a cenovní okrsky rozdělují město každý jiným způsobem (do jiných polygonů). Máme-li porovnat např. sociálně-ekonomickou charakteristiku obyvatelstva s cenou obydlí (pozemků), musíme sjednotit prostorovou základnu - tedy převést údaje jednoho faktoru na polygony vymezené pro druhý faktor.

Nestabilita územního vymezení administrativních jednotek vede často i k potřebě modifikovat tvar areálů těchto jednotek, resp. distribuovat hodnoty z jednoho areálu do druhého. Tuto úlohu označujeme jako problém areálové interpolace neboli problém MAUP (*Modifiable Area Unit Problem*). Problém se již dotýká zpracování dat a ne jejich vymezení.

Cílem je distribuovat hodnoty ze zdrojových areálů do cílového areálu při jejich překrytí. Podstatný je způsob kombinace dílčích hodnot (přenesených ze zdrojových areálů) do výsledné hodnoty pro cílový areál.

Běžně se používají vážené aritmetické průměry a váha se určuje v závislosti na tom, zda jde o veličinu absolutní nebo relativní. Snadné řešení představuje **vážení pomocí ploch areálů** (*areal weighting*).

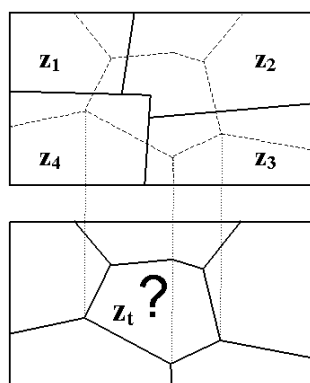
## 6.2.1 Absolutní hodnoty

Je-li veličina  $z$  absolutního charakteru (např. počet, součet), pak se  $z_t$  vypočte jako suma hodnot  $z_{st}$ . Předpokládáme, že hodnoty veličin  $Z$  jsou rovnoměrně rozloženy po celé ploše areálu, a že tedy hodnota části plochy je přímo úměrná velikosti této části ploch

$$z_t' = \sum_s z_{st} = \sum_s \frac{A_{st}}{A_s} * z_s = \sum_s \frac{A_{st} z_s}{A_s}$$

kde

- $z_t'$  odhadovaná hodnota  $z$  v cílovém areálu  $t$
- $z_s$  hodnota zdrojového areálu  $z$
- $z_{st}$  hodnota v místě průniku cílového areálu  $t$  se zdrojovým areálem  $s$
- $A_s$  plocha zdrojového areálu  $z$
- $A_{st}$  plocha průniku cílového areálu  $t$  se zdrojovým areálem  $s$



Obr. 6-1 Schéma překryvné operace

## 6.2.2 Relativní hodnoty

Je-li veličina relativní (např. indexy, poměry, %), pak se  $z_t$  vypočte jako vážený průměr ze  $z_{st}$ . Dále předpokládáme rovnoměrné rozdělení hodnot v areálu a tedy hodnota části areálu se rovná hodnotě areálu (tedy  $z_{st} = z_s$ ), a vahou je v nejjednodušším případě podíl plochy průniku na celkové výsledné ploše.

$$z_t' = \frac{\sum_s z_{st} A_{st}}{\sum_s A_{st}} = \frac{\sum_s z_s A_{st}}{A_t} = \sum_s \frac{A_{st} z_s}{A_t} = \frac{1}{A_t} \sum_s A_{st} z_s$$

kde

- $z_t'$  odhadovaná hodnota  $z$  v cílovém areálu  $t$
- $z_s$  hodnota zdrojového areálu  $z$
- $z_{st}$  hodnota v místě průniku cílového areálu  $t$  se zdrojovým areálem  $s$
- $A_t$  plocha cílového areálu  $t$
- $A_{st}$  plocha průniku cílového areálu  $t$  se zdrojovým areálem  $s$

Sten (1997) popisuje v podstatě dasymetrickou úpravu této metody, využívající pro vážení i rozmístění osídlení (vrstva osídlení).

Složitější postupy využívají pro modelování distribuce sledované veličiny v ploše zpravidla Poissonovo rozdělení. Regresní modely založené na předpokladu Poissonovy distribuce pro  $\mathbf{z}$  se aplikují buď přímo nebo pomocí iteračního postupu označovaného jako EM algoritmus (Flowerdew, Green 1994), používaného pro absolutní hodnoty. EM algoritmus opět předpokládá rozdělení cílového areálu s plochou  $\mathbf{A}_t$  do sady průnikových polygonů s plochami  $\mathbf{A}_{st}$ , pro které se vypočítá neznámá hodnota  $\mathbf{z}$  a následně se součtem určí hledaná hodnota  $\mathbf{z}'_t$ . Iterační povaha EM algoritmu je založena na postupném opakování 2 kroků, označených jako E a M, až do situace, kdy výsledky konvergují.

E krok zahrnuje výpočet očekávané hodnoty  $\mathbf{z}_{st}$ :

$$z_{st}' = \frac{\mu_{st}' z_s}{\sum_s \mu_s'}$$

$\mu$  střední hodnota

M krok provádí výpočet  $\mu_{st}'$  resp.  $\mu_s'$  podle modelu  $\mu_{st} = \mu(\beta, \mathbf{z}_t, \mathbf{A}_{st})$  metodou maximální věrohodnosti (optimalizace neznámých parametrů  $\beta$ ).

První krok E využívá přímo původních dat. Další kroky E již využívají hodnot  $\mu_{st}$  a  $\mu_s$  určených podle postupně zpřesňovaného odhadu  $\beta$  parametrů.

Výhodou modelu je i možnost využití dalších proměnných v regresním vztahu ovlivňujících výsledek  $\mu_s$  resp.  $\mu_{st}$ .

### 6.3 Problém regionalizace

Problém regionalizace řeší seskupení základních jednotek do vyšších celků.

Některé programové produkty mají k tomuto účelu specializované funkce - např. distriktová funkce MapInfra.

Příkladem může být zařazení obcí do správních celků jako jsou okresy. Řešení je možné demonstrovat na příkladu, kdy máme zařadit např. C obvodů do D okresů.

Matematicky můžeme požadavek zapsat maticí:

	1	2	3	..	D	$\sum_d x_{cd}$
1	$x_{11}$	$x_{12}$	$x_{13}$	..	$x_{1D}$	=1
2	$x_{21}$	$x_{22}$	$x_{23}$	..	$x_{2D}$	=1
3	$x_{31}$	$x_{32}$	$x_{33}$	..	$x_{3D}$	=1
..	..	..	..	..	..	=1
C	$x_{C1}$	$x_{C2}$	$x_{C3}$	..	$x_{CD}$	=1
$\sum_c x_{cd}$	$\geq 1$	$\geq 1$	$\geq 1$	$\geq 1$	$\geq 1$	

V matici jsou prvky  $x_{cd}$ , které nabývají hodnoty 1, je-li obvod c v okresu d, nebo hodnoty 0, není-li obvod c v okresu d.

Další podmínky jsou vyjádřeny pomocí sum:

$$\sum_c x_{cd} \geq 1 \quad x_{cd} = 0 \text{ nebo } 1 \quad \text{pro všechna } c \text{ a všechna } d$$

$$\sum_d x_{cd} = 1 \quad \text{pro všechna } d \text{ platí, že každý okres má 1 a více obvodů}$$

$$\sum_c x_{cd} = 1 \quad \text{pro všechna } c \text{ platí, že každý obvod je právě v 1 okresu}$$

Přiřazování obvodů do okresů může probíhat např. na základě požadavku stejného počtu obyvatel v každém okrese:

$$\text{minimalizuj výraz } \sum_d \left( \sum_c p_c x_{cd} - \frac{P}{D} \right)^2$$

$p_c$  ... populace v obvodě  $c$

$P$  .... celková populace ve státu

$D$ .... počet okresů

Vedle minimalizace odchylek od průměru (tedy metody nejmenších čtverců) je možné použít řešení kombinační (pomocí stromu) nebo strategii vyměňování (obvody jsou zařazeny postupně do všech okresů, až je výraz minimální). Tento algoritmus však nemusí vést vždy k nejlepšímu výsledku.

Alternativním testovaným požadavkem bývá vysoká geografická a ekonomická kontinuita.

Konstantní podmínkou je kontinuita plochy okresu, tj. obvody řazené do 1 okresu spolu musí sousedit.

Sousedství můžeme vyjádřit např. pomocí matice sousednosti (někdy také spojitosti). Sousedství je možné definovat různými způsoby (viz např. varianty sousedství (kapitola 6.4.1).

Tab.6-1 Příklad matice sousednosti pro polygony

	obvod A	obvod B	obvod C
obvod A	1	1	0
obvod B	1	1	1
obvod C	0	1	1

Z příkladu vyplývá, že pouze obvod A a C nesousedí.

## 6.4 Vyhlazování areálových dat

Cílem vyhlazování areálových dat je odhalení a zvýraznění efektu 1.řádu, tedy změn ve střední hodnotě (označované zpravidla jako trend), které se ve sledovaném území projevují. Pozorovaná sada dat může vykazovat v území značné rozdíly a někdy lze jen obtížně nalézt generelní průběh hodnot (trend), odhalit a interpretovat texturu (vzor), která se projevuje.

Předpokládejme, že pozorované atributy  $Z(s_i)$  představují realizaci proměnné veličiny  $Z$  v jednotlivých areálech  $s_i$ , složené v nejjednodušším případě z pravidelné složky, která závisí na poloze v území, a náhodné složky, která může ale nemusí záviset na poloze v prostoru. Při vyhlazování dat se tedy potlačuje náhodná složka a posiluje zobrazení pravidelné složky.

Vyhlazení má smysl v případě, že existuje jistá prostorová kontinuita hodnot. Pokud předpokládáme, že náhodná složka nezávisí výrazně na poloze v prostoru, můžeme usuzovat, že náhodné variace se v blízkém okolí mohou vzájemně eliminovat a že tedy můžeme náhodný vliv na hodnoty  $Z$  potlačit nějakou formou váženého průměrování. Zpravidla se přitom využívá jednoduchého lineárního modelu.

K vyhlazování se používají:

- klouzavé průměry (kapitola 6.4.1),
- Bayesovské vyhlazení (kapitola 6.4.2),
- jádrové vyhlazení (modifikovaná varianta jádrového odhadu pro body (kapitola 3.4.2)
- mediánové vyhlazení.

## 6.4.1 Prostorové klouzavé průměry

Metoda klouzavých průměrů představuje jednoduchý způsob vyhlazení 1D (jednorozměrných) dat, tedy např. hodnot v časové řadě nebo údajů zjištěných podél linie. Původní, tj. měřená či pozorovaná hodnota se přitom nahrazuje jistým váženým průměrem ze sousedních hodnot, kde velikost sousedství a váhy přidělované různě vzdáleným hodnotám odpovídají určitému modelu.

V základní variantě postup odpovídá proložení polynomu stupně  $p$ .

$$Z_t = a_0 + a_1 t + a_2 t^2 + \dots + a_p t^p$$

K základním vlastnostem patří, že součet vah je roven 1 a že váhy jsou symetrické kolem prostřední hodnoty.

Doporučuje se pro řezy spojitými povrchy používat lineární, kvadratický a kubický polynom, protože vyšší řády polynomu jsou často umělé, jejich koeficienty nestálé a příliš závislé na vstupních datech. Proto doporučuje vizualizovat raději sekvenci několika polynomů nižšího řádu uplatňovaných v překrývajících se úsecích.

Koeficienty polynomu lze určit např. metodou nejmenších čtverců. Výpočet může být poněkud složitý, proto se často používají již „předdefinované“ vážené klouzavé průměry, které lze použít pro delší řady, protože využívají např. 7, 15 nebo 21 měření v řadě. Uvedme jako příklad postup výpočtu pro Spencerův vyhlazený průměr s 15 členy:

$$Z_x' = [74 * Z_x + 67 * (Z_{x+1} + Z_{x-1}) + 46 * (Z_{x+2} + Z_{x-2}) + 21 * (Z_{x+3} + Z_{x-3}) + 3 * (Z_{x+4} + Z_{x-4}) - 5 * (Z_{x+5} + Z_{x-5}) - 6 * (Z_{x+6} + Z_{x-6}) - 3 * (Z_{x+7} + Z_{x-7})] / 320$$

kde

$Z_x'$  vyhlazená hodnota v místě  $x$   
 $Z_{x+1}, Z_{x-1}, \text{atd.}$  původní hodnota v místě  $x$  a v sousedních místech v řadě

Je zřejmé, že jedním z parametrů, které je nutno zvolit je i délka intervalu výpočtu. Tento nedostatek odstraňuje metoda exponenciálního vyhlazení (Rogalewitz 1993).

**Prostorové klouzavé průměry** pracující ve 2D nahrazují původní hodnoty pro každý areál váženým aritmetickým průměrem hodnot v sousedních areálech.

Na rozdíl od jednorozměrného případu se používá především přímých sousedů a jen u některých variant také vzdálenějších sousedů, tj. areálů bez společné hranice. Přidělované váhy jsou také jednodušší a nesimulují gaussovský úbytek pravděpodobnosti s ohledem na velikost areálů a z toho vyplývající nepřesnost v lokalizaci ujištěných hodnot. Základní vztah pro výpočet nové hodnoty areálu:

$$\mu_i' = \frac{\sum_{j=1}^n w_{ij} * z_j}{\sum_{j=1}^n w_{ij}}$$

kde

$z_j$  původní hodnoty v sousedních areálech  
 $w_{ij}$  váha hodnoty v sousedním areálu  $j$  z místa  $i$   
 $j$  index vymezující sousední areály  
 $\mu_i'$  vyhlazená hodnota  $z$  v areálu  $i$

Alternativou může být pro výpočet poměru i sumace čitatele a jmenovatele a následně výpočet podle vztahu:

$$r_i = \frac{y_i}{n_i} \quad \Rightarrow \quad r_i' = \frac{\sum_{j=1}^n w_{ij} * y_j}{\sum_{j=1}^n w_{ij} * n_j}$$

Protože váhy  $w_{ij}$  nejsou standardizovány do rozsahu  $\langle 0,1 \rangle$ , musí se provést standardizace jejich součtem ve jmenovateli zlomku.

Existuje celá řada variant této metody lišící se způsobem vymezení sousedství a nastavením vah  $w_{ij}$ . Sousedství a hodnoty vah doporučuje popisovat pomocí matice vah  $\mathbf{W}$  o rozměru  $\mathbf{n} \times \mathbf{n}$ , kde  $\mathbf{n}$  je počet areálů. Váha  $w_{ij}$  mezi dvěma areály může být vyjádřena jako:

- 1)  $w_{ij}=1$  Pokud těžiště areálu  $j$  je jedním z  $k$  nejbližších těžišť vůči areálu  $i$   
 $w_{ij}=0$  v ostatních případech
- 2)  $w_{ij}=1$  Pokud těžiště areálu  $j$  je do jisté vzdálenosti  $\delta$  od areálu  $i$   
 $w_{ij}=0$  v ostatních případech
- 3)  $w_{ij}=d_{ij}^\gamma$  Pokud je vzdálenost  $d_{ij}$  mezi těžištěm areálu  $i$  a  $j$  menší než jistá vzdálenost  $\delta$  od areálu  $i$  ( $\gamma < 0$  vyjadřuje strmost vlivu vzdálenosti)  
 $w_{ij}=0$  v ostatních případech
- 4)  $w_{ij}=1$  Pokud areál  $j$  sdílí společnou hranici s areálem  $i$   
 $w_{ij}=0$  v ostatních případech
- 5)  $w_{ij}=l_{ij}/l_i$  kde  $l_{ij}$  je délka společné hranice mezi areálem  $i$  a  $j$ ; resp.  $l_i$  je obvod areálu  $i$

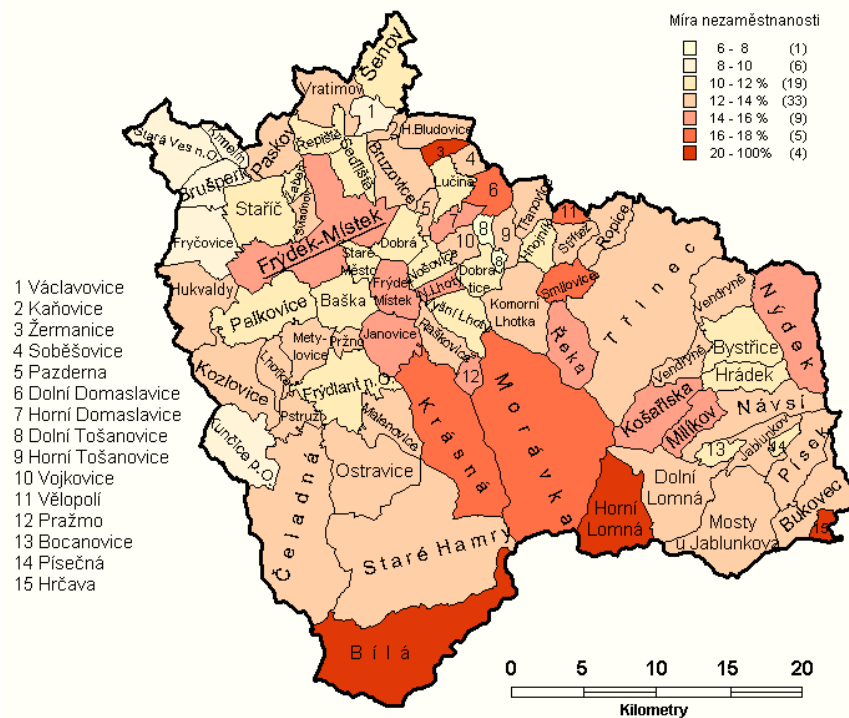
Některé způsoby výpočtu  $w_{ij}$  obsahují parametry  $(k, \delta, \gamma)$ , které je možné optimalizovat.

U rastru se používá sousedství typu věž (rook's case), královna (queen's case) (někdy také král – king) a střelec (bishop's case). U polygonů některé programy (GeoDa) uvažují o sousedství typu věž (společná liniová hranice) a královna (společné hranice liniové i bodové – navíc zahrnuje i polygony, které se dotýkají pouze v 1 uzlu).

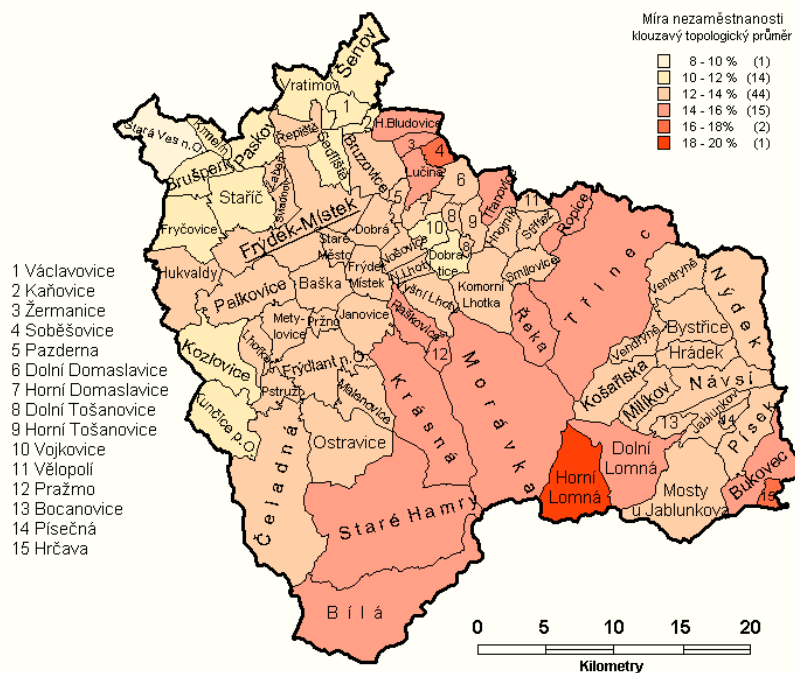
Metoda byla vyzkoušena pro míru nezaměstnanosti v obcích okresu Frýdek-Místek k 31.12.2001, která nejeví příliš zřetelný trend v primární podobě (obr. 6-2). Pro výpočet vah bylo použito jednoduché 4. varianty, tedy topologického způsobu, kdy je vyhlazení prováděno na základě hodnot v sousedních areálech (obr. 6-3). S ohledem na zkoumání míry nezaměstnanosti, tedy poměru, byla dále použita modifikace se sumací čitatele a jmenovatele v sousedních obcích a až následný výpočet míry nezaměstnanosti (obr. 6-4).

Trend u topologického vyhlazení je v obou případech již zřetelný (očekávaná vlastnost topologického vyhlazení), nicméně byly získány dva odlišné vzorky, což vyžaduje vyzkoušení vyhlazení i dalšími metodami, aby se našla vhodná varianta řešení. Je třeba říci, že pozorovaná situace není zcela typická, protože okresní město (s velkým počtem obyvatel) má vysokou míru nezaměstnanosti a skládá se ze 2 nespojitých polygonů, což vede k abnormálnímu zvýšení počtu sousedů a tedy dalšímu zdůraznění vlivu tohoto areálu v území (viz především obr. 6-5).

Uvedených 5 variant neposkytuje vyčerpávající přehled možností vyjádření vah (vazeb) mezi areály, lze využít i různých kombinací výpočtu např. délka společné hranice a vzdálenost mezi těžišti. Vhodnou formou zobecnění vzdálenosti je využití cestovního času mezi areály. Pro některé metody se používají matice vah pro popis vazeb vyšších řádů, nejenom pro nejbližší sousedství. Tímto způsobem lze definovat tzv. prostorové kroky (*spatial lags*), tedy popisovat váhy pro 2. nejbližší sousedství, 3. nejbližší sousedství atd. Tento postup můžeme použít i pro popis kontinuity zkoumaného fenoménu v oblasti. Zde se však již nabízí alternativa ve formě transformace areálových měření na bodová, ať už v nepravidelné síti (využití těžišť areálů nebo center osídlení) nebo v pravidelné síti, a následném uplatnění geostatistických metod. Jinou možností je po provedené transformaci využití jádrového vyhlazení.



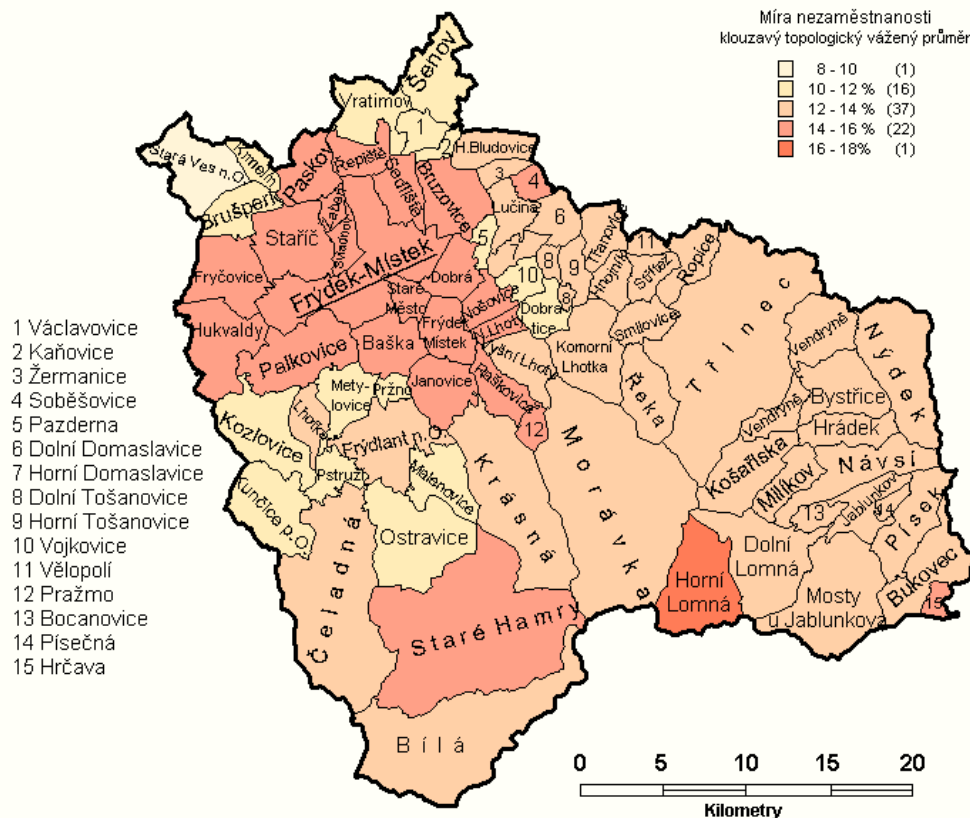
Obr. 6-2 Míra nezaměstnanosti v obcích okresu Frýdek-Místek k 31.12.2001



Obr. 6-3 Topologicky vyhlazená míra nezaměstnanosti v obcích okresu Frýdek-Místek k 31.12.2001

Nevýhodou této metody je skutečnost, že jednotlivé areály mají všechny stejnou váhu, přitom mají rozdílné plochy (což ovšem některé varianty výpočtu  $w_{ij}$  nepřímo zohledňují), ale v případě socioekonomických dat především rozdílnou velikost datové základny.

Další komplikací je výskyt nespojitých areálů (např. některé obce mají nespojitě území). V tom případě výčet sousedních areálů značně narůstá a tyto areály příliš silně ovlivňují své okolí.



Obr. 6-4 Vážený aritmetický průměr míry nezaměstnanosti ze sousedních obcí okresu Frýdek-Místek k 31.12.2001

## 6.4.2 Bayesovo vyhlazení

Bayesovo vyhlazení využívá Bayesova přístupu k výpočtu pravděpodobnosti, kde se apriorní nepodmíněná pravděpodobnost pro sledovaný jev kombinuje se zjištěnou (měřenou) pravděpodobností a vzniká nová, posteriorní pravděpodobnost pro sledovaný jev.

Bayesův přístup je možno při vyhlazování využít tak, že získaná data kombinujeme s vhodným apriorním odhadem, kterým může být např. střední hodnota sledovaného jevu pro celé území. Areály, kde je nízká věrohodnost získaných dat, protože je zde malá datová základna, jsou více regulovány (tj. vyhlazovány) než areály s velkou datovou základnou a tedy vyšší věrohodností získaných dat.

Bayesův odhad hodnot  $z_i$  vychází z obecného vztahu (Bailey, Gatrell 1995):

$$z_i = w_i * r_i + (1 - w_i) * \gamma_i$$

$$w_i = \frac{\phi_i}{\phi_i + \frac{\gamma_i}{n_i}}$$

kde

$w_i$  váhový faktor každého areálu  $i$  (pokud má vzhledem k datové základně areál velkou váhu, bude výsledný odhad odpovídat  $r_i$ , tj. původním poměru v areálu).

$r_i = y_i/n_i$  mapovaný poměr (relativní či intenzivní hodnota)

$y_i$  počet sledovaných jevů v populaci (extenzivní či absolutní hodnota)

$n_i$  velikost příslušné populace

$\gamma_i$  střední hodnota apriorní distribuce  $r_i$  v místě  $i$

$\phi_i$  rozptyl apriorní distribuce  $r_i$  v místě  $i$



Předpokládáme tedy existenci apriorní distribuce sledovaného poměru  $r_i$  pro daný areál, popsané 2 základními charakteristikami  $\gamma_i$ ,  $\phi_i$ . Váhový faktor  $w_i$  je tedy závislý na velikosti dotčené populace (datová základna) a na rozptylu  $\phi_i$ . Hodnota rozptylu je nepřímo úměrná věrohodnosti apriorní distribuce.

Apriorní distribuci a tedy ani její charakteristiky v daném místě  $i$  neznáme a proto ho musíme nahradit vhodným odhadem. Tím může být průměr a rozptyl pro celé sledované území, předpokládáme-li, že charakteristiky jsou v území konstantní. Průměr  $\gamma_i$  a rozptyl  $\phi_i$  můžeme odvozovat metodou maximální věrohodnosti, ta však vyžaduje iterační zpracování. Používají se určitá zjednodušení - buď se využívá předpokladu gama distribuce a odvozují se jeho 2 parametry, což opět vede k numerické aproximaci, nebo se přímo odhadují  $\gamma$  a  $\phi$  na základě momentových charakteristik. V tom případě je

$$\gamma = \bar{r}_i = \frac{\sum y_i}{\sum n_i}$$

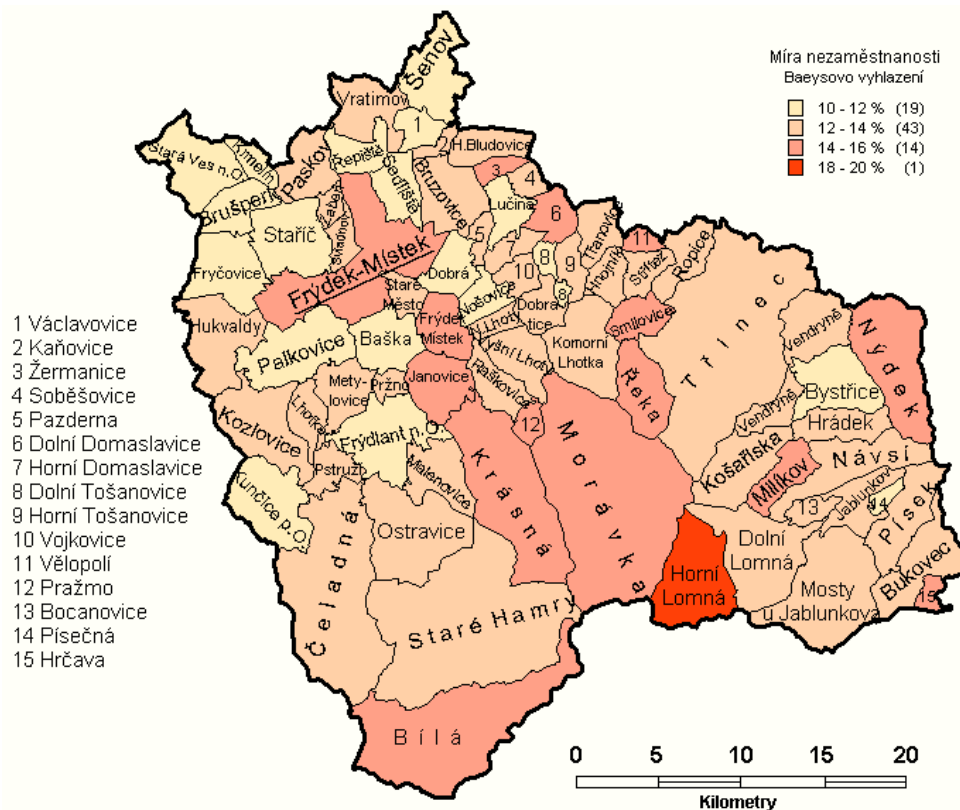
$$\phi = \frac{\sum n_i * (r_i - \gamma)^2}{\sum n_i} - \frac{\gamma}{n}$$

kde  $n$  s pruhem je průměrná populace v území.

Vzorec platí, pokud je  $\phi > 0$ , jinak  $\phi = 0$ .

Původní vzorec lze tedy upravit do podoby:

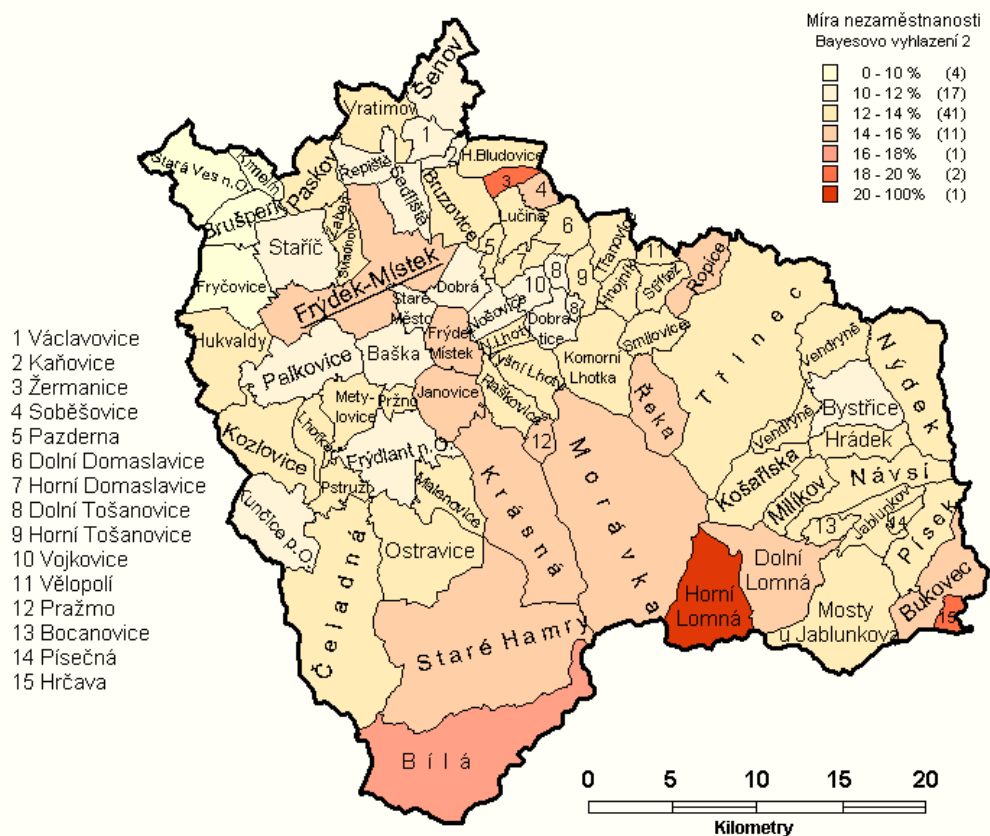
$$z_i = \frac{\phi * r_i + \frac{\gamma^2}{n_i}}{\phi + \frac{\gamma}{n_i}} \quad \text{což lze upravit na} \quad z_i = \gamma + \frac{\phi * (r_i - \gamma)}{\phi + \frac{\gamma}{n_i}}$$



Obr. 6-5 Bayesovo vyhlazení míry nezaměstnanosti v obcích okresu Frýdek-Místek k 31.12.2001

V uvedené základní podobě používá Bayesovo vyhlazení konstantní průměr a rozptyl v celé území. Je zřejmé, že tento předpoklad nemusí vždy vyhovovat a že v řadě případů bude vhodné předpokládat existenci trendu a změn variability v území. Metodu je možné snadno modifikovat nahrazením globálního průměru a rozptylu za statistické charakteristiky sousedství - použije se hodnot sousedních areálů k výpočtu průměru a rozptylu. Jde tedy o adaptivní Bayesovo vyhlazení, které bude vyhlazovat s přihlédnutím k lokálními vývoji hodnot. Pojem sousedství je zde opět použit v obecné rovině - k vymezení sousedství můžeme použít topologie a/nebo vzdálenosti (zpravidla těžišť areálů). Můžeme tak dosáhnout zajímavých výsledků jako je např. vyhlazená míra nezaměstnanosti s ohledem na spádovost území vyjádřenou např. dopravními časy mezi obcemi.

Výsledek adaptivního Bayesova vyhlazení pro studovanou situaci s využitím sousedství pro výpočet lokálních průměrů je na obr. 6-6. Je evidentní, že je míra nezaměstnanosti příliš snížena a že je využití nejbližších sousedů (zvláště pro velmi nehomogenní území) značně omezené. Ukazuje se, že je potřebné upravit metody tak, aby se ovlivňující okolí stanovovalo na základě studia autokorelace jevu v území.



Obr. 6-6 Adaptivní Bayesovo vyhlazení míry nezaměstnanosti v obcích okresu Frýdek-Místek k 31.12.2001

## 6.5 Sledování autokorelace

Předcházející metody se zaměřovaly na sledování prostorového fenoménu v daném území, především na střední hodnoty. Takové hodnocení lze označit za sledování efektu 1.řádu. Oproti tomu se relativně málo metod věnuje vývoji variability v území (především odchylek od střední hodnoty), tedy sledování efektu 2.řádu.

Obecně se posuzuje podobnost charakteristik sledovaných objektů v závislosti na jejich vzdálenosti.

Ke sledování prostorové variability poskytuje nejvhodnější nástroje geostatistika.

Koncept regionalizované proměnné znamená, že předmětem zpracování je spojitá (kontinuální) proměnná. Při aplikaci geostatistických metod na areálová data se používají jednodušší strukturální funkce, kde lze snáze aplikovat zobecněnou vzdálenost. V kapitole 6.4.1 "Klouzavé průměry" je uvedeno několik možností vyjádření sousedství (nejbližšího, ale také sousedství ve druhém, třetím a dalším krocích) pro prostorové průměrování. Tyto míry  $w_{ij}^{(k)}$  definují obecnou vzdálenost mezi areály a zapisují se do samostatných matic pro každý řád  $k$ .

K výpočetně jednoduchým, patří neparametrické testovací charakteristiky pro sledování prostorové autokorelace.

**Smans-Esteve ukazatele** (Smans, Esteve 1996):

$$\sum_{i \neq j} w_{ij} * c_{ij}$$

za podmínky

$C_{ij}=1$  pokud  $i$  a  $j$  patří do stejné kategorie (autor to popisuje jako použití stejné „barvy“)  
 $C_{ij}=0$  v jiných případech

nebo

$$\sum_{i \neq j} w_{ij} * o_{ij}$$

kde

$o_{ij}$       rozdíl pořadí  $z_i$  a  $z_j$  zjištěných hodnot jevu v místě  $i$  a  $j$

Tyto ukazatelé však nevyužívají 2.mocniny odchylek a nelze je tedy považovat za klasické míry autokorelace či kovariance.

Mnohem častěji se jako strukturální funkce používají **Moranovo I kritérium** a **Gearyho C kritérium**, které lze v zobecněné podobě prezentovat jako:

$$I_k = \frac{n * \sum_{i=1}^n \sum_{j=1}^n w_{ij}^{(k)} * (z_i - \bar{z}) * (z_j - \bar{z})}{\left( \sum_{i=1}^n (z_i - \bar{z})^2 \right) * \left( \sum_{i \neq j} \sum w_{ij}^{(k)} \right)} \quad C_k = \frac{(n-1) * \sum_{i=1}^n \sum_{j=1}^n w_{ij}^{(k)} * (z_i - z_j)^2}{2 * \left( \sum_{i=1}^n (z_i - \bar{z})^2 \right) * \left( \sum_{i \neq j} \sum w_{ij}^{(k)} \right)}$$

kde

$w_{ij}^{(k)}$       indikace vzdálenosti mezi areály  $i$  a  $j$  pro krok  $k$  (viz popis sousedství (kapitola 6.4.1))  
 $z_i$               zkoumaná veličina v místě  $i$  ( $z$  s pruhem představuje aritmetický průměr)

Je zřejmé, že Moranovo I kritérium představuje analogii kovariační funkce, zatímco Gearyho C kritérium odpovídá semivariogramu.

Shluková textura má hodnotu C mezi 0 a 1, hodnota I je větší než očekávaná pro náhodnou distribuci. Rozptýlená textura (dispersed) má hodnotu C mezi 1 a 2, hodnota I je menší než očekávaná pro náhodnou distribuci. Očekávaná hodnota Moranova I kritéria pro náhodnou distribuci je  $(-1)/(n-1)$ ; očekávaná hodnota pro Gearyho C je 1 (nezáleží na velikosti souboru) (Lee, Wong).

Je třeba upozornit, že ani jedno z kritérií neposkytuje hodnoty přesně v uvedeném rozmezí. Způsobují to především reálné problémy s vyjádřením matice vzdáleností a reálných vazeb. Bailey, Gatrell (1995) navrhuje korekční člen  $I'$ , kterým se I hodnota dělí a tím se získá předpokládaných rozsah hodnot (podobně pro C).

Následně se vykreslují I nebo C kritéria proti hodnotě kroku, pro který se provádí interpretace.

Některí autoři varují před používáním statistických charakteristik uvedených indikátorů prostorové autokorelace pro adjustovaná data, založená na věkové standardizaci z populace. Populace obyvatel je totiž zpravidla prostorově autokorelována a to vede k promítnutí této závislosti do adjustovaných dat o sledovaném jevu. Často by tedy jev vykazoval prostorovou autokorelaci, i když ve skutečnosti by byl konstantní v celé oblasti. Jako vhodné řešení se doporučuje testování pravděpodobnosti založené na simulaci.

Namísto korekčního činitele se v poslední době uplatňují spíše techniky randomizace při vzorkování, kdy jsou hodnoty atributů (ty známé hodnoty) přidělovány náhodně různým nositelům (objektům v území). Tedy namísto náhodné volby hodnoty atributu se pro známé hodnoty volí náhodně různá možná místa jejich výskytu. Na jejich základě se pak počítají upravené odhady rozptylu, které lze společně s odhadnutou střední hodnotou použít k testování významnosti zjištěných odchylek od náhodného (očekávaného) stavu (Lee, Wong 2001).

Dykes (1994) uvádí jako alternativu pro sledování prostorové autokorelace u areálových dat zjišťování prostorové závislosti distribuce pravděpodobnosti sledovaných hodnot, která se zapisuje do matice podle směrů a kroků (využity techniky digitálního zpracování obrazu, zřejmě analogie povrchu semivariogramu podle Meer 1992). Namísto stupně šedi v obraze se definují třídy a zaznamenají se výskytu pro každou třídu, které jsou dále standardizovány a pak porovnány s očekávanou hodnotou.

Výsledek může být zobrazen jako soustava 2,5D povrchů ve 3D zobrazení, kde každý povrch odpovídá určitému kroku vzdálenosti v autokorelační analýze.

Moranovo I a Gearyho C nerozliší, zda se projevuje pozitivní autokorelace mezi vysokými nebo naopak nízkými hodnotami. Proto se používá charakteristika  $G(d)$ , jako analogie Ripleyovy  $K$  funkce.

$$G_d = \frac{\sum_{i=1}^n \sum_{j=1}^n w_{ij}(d) * z_i * z_j}{\left( \sum_{i \neq j} z_i * z_j \right)}$$

Váha  $W_{ij}$  nabývá hodnoty 1, pokud jsou areály  $i$  a  $j$  od sebe vzdáleny méně než  $d$  (Lee, Wong).

$G$  hodnota bude větší, pokud jsou blízké hodnoty vysoké a naopak nižší v případě nízkých okolních hodnot.

Moranovo I a Gearyho C měří prostorovou autokorelaci pro celé území, jde tedy o globální stanovení míry autokorelace či variability. V některých případech ale potřebujeme rozlišit lokální situaci v autokorelaci, zjistit, zda území se chová jako homogenní (stále stejná hodnota autokorelace v území) nebo naopak heterogenní (hodnoty autokorelace v jednotlivých místech se významně liší). K tomu slouží např. LISA.

Lokální míra prostorové asociace LISA (Local Indicators of Spatial Association) odpovídá lokální verzi Moranova nebo Gearyho kritéria, případně  $G$  ukazatele. Tyto míry vyjadřují úroveň vazeb s okolím, blízkost (resp. různost pro Gearyho  $k$ ) hodnot u prostorově blízkých jednotek.

Lokální Moranovo I pro jednotku  $i$  je definováno jako (Lee, Wong 2001):

$$I_i = r_i * \sum_j (w_{ij} * r_j)$$

kde  $r_i$  a  $r_j$  jsou standardizované hodnoty veličiny  $z$ .

Podobně jako u globálního Moran I, vysoká hodnota lokálního Moranova I ukazuje na shlukování podobných hodnot (ať již jsou vysoká nebo nízká). Pouhý výpočet tohoto ukazatele nám ale sám o sobě nic neřekne, vysoká hodnota mohla vzniknout náhodou. Proto podobně jako v předchozích případech je potřebné stanovit odhadovanou střední hodnotu a rozptyl a spočítat pravděpodobnost takového výsledku (resp. testovat, zda odchylka proti náhodnému stavu mohla vzniknout náhodně).

Očekávaná střední hodnota (Anselin 1995) je  $E[I_i] = (-w_i)/(n-1)$ . Složitěji se počítá rozptyl. Pro každou územní jednotku vypočteme tyto očekávané hodnoty a můžeme tedy i v daném místě prověřit dosaženou pravděpodobnost (nenáhodnost) výsledku.

Existuje také lokální verze Gearyho kritéria  $C$ , která je však obtížně interpretovatelná (Lee, Wong 2001) díky jejím distribučním vlastnostem.

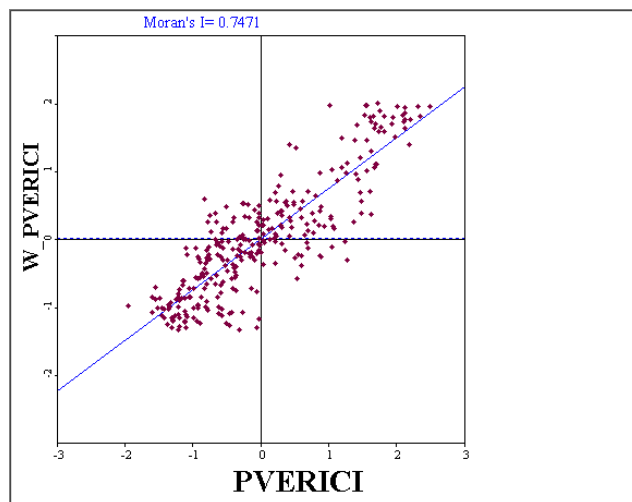
$$C_i = c_i * \sum_j w_{ij} * (r_i - r_j)^2$$

Podobně lze používat lokální verzi  $G$  indikátoru.

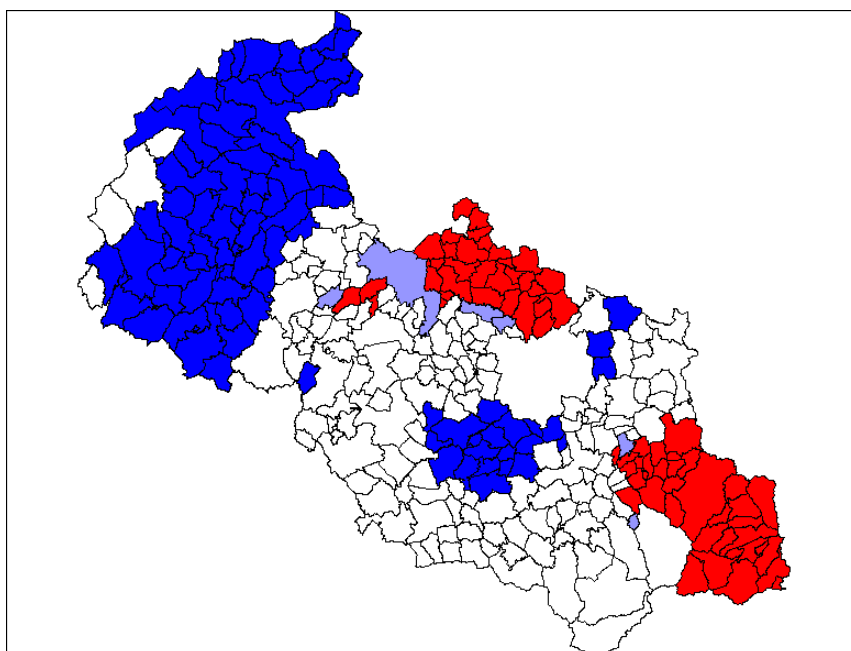
Pro explorační účely je možné sledovat, které oblasti mají neobvyklou hodnotu prostorové autokorelace. K jejímu studiu je možné použít Moranův diagram (Obr. 6-7), ve kterém se zobrazuje hodnota  $Wz$  (průměrná hodnota v sousedství  $z$ ) v závislosti na  $z$  ( $z$  představují standardizované

hodnoty sledované proměnné). V Moranově diagramu se zobrazuje regresní přímka, jejíž sklon odpovídá hodnotě Moranova I kritéria (Spurná 2008).

Na základě výpočtu LISA můžeme provést kategorizaci sledovaných areálů podle charakteru prostorové závislosti do čtyř skupin, které odpovídají čtyřem kvadrantům v Moranově diagramu (obr. 6-8). Prostorové shluky vykazující nadprůměrné či podprůměrné hodnoty proměnné v určité jednotce souhlasně s jejím okolím se v grafu nalézají v pravém horním (hot spots, hodnota vysoká-vysoká) a levém dolním (cold spots, hodnota nízká-nízká) kvadrantu. To svědčí o vysoké autokorelaci. Naopak areály identifikované v levém horním (LH) nebo pravém dolním (HL) kvadrantech jsou charakteristická existencí nízké hodnoty obklopené vysokými a naopak.



Obr. 6-7 Moranův diagram pro podíl věřících v populaci obcí Moravskoslezského kraje (data SLDB 2001, program GeoDa)



Obr. 6-8 Vyznačení shluků hodnot na hladině významnosti 0.05 pro podíl věřících (červená HH, sytě modrá LL, světle modrá LH nebo HL)

## 6.6 Multivariační techniky v prostorových aplikacích

Multivariační metody (metody vícerozměrné statistické analýzy) jsou spojeny se statistickou analýzou vícerozměrných dat. Pomocí vícerozměrné statistické analýzy se snažíme o popis vztahů mezi proměnnými a toto zkoumání probíhá pro všechny vztahy současně. S rostoucím počtem proměnných roste i složitost úlohy.

Uplatnění multivariačních metod pro prostorové analýzy se nachází především v :

- 1) redukci množství dat a průzkumu multidimenzionálního atributového prostoru s cílem identifikovat malý počet zajímavých subdimenzí (resp. kombinací atributů), které pak mohou být zkoumány z prostorového hlediska (uplatnění klasických multivariačních metod a následně vizualizace výsledků a jejich interpretace).
- 2) Průzkumu prostorových textur (vzorů) a vztahů
- 3) Prostorové klasifikaci a diskriminaci.

Klasické multivariační metody pracují s kvantitativními (přesněji poměrovými) daty, existují však i metody či jejich modifikace pro kategorizovaná, nominální či pořadová data. Problémem je však především potřeba standardizace dat pro většinu metod (veličiny v modelu mají typicky různé rozsahy hodnot, různé typy distribuce) a správný způsob provedení standardizace, aby byly splněny metodické požadavky příslušné metody a přitom byla zachována i jistá variabilita veličin, která je předmětem našeho zájmu a podle které dělíme a organizujeme multidimenzionální atributový prostor.

Multivariačních metod je celá řada, k těm základním patří shluková analýza (kapitola 6.6.1), analýza hlavních komponent (kapitola 6.6.2), faktorová analýza, diskriminační analýza, kanonická korelační analýza.

### 6.6.1 Shluková analýza

Shluková analýza je společný název pro celou řadu metod, jejichž cílem je využití informací z analýzy vícerozměrných dat k rozřídění množiny objektů do několika relativně homogenních podsouborů, označených jako shluky (*clustery*). Objekty uvnitř shluků mají být co nejvíce podobné a objekty patřící do různých shluků co nejvíce rozdílné. Podobnost mezi objekty je uplatněna jako kritérium pro tvorbu shluků objektů. Nejdříve se určí znaky definující podobnost a ty se dále sdružují do podobnostních měr. Meziobjektová příbuznost se měří různými prostředky, které se dají zpravidla zahrnout do jedné ze tří elementárních kategorií a to míry korelace, míry vzdálenosti a míry asociace. Korelační a vzdálenostní míry jsou míry metrických dat, zatímco asociční míry jsou určeny spíše pro nemetrická data.

#### Korelační míry

K základní korelačním mírám patří běžně používané korelační koeficienty. Pearsonův párový korelační koeficient  $r$  se používá pro poměrová data a je vyžadována normalita distribucí obou porovnávaných veličin (zde charakteristiky 2 objektů, které se mají shlukovat). Objekty jsou si tím podobnější, čím je jejich párový korelační koeficient větší a bližší jedničce.

Dalším používanou mírou je Spearmanův korelační koeficient, který nevyžaduje splnění žádných předpokladů a pro výpočet korelace používá rozdíl v pořadí členů párů ve variačních řadách. Výpočet lze najít např. v Kaňok 1996 či Meloun et al (). S výhodou ho lze uplatnit i pro pořadová (ordinální) data.

#### Míry vzdálenosti

Míry vzdálenosti představují nejvíce používané míry (Meloun et al. 2005). Vzdálenosti jsou měřeny v prostoru, jehož souřadnice jsou ale představovány hodnotami měřených znaků objektů, nikoliv klasickými souřadnicemi. Nezbytnou podmínkou je standardizace těchto znaků (sjednocení měřítek), jinak dochází ke značně odchýleným výsledkům. Největší problémy se pak vyskytují u čtverce Euklidovské vzdálenosti. K běžným mírám vzdálenosti patří:

- Euklidovská vzdálenost (geometrická metrika) – vypočte se pomocí Pythagorovy věty z rozdílu souřadnic

$$d_E(x_k, x_l) = \sqrt{\sum_{j=1}^m (x_{k,j} - x_{l,j})^2}$$

- čtverec Euklidovské vzdálenosti – využívá se pro Wardovu metody shlukování
- Manhattanská vzdálenost (Hammingova metrika) – pravoúhlá vzdálenost. Meloun et al. (2005) upozorňují, že před použitím této vzdálenosti se musíme přesvědčit, že znaky spolu nekorelují. Jestliže tento předpoklad není splněn, shluky jsou nesprávné.

$$d_H(x_k, x_l) = \sum_{j=1}^m |x_{kj} - x_{lj}|$$

- zobecněná Minkovského metrika – vyšší hodnota z zdůrazňuje odchylky mezi vzdálenými objekty

$$d_M(x_k, x_l) = \sqrt[z]{\sum_{j=1}^m |x_{kj} - x_{lj}|^z}$$

- tětiová vzdálenost (chord distance). V případě, že se použijí tři znaky, je tětiová vzdálenost přímkou vzdáleností dvou bodů na povrchu koule s jednotkovým poloměrem a počátkem v těžišti.

$$d_{CH}(x_k, x_l) = \sqrt{2 \left[ 1 - \frac{\sum_{j=1}^m x_{kj} x_{lj}}{\sum_{j=1}^m x_{kj}^2 \sum_{j=1}^m x_{lj}^2} \right]}$$

- Mahalanobisova metrika – na rozdíl od předchozích využívá popis závislosti mezi znaky (ve výpočtu je zahrnuta kovarianční matice C). Výsledek si můžeme představit jako výpočet vzdálenosti bodů v prostoru, jehož osy nemusí být ortogonální. Vysoce korelovaná selekce znaků může skrytě převážit celý soubor znaků shlukování.

$$d_{Ma}(x_k, x_l) = \sqrt{(x_k - x_l)^T C^{-1} (x_k - x_l)}$$

### Míry asociace

Míry asociace se využívají pro výčtová (nominální) data. K základním koeficientům podobnosti patří (Meloun et al. 2005) Sokalův-Michenerův koeficient asociace, Russelův-Raoův koeficient asociace, Jaccardův koeficient, Hamannův koeficient asociace, korelační koeficient, Rogersův a Tanimotův koeficient asociace, Sørensenův koeficient asociace, míra genetické vzdálenosti a Growerův koeficient podobnosti.



Roztříděním do několika podsouborů rozumíme klasifikaci, která vede k vytvoření systému tříd. Na závěr shlukovací analýzy se proto provádí charakterizace (popis) jednotlivých tříd (tj. shluků) a interpretace. Tímto způsobem lze i významně snížit dimenzionalitu úlohy tak, že původní sadu proměnných nahradíme příslušností k nové třídě. Shlukovací metody jsou úspěšné především v situacích, kdy objekty mají tendenci se seskupovat do přirozených tříd, než v případě náhodného rozmístění objektů v atributovém prostoru.

Doporučuje se před zahájením shlukovací analýzy prozkoumat znaky, které se hodljají použít, a vypustit nevýznamné znaky. Rovněž se doporučuje identifikovat a vyloučit odlehle hodnoty, které mají na shlukovací proces neblahý vliv a zhoršují strukturu dat (Meloun et al. 2005).

Shlukovací analýzu provádíme zpravidla na množině objektů, kde každý objekt je popsán řadou znaků (veličin). Takový postup označíme jako Q-techniku shlukování. Oproti tomu R-technika shlukování vychází z analýz množiny znaků, charakterizovaných prostřednictvím objektů. Podobný duální přístup se uplatňuje i v dalších multivariačních metodách např. faktorové analýze.

Shlukovací analýzy je možné dělit podle různých kritérií, např. zda je na začátku určen počet shluků nebo se má v průběhu řešení nalézt optimální počet shluků, nebo podle výsledné struktury skupin objektů, která je uvedena v tabulce 6-2.

Tab. 6-2 Základní rozdělení metod shlukové analýzy

Skupina	Metoda	Poznámka
Hierarchické	aglomerační (sdružovací)	Postupným seskupováním vytváří stromovou strukturu od jednotlivých objektů až po 1 shluk
	divizní (rozdělovací)	rozdělují počáteční celkový shluk do hierarchického systému dílčích skupin či objektů
Nehierarchické	optimalizační	
	analýzy modů	

### Hierarchické shlukovací postupy

Hierarchická struktura se zobrazuje pomocí dendrogramu. U aglomeračního shlukování se dva objekty s nejmenší vzdáleností spojí, vytvoří shluk a provede se přepočítání vzdáleností všech objektů k novému shluku. Následně se opět hledá nejbližší dvojice. Postup se opakuje tak dlouho, dokud nevznikne ze všech dat jeden shluk (případně dokud není dosaženo požadovaného počtu shluků). V případě divizního shlukování se vychází z celkového shluku (všechna data) a postupně se shluk dělí.

Při shlukování se volí vhodná metoda vyjádření vzdáleností a shlukovací procedura.

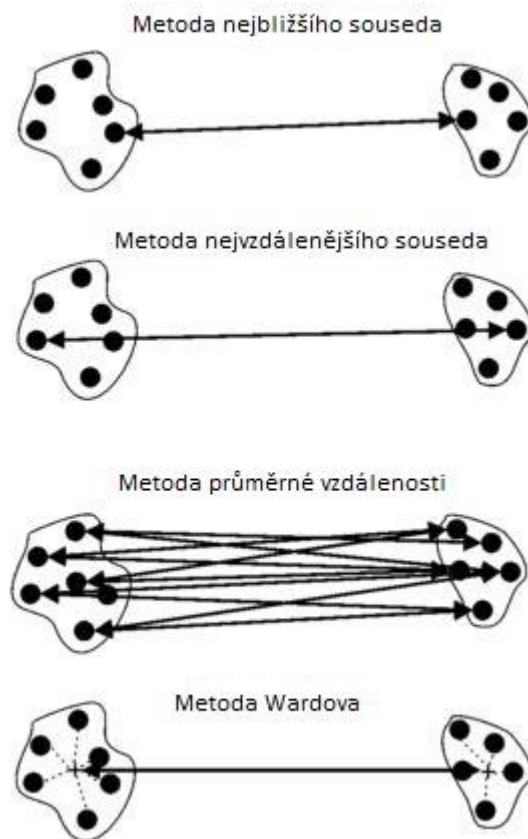
K běžným shlukovacím procedurám patří:

- metoda nejbližšího souseda – pár pro shlukování se vybere podle nejmenší vzdálenosti
- metoda nejvzdálenějšího souseda - pár pro shlukování se vybere podle největší vzdálenosti
- metoda průměrné vzdálenosti – vychází se z průměrné vzdálenosti všech objektů v 1.shluku ke všem objektům ve 2.shluku.
- Wardova metoda – kritériem je minimalizace heterogenity shluků. V každém kroku se spočítá přírůstek součtu čtverců odchylek, vzniklý sloučením shluků. Spojí se ty shluky, které mají minimální hodnotu přírůstku (Meloun et al. 2005).

$$VSS = \sum_{j=1}^m \sum_{i=1}^k (x_{ij} - \bar{x}_j)^2$$

- metoda těžiště
- metoda mediánová

Jejich základní přehled uvádí Turčan (2002), podrobnější popis vybraných metod a jejich algoritmizaci lze nalézt v (Lukasová, Šarmanová 1985) či Meloun et al. (2005).



Obr. 6-9 Nejčastěji užívané metriky shlukování (Meloun et al. 2005)

Příkladem prostorové aplikace shlukové analýzy, která však nevyužívá prostorového složky popisu dat, je analýza vztahu mezi cenou práce a nezaměstnaností (Flek 2000). Ke zpracování byla použita data o míře nezaměstnanosti a průměrných mzdách v okresech ČR ze I. až IV. čtvrtletí každého roku (postupně zpracovány roky 1992-99). Cílově byly vymezeny 3 shluky (metoda K-means) - 1.skupina zahrnuje okresy s nejnižší mírou nezaměstnanosti a nejvyššími průměrnými mzdami, 3.skupina s vysokou mírou nezaměstnanosti a nízkými mzdami, 2.skupina ostatní okresy. Provedená územní interpretace ukazuje na časový vývoj v diferenciaci situace - v roce 1992 byly primární rozdíly mezi českými a moravskými okresy, které však postupně ustoupily do pozadí a byly nahrazeny známou diferenciací Střední Čechy x SZ Čechy a S Morava (především okres Bruntál). Nevyhraněná 2.skupina zahrnuje především okresy východočeské, kde lze do budoucna očekávat nárůst ekonomických a sociálních problémů.

K hierarchickým aglomeračním metodám patří i metoda **minimálního překlenovacího stromu** (*minimum spanning tree*). Překlenovací strom je tvořen soustavou linek (hran dle teorie grafů) mezi všemi pozorováními, které jsou tvořeny tak, aby každé pozorování (uzel) byl spojen s libovolným jiným a přitom v síti nevznikaly smyčky (spojitý a acyklický graf je požadavkem teorie grafů na strukturu typu strom). Délka linky je vyjádřením nepodobnosti obou spojených pozorování. Minimální překlenovací strom je strom uvedených podmínek s minimální délkou. Prakticky jde o minimální kostru grafu. Bailey, Gatrell (1995) popisují prostorovou aplikaci této metody, kdy předpokládají, že sada zkoumaných objektů má prostorovou lokalizaci a doporučují vytvořený graf (minimální překlenovací strom) zobrazit na mapě, aby se mohla prověřit hypotéza, že pozorování blízká v

atributovém prostoru jsou také blízka v geografickém prostoru, a indikovaly se příslušné odchylky. Uvedená metoda se však využívá i ke stanovení části grafu (např. uliční sítě), která je dosažitelná v rámci stanoveného limitu (časového, vzdálenostního).

Řada prostorových aplikací může vyžadovat, aby vytvořené shluky byly prostorově spojitě, tedy aby se vytvářely shluky z geograficky blízkých objektů. Tradiční multivariační metody tento požadavek nezohledňují. Nejjednodušší cestou je zařazení souřadnic objektů do sady pozorování jako dvě další, nové proměnné, je ale zjevné, že optimalizace vzdálenosti podél obou souřadnicových os bude probíhat do jisté míry nezávisle (to lze řešit zavedením prostorové indexace) a celkově samozřejmě jsou souřadnice postaveny na roveň ostatních atributů a není tedy zajištěno, že vzniknou geograficky homogenní shluky.

Výhodnější možnost představuje úprava algoritmu tvorby shluků s přihlédnutím k prostorovým vztahům. V řadě případů - např. administrativního dělení území - můžeme vyžadovat, aby areály spojené do jednoho shluku spolu sousedily. Nabízí se využití matice sousednosti, která obsahuje informaci o sousedství mezi jednotlivými páry areálů. Můžeme tedy např. při aglomerační metodě tvorby shluků připojovat do shluku jen takové areály, které sousedí s některým z již do shluku zahrnutých areálů.

### **Prostorové hierarchické shlukování**

Metoda hierarchického prostorové shlukování je popsána např. v Carvalho et al. (2009).

Algoritmus zahrnuje následující kroky:

1) Určení sousedství jednotlivých areálů – zpravidla se používají topologická sousedství typu královna.

2) Spočítá se vektor vzdáleností (ne geografických, ale mezi vektory proměnných/indikátorů) mezi všemi páry tvořenými sousedními areály a vybuduje se matice blízkosti (symetrická).

3) nalezne se pár sousedů, který má nejmenší vzdálenost mezi nimi. Tento pár se seskupí do jednoho shluku.

4) V nejjednodušším případě je pro definici nového shluku nutné kombinovat seznamy sousedů. Proto bude nový seznam sousedů vytvořen spojením seznamu sousedů města A a seznamu sousedů města B (provedeme sjednocení obou relací).

5) Pro nových  $N-1$  shluků musí být aktualizována matice blízkosti. Aktualizace matice blízkosti (nebo vzdálenosti) závisí na metodě shlukování. Například pro metodu nejbližšího souseda je vzdálenost mezi dvěma shluky I a J minimální vektor vzdáleností mezi všemi dvojicemi vektorů proměnné ve dvou shlucích. Na druhou stranu pro metodu nejvzdálenějšího souseda je vzdálenost mezi dvěma shluky maximální vektor vzdáleností mezi všemi páry vektorů.

6) Opakujeme kroky 3 až 5, dokud zbude jen jeden shluk, který bude obsahovat všech areály.

Výsledkem procesu je strom popisující postupně shlukování.

Algoritmus, který Carvalho a kolektiv v dokumentu používá, vykazuje značné rozdíly proti tradičním hierarchickým shlukovacím algoritmům. Díky tomu se nutně nemusejí shodovat s tradičním shlukováním (neprostorovým).

### **Nehierarchické shlukovací metody**

#### **Metoda nejbližších těžišť K-means**

Meloun a kol. (2005) uvádí, že metoda nejbližších těžišť poskytuje pouze jediné řešení pro zadaný počet požadovaných shluků. Počet shluků musí být předem zadán uživatelem. Postup je založen na nejbližším těžišti, kdy je objekt zařazen do shluku s nejmenší vzdáleností mezi objektem a těžištěm

shluku. Konkrétní technika zařazení objektu závisí na dostupné informaci. Jsou-li těžiště shluků známá, mohou být specifikována v datech a zařazení objektu je založeno na nich. Jinak jsou těžiště shluků určována iteračním výpočtem z dat.

Princip metody nejbližších těžišť (K-means) spočívá v rozdělení  $n$  objektů o  $m$  znacích do  $k$  shluků tak, že mezishluková suma čtverců je minimalizována. Jelikož počet možných uspořádání je enormně veliký, nelze očekávat vždy jediné a nejlepší řešení. Algoritmus nalezne vždy spíše optimum lokální než globální. Jde o takové uspořádání shluků, kdy přemístění objektu z jednoho shluku do druhého nezpůsobí snížení sumy čtverců. Algoritmus pracuje iterativně, startuje vždy z jiného počátečního uspořádání. Nakonec vybere vhodné řešení ze všech možných dosažených uspořádání shluků.

## 6.6.2 Analýza hlavních komponent

Často využívanou multivariační metodou je analýza hlavních komponent. Jejím cílem je redukce původního počtu popisovaných proměnných novými veličinami (umělými), označenými jako komponenty, které shrnují informaci o původních proměnných za cenu minimální ztráty informace. Tyto komponenty jsou vzájemně nezávislé a jsou seřazeny podle svého příspěvku k vysvětlení celkového rozptylu pozorovaných proměnných.

Metoda je citlivá na změnu měřítka, proto se provádí normalizace původních proměnných.

Vlastnosti hlavních komponent jsou takové, že 1.komponenta vysvětluje největší množství rozptylu, 2.menší a podíly vysvětleného rozptylu se u dalších komponent zpravidla rychle snižují.

Hlavním komponentám se snažíme při interpretaci přiřadit nějaký reálný význam. Mají charakter faktorů, které stojí v pozadí a reprezentují zobecněné vlivy, vyvolávají variabilitu a ovlivňují strukturu závislosti proměnných. Při interpretaci využíváme především korelace s původními proměnnými.

Analýza hlavních komponent může být chápána jako transformace z původního do nového souřadnicového systému, jehož osy jsou tvořeny hlavními komponentami. Osy procházejí směry maximálního rozptylu, protože podmínka nezávislosti komponent vede ke kolmosti os.

## 6.7 Regresní modely

Častým cílem modelování je hledání vztahu mezi hodnotami atributů  $\mathbf{z}_i$ , prostorovým uspořádáním areálů  $A_i$  a také hodnotami dalších atributů  $\mathbf{x}_i^T(x_{i1}, x_{i2}, \dots, x_{ip})$  zaznamenanými v areálu  $A_i$ .

Z hlediska uplatnění prostorové složky popisu dat a konstrukce se rozdělují modely na:

- Neprostorové (kapitola 6.7.1)
- Prostorové (kapitola 6.7.2)
- Generalizované (kapitola 6.7.3)

### 6.7.1 Neprostorový regresní model

Neprostorový regresní model využívá jednoduchý standardní regresní model, kde se nevyužívají prostorové vztahy mezi areály. Z hlediska zpracování představují sady atributů naměřených v jednotlivých místech sady proměnných (vektory)  $X$ , které vystupují v regresní rovnici a na základě které se vypočítává nová hodnota  $Y$ . Regresní model obecně vychází ze vztahu:

$$Y = X * \beta + \varepsilon$$

kde

- |       |   |
|-------|---|
| $Y_i$ | vektor $n \times 1$ náhodných závisle proměnných v místě $i$ (výsledná hodnota) |
| $X_i$ | matice $n \times p$ nezávisle proměnných zjištěných v místě (areálu $i$ )       |
| $n$   | počet sad pozorování (počet míst, resp. vzorků)                                 |

- $p$  počet proměnných zjištěných u 1 místa ( a uplatněných zde ve výpočtu)  
 $\beta$  vektor  $p \times 1$  parametrů  
 $\varepsilon$  vektor  $n \times 1$  náhodných proměnných, které reprezentují odchylky od trendu. Odchylky od trendu se označují jako fluktuace, disturbance, někdy jsou interpretovány jako chyby. Střední hodnota  $\varepsilon$  musí být nulová, střední hodnota ( $\varepsilon^* \varepsilon^T$ ) je rovna rozptylu.

Tento základní model se uplatňuje ve všech lineárních regresních modelech. Obsahuje 2 základní komponenty. První představuje odhad průměrné hodnoty  $\mathbf{Y}$  na základě známých veličin  $\mathbf{X}$  a samozřejmě odhad parametrů  $\beta$ , které vztahy mezi  $\mathbf{Y}$  a  $\mathbf{X}$  popisují. Druhá komponenta sleduje, jak jsou získané hodnoty  $\mathbf{Y}$  rozloženy kolem střední hodnoty, sleduje odchylky a tedy popisuje distribuci fluktuací.

Standardní model předpokládá konstantní rozptyl fluktuací  $\varepsilon_i$  v celém území, nezávislost hodnoty rozptylu v jednotlivých místech a žádnou prostorovou autokorelaci. K výpočtu hodnot  $\beta$  a  $\varepsilon$  (k optimalizaci modelu) se díky lineárnímu modelu používá metoda nejmenších čtverců. Odhad metodou nejmenších čtverců je však dobrý jen v případě, že rozdělení reziduí je normální.

Standardní vícenásobný regresní model je popisován v řadě statistických učebnic, včetně výpočtu jeho spolehlivosti (koef. determinance apod.) viz např. (Víšek 1998).

Jak již bylo uvedeno, tento základní model předpokládá jen existenci variací 1.řádu a nepřítomnost variací 2.řádu. V řadě případů to však neplatí a rezidua jsou např. prostorově korelována a rozptyl v území se mění. Úprava nekonstantních rozptylů v jednotlivých areálech tak, aby se dosáhlo jeho vyrovnání, může být prováděna v zásadě váženou regresí nebo pomocí transformace měřených hodnot. Transformace jsou doporučovány zvláště pro případy, kdy zájmová proměnná tvoří počet nebo podíl (absolutní a relativní hodnoty).

V případě počtu (absolutní data) se předpokládá Poissonova distribuce pro tato data a jako transformace se využívá odmocninná transformace  $\sqrt{y_i}$ , která stabilizuje rozptyl (je pak konstantní), nebo logaritmická transformace  $\text{LN}(y_i)$ , která má výhodu v linearizaci násobných vztahů. Sřetáváme se zde s problémem úpravy dat, aby vyhověly současně 2 podmínkám, tj. požadavku konstantního rozptylu v celé oblasti a požadavku lineárního vztahu mezi  $\mathbf{X}$  a  $\mathbf{Y}$  (mezi středními hodnotami  $\mathbf{Y}$  a nezávislými proměnnými ( $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p$ )).

V případě podílu z celku (relativní data) je situace složitější, protože se musí dodržovat rozsah podílu  $<0,1>$ . Jako vhodné uvádí Bailey, Gatrell (1995) úhlovou transformaci a logitovou transformaci, ale upozorňují, že výsledky nejsou zcela uspokojivé, navrhují další úpravu transformace a použití vážené regrese a především generalizovaných lineárních modelů, kde se již počítá s efekty 2.řádu a prostorovou korelací reziduí. Řešení odvozování parametrů pomocí metody maximální věrohodnosti někteří autoři nedoporučují z praktických důvodů (výpočetní náročnost, tvorba samostatného modelu pro každý výpočet).

Úhlová transformace:

$$\frac{1}{\sin \sqrt{y_i}}$$

Logitová transformace:

$$\log\left(\frac{y_i}{1 - y_i}\right)$$

Logistická distribuční funkce se používá z důvodu dobré aproximace normální distribuce (max. odchylka 0.11). Funkce má zrychlující se růst až do inflexního bodu, pak se růst zpomaluje až se v nekonečno zastaví.

Příklad aplikace jednoduchého regresního modelu s logistickou funkcí pro oblast trhu práce uvádí Nijkamp, Rietveld (1987). Základní model epidemického šíření infekční choroby v populaci byl použit pro modelování šíření informací o inovaci nového produktu či procesu. Rychlost difúze závisí na řadě věcí např. na komunikační síti, četnosti kontaktů a ochotě příjemce informace aplikovat inovaci. Pokud budeme předpokládat stejnou pravděpodobnost kontaktů mezi lidmi, dále že každý se může stát potenciálním příjemcem a že frekvence kontaktů ani chování příjemců se nemění během působení, můžeme použít následující model:

$$p(t) = \frac{1}{1 + e^{(a-b*t)}}$$

kde

t	čas
p(t)	podíl příjemců v čase t
a	parametr určující počáteční bod
b	parametr určující poměr růstu populace (příjemců). Musí platit $b > 0$ .

Je třeba dodat, že základní tvar logistické křivky má v čitateli místo jedničky **k** (**k** je hodnota, ke které se funkce blíží, ale nepřekročí ji). Parametr **b** znamená rychlost, s jakou se křivka blíží limitě.

Základní předpoklady jsou značně jednoduché a nijak se např. nezohledňuje prostorový aspekt. Model je možno rozšířit zahrnutím role vzdálenosti. Předpokládáme, že kontakty mezi blízko žijící osobami jsou častější a podobně i podíl příjemců nezávisí jen na čase, ale i na vzdálenosti, kde byla inovace přijata. Vzdálenost lze implementovat náhradou  $a = a_1 + a_2 d$  ( $a_2 > 0$ ), podobně parametr  $b = b_1 + b_2 d$ .

**Vážená metoda nejmenších čtverců** je vhodná pro případy heteroskedacity. Každá primární hodnota v regresním vztahu je vážena inverzní hodnotou příslušného rozptylu, takže proměnné s větším rozptylem získávají v regresní rovnici menší váhu.

## 6.7.2 Prostorový regresní model

Prostorový regresní model umožňuje zahrnout do jednoduchého regresního modelu předpoklad prostorové korelace reziduí. V tomto případě se místo vektoru odchylek  $\epsilon$  do rovnice dosazuje vektor odchylek **U**, který má popsánu kovarianční matici **C** (a střední hodnota **U** je nulová).

Kovarianční matici **C** lze odvodit pomocí strukturální analýzy pro uvedenou oblast, za předpokladu stacionárního procesu 2.řádu s nulovou střední hodnotou.

Předpoklad stacionarity 2.řádu je obtížné dodržet pro areálová data. Dokonce i když je zdrojový proces stacionární, agregace do nepravidelných oblastí způsobuje, že rozptyl i kovariance není stálá pro celou oblast. Navíc nemáme jednoduché měření vzdálenosti/blízkosti mezi areály (viz typy sousedství v kapitola 6.4.1). Proto se vedle klasické strukturální analýzy uplatňují i jiné postupy, např. interakční schéma. Do základního regresního modelu jsou zahrnuty vztahy mezi proměnnými a jejich hodnotami v sousedství a nepřímo se tak odvozuje matice **C**. Výhodou tohoto postupu je, že není vyžadována stacionarita 2.řádu ani využití Euklidovských vzdáleností. Matice **U** se pak vyjadřuje ve tvaru:

$$Y = X\beta + U \quad U = \rho WU + \epsilon$$

kde

<b>W</b>	standardizovaná matice blízkosti (vazeb mezi areály)
$\epsilon$	vektor nezávislých náhodných odchylek
$\rho$	vektor autokorelačních parametrů

Matici blízkosti lze standardizovat (tedy stabilizovat rozptyl) pomocí C schématu (globální standardizace), W schématu (řádková standardizace, suma na řádku se rovná 1) nebo S schématu (leží mezi C a W). Běžně se doporučuje pro tento případ řádkovou standardizaci (tedy W schéma).

Ukázka standardizace popisuje LeSage (1998).

$$W = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \end{pmatrix} \quad C = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1/2 & 1/2 \\ 0 & 0 & 1/2 & 0 & 1/2 \\ 0 & 0 & 1/2 & 1/2 & 0 \end{pmatrix}$$

Při plné závislosti na okolí by se cílová hodnota Y spočítala z okolních hodnot y

$$\begin{pmatrix} y_1^* \\ y_2^* \\ y_3^* \\ y_4^* \\ y_5^* \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.5 & 0.5 \\ 0 & 0 & 0.5 & 0 & 0.5 \\ 0 & 0 & 0.5 & 0.5 & 0 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{pmatrix}$$

$$\begin{pmatrix} y_1^* \\ y_2^* \\ y_3^* \\ y_4^* \\ y_5^* \end{pmatrix} = \begin{pmatrix} y_2 \\ y_1 \\ 1/2y_4 + 1/2y_5 \\ 1/2y_3 + 1/2y_5 \\ 1/2y_3 + 1/2y_4 \end{pmatrix}$$

FAR – first-order spatial regression model:

$$y = \rho C y + \varepsilon$$

Nakonec se do modelu přidají tradiční vysvětlovací proměnné (finální odhad Y závisí jak na prostoru tak i na nezávislých proměnných Xi)

$$y = \rho C y + X\beta + \varepsilon$$

To je i model., který používá GeoDa a označuje ho jako smíšený regresně-prostorový autoregresní model (spatial lag model).

$$Y = X\beta + \rho WY + \varepsilon$$

Po dosazení a úpravě získáme:

$$Y = X\beta + \rho WY - \rho WX\beta + \varepsilon$$

kde

člen  $X\beta$  vyjadřuje celkový trend

člen  $\rho WY$  vyjadřuje vliv okolí

člen  $\rho W X\beta$  vyjadřuje místní (lokální) trend

Tento model se označuje jako autokorelační chybový model a patří do skupiny tzv. simultánně autoregresních modelů (SAR), kde vystupuje jediný interakční parametr  $\rho$  ve významu autokorelačního koeficientu. Právě odhad tohoto parametru vyžaduje zvláštní iterační postupy, protože odhad metodou nejmenších čtverců pro tento parametr je vychýlený a nekonzistentní. Tiefelsdorf

(2000) upozorňuje na problémy s odvozením členu  $\rho WY$  v rovnici, protože je závislý jak na nezávislé proměnné  $X$ , tak i na hodnotách odchylek  $U$ .

Model SAR může být upraven vypuštěním členu místního trendu ( $\rho WX\beta$ ) nebo dokonce i vyloučením členu popisujícího globální trend. Pak tedy  $Y = \rho WY + \varepsilon$ . Další varianty jsou možné.

Program GeoDa rovněž implementuje tzv. prostorový autoregresní model s hodnocením závislosti chyb (spatial error model). Ten hodnotí prostorovou závislost chyb modelu a optimalizuje model vůči chybám. Zavádí nový parametr  $\lambda$  jako koeficient prostorové korelace chyb modelu.

$$Y = X\beta + \varepsilon \quad \varepsilon = \lambda W\varepsilon + u$$

kde  $\varepsilon$  je vektor prostorově autokorelovaných chyb  
 $u$  je vektor chyb

Dalším typem prostorových regresních modelů je **podmíněný autoregresní model (CAR)**:

$$Y = X\beta + U$$

$$U = \rho W^{(1)}U + \rho W^{(2)}U + \dots + \varepsilon$$

kde

$W^{(k)}$  je matice blízkosti pro řádek  $k$  (pro různé prostorové kroky, tedy sousedství 1.řádu, 2.řádu atd.)

Příklad: Odhad kriminality z příjmů a hodnoty domů (LeSage 1998) – srovnajte výsledek běžné regresní analýzy a prostorové regresní analýzy. Rozdíl 10% (z 0,55 na 0,65) ve vysvětlení variability je vysvětleno prostorovou závislostí.

Ordinary Least-squares Estimates

```
Dependent Variable =      Crime
R-squared          =      0.5521
Rbar-squared       =      0.5327
sigma^2            =     130.8386
Durbin-Watson     =      1.1934
Nobs, Nvars        =      49,      3
```

```
*****
Variable           Coefficient      t-statistic    t-probability
constant           68.609759      14.484270     0.000000
income             -1.596072     -4.776038     0.000019
house value        -0.274079     -2.655006     0.010858
```

Spatial autoregressive Model Estimates

```
Dependent Variable =      Crime
R-squared          =      0.6518
Rbar-squared       =      0.6366
sigma^2            =     95.5032
log-likelihood     =     -165.41269
Nobs, Nvars        =      49,      3
```

```
*****
Variable           Coefficient      t-statistic    t-probability
constant           45.056251      6.231261     0.000000
income             -1.030641     -3.373768     0.001534
house value        -0.265970     -3.004945     0.004331
rho                0.431381      3.625340     0.000732
```

Základní 2 přístupy k modelování prostorové regrese zahrnují:

- Model prostorové expanze (Casetti in LeSage)
- geograficky vážená regrese

Model prostorové expanze využívá:

$Y$  závisle proměnná,  $X$  nezávisle proměnná,  $Z$  poloha (zpravidla souřadnice  $Z_x$  a  $Z_y$ )

Parametry v modelu jsou funkcí polohy (souřadnic). Odhaduje se vektor parametrů  $\beta$  (je jich  $2k$ , protože  $k$  je počet nezávisle proměnných a máme 2 souřadnice)



$$y = X\beta + \varepsilon$$

$$\beta = ZJ\beta_0$$

Model vykazuje změny napříč územím (jiné parametry, tedy jiné vztahy, se ukazují v různých částech území).

### Geograficky vážená regrese (GWR)

Při použití regresní analýzy můžeme získat jiný funkční vztah mezi dvěma proměnnými v závislosti na charakteru určitého regionu atd. Například vztah mezi cenou bytu a hustotou zalidnění v jeho okolí může být prostorově velmi odlišný, neboť na cenu bytu má vliv množství dalších charakteristik, které uvedený vztah modifikují (Spurná 2008).

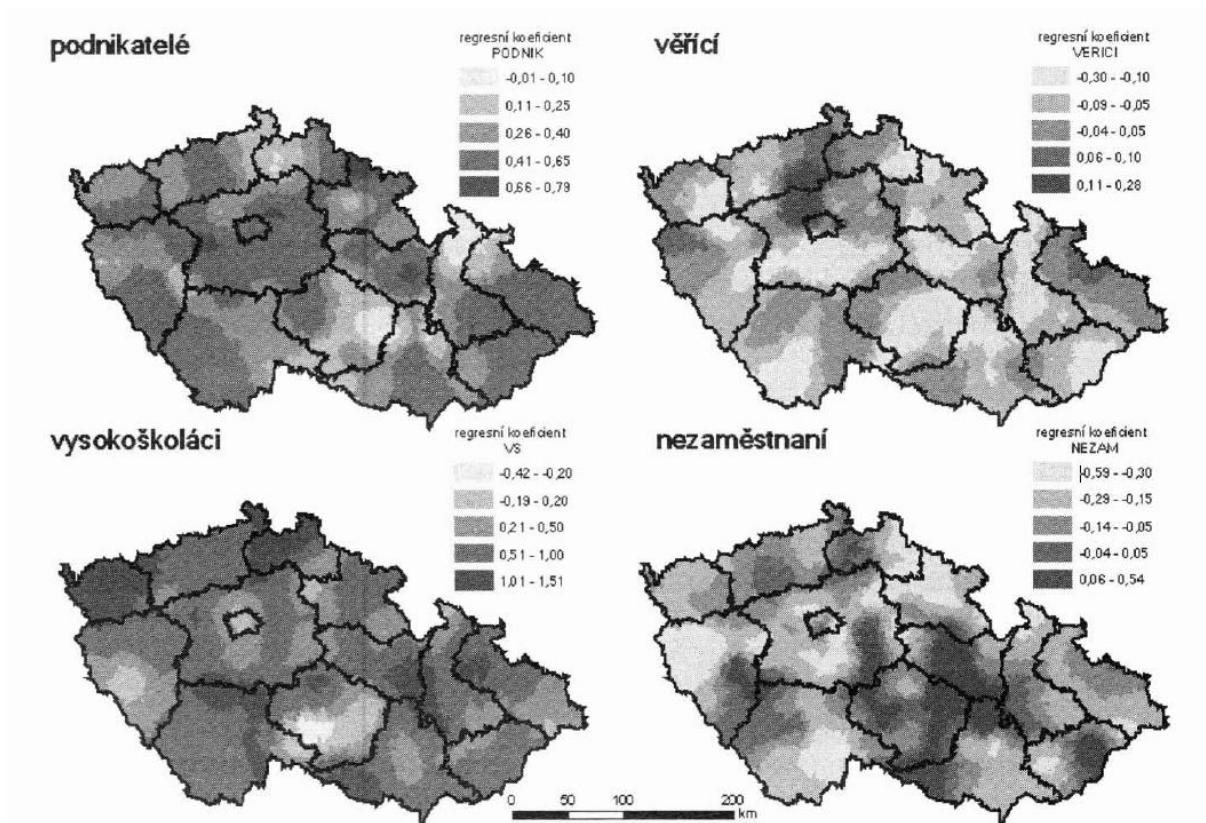
Prostorové regresní metody nám dovolují odhalit např. přítomnost tzv. Simpsonova paradoxu, který označuje obrácení závislosti na lokální úrovni, kdy parciální vztahy jsou silnější než vztah původní. I když je Simpsonův paradox obvykle demonstrován na neprostorových datech, kde dochází k agregaci subpopulací, příklady z geografické literatury (Spurná 2006) dokazují, že platí rovněž pro agregovaná prostorová data, a poukazují tak na nebezpečí spojené s analýzou dat za větší území. Může se totiž stát, že nalezený závěr bude pravým opakem reálných vztahů, které platí na lokální úrovni.

Metoda GWR tedy předpokládá možnost existence prostorových odlišností ve vztazích dvou a více proměnných a poskytuje způsob, jak tyto odchylky měřit. V rámci GWR je regresní analýza provedena pro každý regresní bod zvlášť, čímž jsou získány lokální regresní parametry. Vynesením výsledných odhadů lokálních regresních parametrů do mapy je následně přehledně znázorněn charakter zkoumaného vztahu.

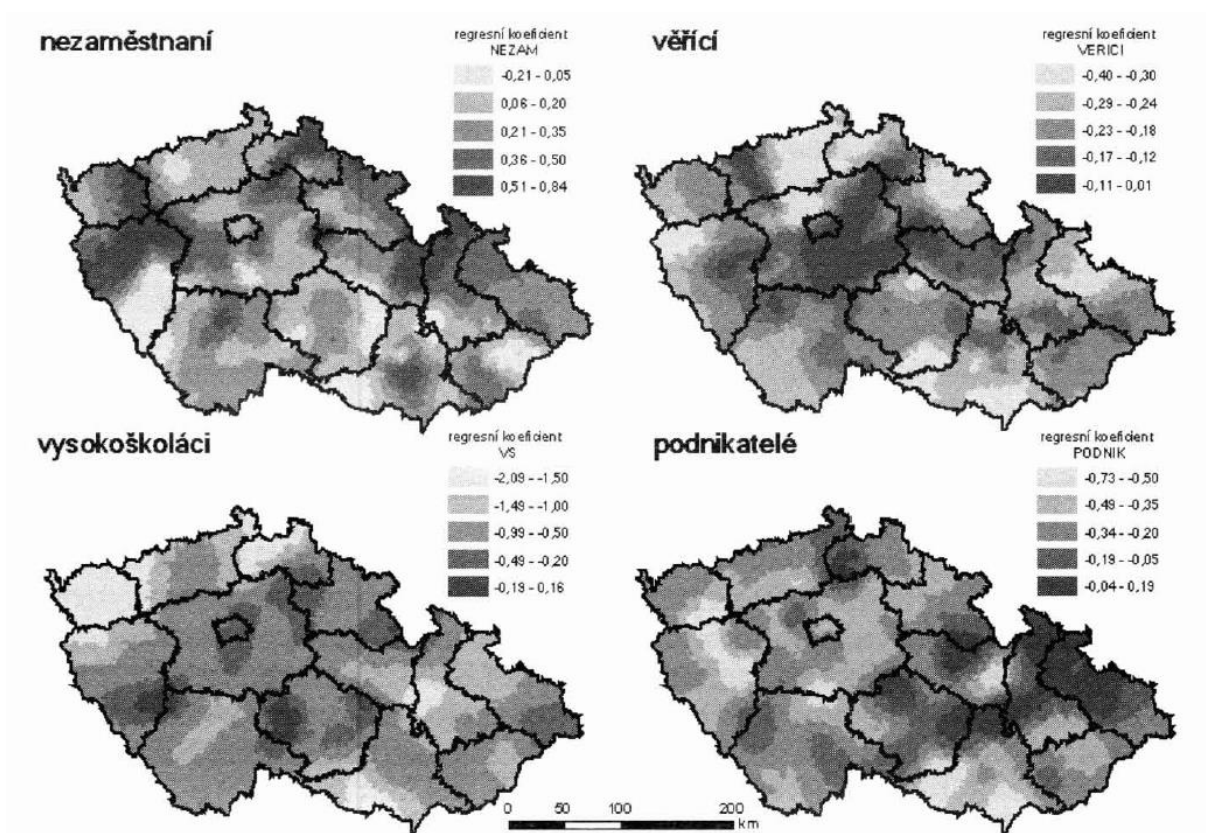
Metoda využívá jadrových odhadů – ve fixní i adaptivní podobě.

Tab. 6-3 Srovnání výsledků neprostorové mnohonásobné lineární regrese a geograficky vážené regrese (Spurná 2008).

Vysvětlující proměnné	Mnohonásobná lineární regrese			Geograficky vážená regrese				
	r <sup>2</sup>	Beta	sig. 1	r <sup>2</sup>	b <sub>min</sub>	b <sub>medián</sub>	b <sub>max</sub>	sig. 2
<b>ODS</b>	0,369	0,306	0,000	0,492	-0,014	0,409	0,792	0,060
PODNIK		-0,336	0,000		-0,299	-0,067	0,275	0,000
VERICI		0,211	0,000		-0,416	0,628	1,505	0,000
VS		-0,150	0,000		-0,585	-0,130	0,542	0,000
<b>KSČM</b>	0,318	0,284	0,000	0,494	-0,217	0,228	0,841	0,000
NEZAM		-0,296	0,000		-0,396	-0,206	0,006	0,000
VERICI		-0,228	0,000		-2,093	-0,767	0,159	0,000
PODNIK		-0,200	0,000		-0,725	-0,295	0,185	0,000



Obr. 6-9 Výsledky GWR pro volební preference ODS v roce 2002 (Spurná 2008)



Obr. 6-10 Výsledky GWR pro volební preference KSČM v roce 2002 (Spurná 2008)

Výsledky regresního modelu mohou být srovnány s výsledky Chowova testu, který měří stabilitu regresních koeficientů v jednotlivých subregionech dat proti celé populaci (Guajarati 2003 in Kouba

2007). Standardní regresní model předpokládá homogenitu regresních koeficientů pro celou populaci. Chowův test (na bázi F testů) ověřuje srovnáním teoreticky zdůvodněných oblastí dat s celkovou populací, zda a nakolik je tento předpoklad oprávněný. Jeho interpretace využívá F pravděpodobnostní distribuci se zvolenou mírou pravděpodobnosti (zde často  $p < 0.01$ ). Pokud je vypočtené F nižší než kritické F, musíme přijmout nulovou hypotézu, že regresní parametry jsou stabilní pro celou populaci. V opačném případě, kdy F je vyšší než kritické F, jsou regresní koeficienty v území nestabilní a můžeme hovořit o nehomogenním území (o přítomnosti různých „prostorových režimů“) (Kouba 2007).

Praktické výsledky ukazuje např. Kouba (2007) ve svém příspěvku k regionalizaci českého stranického systému. V případě Moravy nebyla prostorová nehomogenita v regresních vztazích potvrzena (tedy neprojevil se prostorové rozdíly), v případě Čech se prokázalo, že u sudetských okresů zde získává KSČM zhruba o 2% hlasů více ve srovnání se zbytkem země, naopak ODS je relativně úspěšnější v nesudetských okresech.

### 6.7.3 Generalizované lineární modely

Generalizovaný lineární model zahrnuje jak standardní regresní model (předpoklad normální distribuce odchylek), tak i možnost modelovat distribuci odchylek pomocí binomické či Poissonovy distribuce a také využití log-lineárních modelů. K odvozování parametrů se používá metoda maximální věrohodnosti, implementovaná pomocí techniky označované jako procedura nejmenších čtverců s iteračním vícenásobným vážením.

Generalizovaný lineární model se skládá opět ze dvou základních komponent - pro komponentu odchylek se používá distribuce normální, biomická či Poissonova, pro komponentu trendu se používá monotónní spojovací hladká funkce  $g()$ .

$$g(\mu_i) = \mathbf{X}_i^T \boldsymbol{\beta}$$

kde

$\mu_i$	střední hodnota $\mathbf{Y}$ v místě $\mathbf{i}$
$\mathbf{X}_i^T$	transponovaný vektor nezávislých proměnných $\mathbf{X}$ v místě $\mathbf{i}$
$\boldsymbol{\beta}$	vektor parametrů

Spojovací funkce  $g()$  umožňuje nelineární vztah mezi střední hodnotou závislé veličiny  $X$  a nezávislých proměnných  $X$  (resp.  $\mathbf{X}\boldsymbol{\beta}$ ). Distribuce chyb se modeluje zvlášť, což samozřejmě nepředstavuje zcela komplexní a obecné řešení.

Např. pro data typu počet (absolutní data) se používá Poissonova distribuce chyb a logaritmičká spojovací funkce:

$$\log(\mu_i) = X_i^T \boldsymbol{\beta}$$

Pro podíly (relativní data) se používá binomická distribuce chyb a logitová spojovací funkce:

$$\log\left(\frac{\mu_i}{1 - \mu_i}\right) = X_i^T \boldsymbol{\beta}$$

Např. tedy

$$\log\left(\frac{\mu_i}{1 - \mu_i}\right) = \beta_1 + \beta_2 X_i$$

Pomocí vážené regrese se vypočte první odhad parametrů  $\beta$  a ten se postupně zpřesňuje iteračním postupem opět s využitím vážené regrese. Způsob ocenění kvality modelu je popsán např. v Bailey, Gatrell 1995. Cressie (1993) využívá jako spojovací funkci druhou odmocninu a doporučuje použít Freeman-Tukey transformaci pro tento účel, protože vykazuje větší stabilitu.

Nezanedbatelnou výhodou generalizovaného lineárního modelu je možnost využití i nominálních nezávislých proměnných (resp. binárních ve tvaru indikátorů), tak i možnost využití kategorizovaná data.

Regresní modely jsou socioekonomické oblasti poměrně populární, relativně málo však nacházíme příklady prostorových regresních modelů. Přitom právě uplatnění těchto principů představuje významné vylepšení modelů. V oblasti socioekonomických aplikací se vzhledem k vazbě na distribuci obyvatelstva předpokládá výrazný autokorelační efekt, který lze jen stěží do běžných regresních modelů zahrnout. V případě prostorových regresních modelů je přímo implementován a dokonce se nemusí ani popisovat struktura pole.

## Seznam literatury

Amendola et al. (1998)

Amendola, A., Caroleo, F., Coppola, G.: Regional differences in the labour market and economic growth in Italy. [CD-ROM] In *Sbor. ref. 10. EALE konference, Blankenberge (Belgie), 1998, 36 s.*

Aronoff (1989)

Aronoff, S.: *Geographic Information Systems: A Management Perspective. Ottawa, WDL Publications, 1989, ISBN 0-9218404-00-8.*

Bachi (1999)

Bachi, R.: New methods of geostatistical analysis and graphical presentation. Distribution of populations over territories. *New York, Kluwer Academic / Plenum Publishers, 1999, 478 s., ISBN 0-306-45544-7.*

Bailey (1994)

Bailey, T.: A review of statistical spatial analysis in geographical information systems. In *Fotheringham S., Rogerson P. (ed.): Spatial Analysis and GIS. Taylor&Francis Ltd., 1994, s. 13-44, ISBN 0-7484-0101-0.*

Bailey, Gatrell (1995)

Bailey, T., Gatrell, A.: *Interactive spatial data analysis. Essex, Longman Scientific & Technical, 1995, 413 s.*

Bezák (1991)

Bezák, A.: Migračné toky a regionálna štruktúra Slovenska: hierarchická regionalizácia. *Geografický časopis, ročník 43, 1991, číslo 3, s. 193-202.*

Bezák (2001)

Bezák, A.: O regionálnych trhoch práce, nových krajoch a tokoch nezamestnaných. *Geografický časopis, ročník 53, 2001, číslo 4, s. 295-305.*

Birkin et al. (1996)

Birkin, M., Clarke, G., Clarke, M., Wilson, A.: *Intelligent GIS - location decision and strategic planning. Glasgow, Geoinformation International, 1996, 292 s., ISBN 1-899761-25-X.*

Bithel (1996)

Bithel, J.: Statistical methods for analysing point-source exposures. In *Elliott P., Cuzick J., English D., Stern R.(ed.): Geographical & Environmental Epidemiology. New York, Oxford University Press, 1996, s. 221-230, ISBN 0-19-262235-8.*

Brabec et al. (1970)

Brabec, F., Knap, J., Novotný, M., Zach, J.: *Matematické a grafické metody v geografii. Skriptum, Praha, Univerzita Karlova, SPN., 1970, 174 s.*

Brázdil et al. (1992)

Brázdil, R., Kolář, M., Prošek, P., Tarabová, Z., Wokoun, R.: *Statistické metody v geografii - cvičení. Skriptum, Brno, Masarykova univerzita, Přírodovědecká fakulta, 1992, 177 s.*

Bracken (1994a)

Bracken, I.: A surface model approach to the representation of population-related social indicators. In *Fotheringham S., Rogerson P. (ed.): Spatial Analysis and GIS. Taylor&Francis Ltd., 1994, s. 247-259, ISBN 0-7484-0101-0.*

Bracken (1994b)

Bracken, I.: Towards improved visualization of socio-economic data. In *Hearnshaw H., Unwin D. (ed.): Visualization in Geographical Information Systems. UK, John Wiley & Sons Ltd., 1994, s.76-84, ISBN 0-471-94435-1.*

Brewer (1994)

- Brewer, C.: Color Use Guidelines for mapping and visualization. In *MacEachren, A., Taylor, F. (ed.): Visualization in modern cartography. Elsevier Science Ltd., 1994, s. 123-147, ISBN 0-08-042415-5.*
- Brunsdon (1995)  
 Brunsdon, Ch.: Estimating probability surfaces for geographical point data: an adaptive kernel algorithm. *Computers&Geosciences, Elsevier Science Ltd., 21, 7, 1995, s. 877-894, ISSN 0098-3004.*
- Burrough, McDonnell (1998)  
 Burrough, P., McDonnell, R.: Principles of Geographical Information Systems. *Oxford, Oxford University Press, 1998, 336 s., ISBN 0-19-823365-5.*
- Carcach (2000)  
 Carcach, C.: The spatial analysis of ambulance calls for drug overdose in Adelaide. Searching for a link between drug use and property crime. In *Sbor. ref. konference Crime Mapping: Adding Value to Crime Prevention and Control, 21-22.září 2000, Australian Mineral Foundation, Adelaide, 10 s.*  
 Dostupný na WWW: <http://www.aic.gov.au/conferences/mapping/carcach.pdf>
- Carsjens, Knaap (1996)  
 Carsjens, G., Knaap, W.: Multi-criteria Techniques Integrated in GIS Applied for Land Use Allocation Problems. In *Sbor. ref. konference JEC, Volume 1, 1996, s. 575-578, ISBN 90-5199-268-8.*
- Carvalho et al. (2009)  
 CARVALHO, Alexandre Xavier Ywata; ALBUQUERQUE, Pedro Henrique Melo; ALMEIDA JUNIOR, Gilberto Rezende de. GUIMARÃES, Rafael Dantas. Clusterização espacial hierárquica. *Rev. Bras. Biom [online]. São Paulo : v.27, n.3, p.412-443, 2009. Dostupné z WWW: <[http://www.fcav.unesp.br/RME/fasciculos/v27/v27\\_n3/A6\\_Alexandre.pdf](http://www.fcav.unesp.br/RME/fasciculos/v27/v27_n3/A6_Alexandre.pdf)>.*
- Cressie (1993)  
 Cressie, N.: Statistics for spatial data. *USA, John Wiley & Sons Inc., 1993, 900 s., ISBN 0-471-00255-0.*
- DeMers (1997)  
 DeMers, M.: Fundamentals of Geographic Information System. *USA, John Wiley & Sons Inc., 1997, 486 stran, ISBN 0-471-14284-0.*
- Diamond (1996)  
 Diamond, I.: Population counts in small areas. In *Elliott P., Cuzick J., English D., Stern R.(ed.): Geographical & Environmental Epidemiology. New York, Oxford University Press, 1996, s. 96-105, ISBN 0-19-262235-8.*
- Dingemans (1996)  
 Dingemans, P.: Aggregation of Point Information in Thematic Maps. In *Sbor. ref. konference JEC, Volume 1, 1996, s. 721-731, ISBN 90-5199-268-8.*
- Dohnal, Pour (1997)  
 Dohnal, J., Pour, J.: Architektury informačních systémů. *Ekopress, 1997, 301 s., ISBN 80-86119-02-5.*
- Dorling (1994)  
 Dorling, D.: Cartograms for visualizing human geography. In *Hearnshaw, H., Unwin, D. (ed.): Visualization in Geographical Information Systems. UK, John Wiley & Sons Ltd., 1994, s. 85-102, ISBN 0-471-94435-1.*
- Dudorkin (1997)  
 Dudorkin, J.: Operační výzkum. *Skriptum, Praha, ČVUT, fakulta elektrotechnická, 1997, 296 s., ISBN 80-01-01571-8.*
- Dykes (1994)

- Dykes, J.: Area-value data: New visual emphases and representations. In *Hearnshaw, H., Unwin, D. (ed.): Visualization in Geographical Information Systems. UK, John Wiley & Sons Ltd., 1994, s. 103-114, ISBN 0-471-94435-1.*
- Eastman, Jiang (1996)  
Eastman, J., Jiang, H.: Fuzzy measures in multicriteria evaluation. In *Sbor. ref. konference Second International Symposium on Spatial Accuracy Assessment in Natural Resources and Environmental Studies, Fort Collins, Colorado, 21.-23.5.1996, s.527-534.*
- English (1996)  
English, D.: Geographical epidemiology and ecological studies. In *Elliott P., Cuzick J., English D., Stern R.(ed.): Geographical & Environmental Epidemiology. New York, Oxford University Press, 1996, s. 3-13, ISBN 0-19-262235-8.*
- Eurostat (1992)  
Eurostat: Labour Force Survey. Methods and Definitions. *Luxembourg, Office for Official Publications of the European Communities, 1992.*
- Fischer, Nijkamp (1987a)  
Fischer, M., Nijkamp, P.: Spatial labour market analysis: relevance and scope. In *Fischer, M., Nijkamp P. (ed.): Regional Labour Markets. Elsevier Science Publishers, 1987, s.1-33, ISBN 0-444-70323-3.*
- Fischer, Nijkamp (1987b)  
Fischer, M., Nijkamp, P.: Labour market theories: perspectives, problems and policy implications. In *Fischer, M., Nijkamp P. (ed.): Regional Labour Markets. Elsevier Science Publishers, 1987, s.37-52, ISBN 0-444-70323-3.*
- Flek et al. (2000)  
Flek, M. a kolektiv: Cena práce a zaměstnanost v procesu transformace a restrukturalizace české ekonomiky. *Zpráva o výsledcích řešení úvodní etapy výzkumného úkolu, 1993, Výzkumný ústav práce a sociálních věcí, Praha.*
- Flowerdew, Green (1994)  
Flowerdew R., Green M.: Areal interpolation and types of data. In *Fortheringham, S., Rogerson (ed.): Spatial Analysis and GIS. Taylor&Francis Ltd., 1994, s. 121-145, ISBN 0-7484-0101-0.*
- Fortheringham et al. (1997)  
Fortheringham, S., Charlton, M., Brunson, Ch.: Measuring spatial variations in relationships with geographically weighted regression. In *Fischer M., Getis A. (ed.): Recent Developments in Spatial Analysis. Spatial Statistics, Behavioural Modelling and Computational Intelligence. Springer-Verlag Berlin, 1997, s. 60-82, ISBN 3-540-63180-1.*
- Fortheringham et al. (2000)  
Fortheringham, S., Brunson, Ch., Charlton, M.: Quantitative Geography. Perspectives on Spatial Data Analysis. *London, SAGE Publications, 2000, s. 270, ISBN 0-7619-5948-3.*
- Fox (1975)  
Fox, W.: Some practical aspects of time series analysis. In *McCammon, R.:(ed.) Concepts in Geostatistics, New York, 1975, s.70-89, ISBN 3-540-0682-9.*
- Friedmannová (2001)  
Friedmannová, L.: Automatizovaná tvorba barevných škál pro WWW kartografickou vizualizaci v GIS. In *Sbor. ref. konference GIS Ostrava 2001, Ostrava, VŠB - TUO, 2001, 6 s., ISBN 1213-2454.*
- Gatrell (1994)  
Gatrell, A.: Density estimation and the visualization of point patterns. In *Hearnshaw H., Unwin D. (ed.): Visualization in Geographical Information Systems. UK, John Wiley & Sons Ltd., 1994, s. 65-75, ISBN 0-471-94435-1.*
- Goodchild (1988)

- Goodchild, M.: Modelling error in raster-based spatial data. In *Sbor. ref. konference 3rd International Symposium on Spatial Data Handling, 1998*, s. 97-106.
- Goodchild et al. (1994a)  
 Goodchild, M., Buttenfield, B., Wood, J.: Introduction to visualizing data validity. In *Hearnshaw H., Unwin D. (ed.): Visualization in Geographical Information Systems. UK, John Wiley & Sons Ltd., 1994*, s. 141-149, ISBN 0-471-94435-1.
- Goodchild et al. (1994b)  
 Goodchild, M., Chih-Chang, L., Leung, Y.: Visualizing fuzzy maps. In *Hearnshaw H., Unwin D. (ed.): Visualization in Geographical Information Systems. UK, John Wiley & Sons Ltd., 1994*, s. 158-167, ISBN 0-471-94435-1.
- Goodchild, Janelle (2004)  
 Goodchild M., Janelle D.: Spatially integrated social science. *Oxford University Press, 2004*, 456 stran, ISBN 0-19-515270-0.
- Hamilton (1987)  
 Hamilton, F.: Industrial organisation and regional labour markets. In *Fischer, M., Nijkamp, P.(ed.): Regional Labour Markets. Elsevier Science Publishers, 1987*, s. 289-312, ISBN 0-444-70323-3.
- Hanousek, Charamza (1992)  
 Hanousek, J., Charamza, P.: Moderní metody zpracování dat - matematická statistika pro každého. *Praha, Grada, 1992*, 210 s., ISBN 80-85623-31-5.
- Hansen (1996)  
 Hansen, S.: Comparing the Accessibility for patterns for public versus private transport networks. In *Sbor. ref. konference JEC 1996, Volume 1, 1996*, s. 703-707, ISBN 90-5199-268-8.
- Hansen (1997)  
 Hansen, S.: Integration of Incompatible Zonal Systems using Multiple Overlay Techniques. In *Sbor. ref. konference JEC 1997, Volume 2, 1997*, s. 266-270, ISBN 90-5199-331-5.
- Hanuš (1992)  
 Hanuš, F.: Systémová a operační analýza. Vybrané modely a metody řešení na osobních počítačích. *Skriptum, Praha, ČVUT, fakulta strojní, 1992*, 196 s., ISBN 80-01-00760-X.
- Havlík (1981)  
 Havlík, V.: Geografie v územním plánování. *Skriptum, Praha, ČVUT, fakulta architektury, 1981*, 83 s.
- Hearnshaw, Unwin (1994)  
 Hearnshaw, H., Unwin D.: Visualization in Geographical Information Systems. *UK, John Wiley & Sons Ltd., 1994*, 243 s., ISBN 0-471-94435-1.
- Helísek (2000)  
 Helísek, M.: Makroekonomie. Základní kurs. *Slaný, Melandrium, 2000*, 320 s., ISBN 80-86175-10-3.
- Heyne (1991)  
 Heyne, P.: Ekonomický styl myšlení. *Praha, VŠE, 1991*, 509 s., ISBN 80-7079-781-9.
- Horák (2002a)  
 Horák, J.: Implementace geoinformačních nástrojů v činnosti úřadů práce. [CD-ROM] In *Sbor. ref. konference GIS Ostrava 2002, Ostrava, 2002, ISSN 1213-2454. Dostupné na WWW: <[http://gis.vsb.cz/Publikace/Sborniky/GIS\\_Ova/GIS\\_Ova\\_2002/Sbornik/Referaty/horak2.htm](http://gis.vsb.cz/Publikace/Sborniky/GIS_Ova/GIS_Ova_2002/Sbornik/Referaty/horak2.htm)>*
- Horák (2002b)  
 Horák, J.: Využití pravděpodobnostního mapování v analýze trhu práce. *Ostrava, Sborník vědeckých prací VŠB-TU Ostrava, řada HGF, ročník 48, 2002, 1, s. 131-139, v tisku.*
- Horová, Zelinka (2000)



- Horová, I., Zelinka, J.: Základy a aplikace jádrových odhadů. In *Kupka, K. (ed.): Analýza dat. Pardubice, Trilobyte Ltd., 2000, s. 141-167, ISBN 80-238-6590-0.*
- Hůrský (1969)  
Hůrský, J.: Metody grafického znázornění dojížděky do práce. *Praha, Rozpravy československé akademie věd, řada matematických a přírodních věd, ročník 79, 1969, sešit 3, Academia, 85 s.*
- Ivanička (1980a)  
Ivanička, K.: Základy teórie a metodológie socioekonomickej geografie. *Bratislava, Slovenské pedagogické nakladateľstvo, 1980, 448 s.*
- Ivanička (1980b)  
Ivanička, K.: Prognóza ekonomicko-geografických systémov. *Bratislava, Alfa, 1980, 275 s.*
- Ivanička (1983)  
Ivanička, K.: Základy teórie a metodológie socioekonomickej geografie. *Bratislava, SPN, 1983, 448 s.*
- Jánošíková, Kubáni (2000)  
Jánošíková, L., Kubáni, A.: Dopravná dostupnosť obcí. *Komunikácie, 2000, 4, Vydavateľstvi Žilinské univerzity, 7 s., ISSN 1335-4205.*
- Johansson, Kalsson (1987)  
Johansson, B., Kalsson, Ch.: Processes of industrial change: scale, location and type of job. In *Fischer, M., Nijkamp, P. (ed.): Regional Labour Markets. Elsevier Science Publishers, 1987, s.139-165, ISBN 0-444-70323-3.*
- Johnston et al.(1994)  
Johnston, R., Gregory, D., Smith D.: The dictionary of human geography. *Cambridge, Blackwell, 1994.*
- Jong, Eck (1997)  
Jong, T., van Eck R.: Threshold surfaces as an alternative for locating service centres with GIS. In *Sbor. ref. konference JEC 1997, Volume 2, 1997, s. 799-808, ISBN 90-5199-331-5.*
- Jordan (1995)  
Jordan, P.: Functional regions in East-Central Europe defined on the basis of the frequency of public bus traffic. *Geografický časopis, ročník 47, 1995, 1, s. 9-15.*
- Journal (1991)  
Journal, A.G.: Fundamentals of Geostatistics in five lessons. American Geophysical Union 1989 - kat. ložiskové a průzkumné geologie HGF VŠB Ostrava, 1991. 38 stran.
- Journal, Huijbregts  
Journal, A.G, Huijbregts: Mining geostatistics. Academic Press, 600 stran.
- Jurečka et al. (1999)  
Jurečka, V. a kolektiv: Základy ekonomie. *Skriptum, VŠB-TU Ostrava, Ekonomická fakulta, Ostrava, 1999, 260 s., ISBN 80-7078-660-4.*
- Jurečka, Březinová (1999)  
Jurečka, V., Březinová O.: Mikroekonomie. *Skriptum, Ostrava, VŠB-TU Ostrava, Ekonomická fakulta, 1999, 300 s., ISBN 80-7078-472-5.*
- Kala (2000)  
Kala, F.: Nerovnoměrnost mezd - Česká republika. Úroveň a struktura mzdové diferenciace zaměstnanců v roce 2000 a dynamika za uplynulé tříleté období. *Zpráva o výsledcích projektu PHARE ACE „Determinanty individuálních výdělků a podnikových mzdových struktur v ČR a SR“, Trexima s.r.o., Zlín, 2001.*
- Kaňok (1992)  
Kaňok, J.: Kvantitativní metody v geografii - 1.díl. Grafické a kartografické metody. *Skriptum, Ostrava, Ostravská univerzita, fakulta přírodovědecká, 1992, 236 s., ISBN 80-7042-700-0.*
- Kaňok (1996)

- Kaňok M.: Statistické metody v řízení. *Skriptum, Praha, ČVUT, 1996, 210 s.*
- Kaňok (1998 )  
Kaňok, J.: Tvorba stupnic pro kartogramy a kartodiagramy. *Folia geographica, Prešov, 1998, 2, s. 224-227.*
- Kaňok (1999 )  
Kaňok, J.: Tématická kartografie. *Skriptum, Ostrava, Ostravská univerzita, fakulta přírodovědecká, 1999, 318 s., ISBN 80-7042-781-7.*
- Kerr (1950)  
Kerr, C.: Labor Markets: Their Character and Consequences. *American Economic Review, 40, 1950, s. 278-291.*
- Kouba K. (2007)  
Prostorová analýza českého stranického systému. Institucionalizace a prostorové režimy. *Sociologický časopis, 2007, Vol. 43, No. 5, pp. 1017-1037.*
- Kraak, Ormelling (1996)  
Kraak, M., Ormelling, F.: Cartography. Visualisation of spatial data. *Londýn, Longman Scientific & Technical, 1996, 222 s., ISBN 0-582-25953-3.*
- Kusendová (1996a)  
Kusendová, D.: Netradičné formy kartografických modelov a ich použitie v geografii. *Bratislava, Kartografické listy, 1996, 4, s. 89-100.*
- Kusendová (1996b)  
Kusendová, D.: Analýza dostupnosti obcí Slovenska. In *Sbor. ref. konference Aktivity v kartografii '96, Kartografická spoločnosť SR a Geografický ústav SAV, Bratislava, s. 29-49.*
- Kafka (1996)  
Kafka, Š.: Využití GIS při analýze dopravní obslužnosti na okrese Kutná Hora. *Závěrečná práce PGS GIS, Ostrava, VŠB-TU Ostrava, 1996, 20 s.*
- Laurini, Thompson (1994)  
Laurini, R., Thompson, D.: Fundamentals of Spatial Information Systems. *Londýn, The APIC series, number 37, Academic Press, 1994, 680 s, ISBN 0-12-438380-7.*
- Leakha, Brett (2000)  
Leakha, M., Brett, A.: Visualising the spatio-temporal patterns of motor vehicle theft in Adelaide, South Australia. In *Sbor. ref. konference Crime Mapping: Adding Value to Crime Prevention and Control, 21-22.září 2000, Australian Mineral Foundation, Adelaide, 20 s. Dostupné na WWW: <http://www.aic.gov.au/conferences/mapping/carcach.pdf>.*
- Lee, Wong (2001)  
Lee, J., Wong, D.: Statistical analysis with ArcView GIS. *USA, John Wiley & Sons Inc., 2001, 192 s., ISBN 0-471-34874-0.*
- Lever (1987)  
Lever, W.: New trends in the supply and demand patterns of labour in western economies. In *Fischer, M., Nijkamp, P. (ed.): Regional Labour Markets. Elsevier Science Publishers, 1987, s.249-267, ISBN 0-444-70323-3.*
- Levine (2007)  
**Levine ...**
- LeSage (1998)  
LeSage J. P.: Spatial Econometrics. 1998. 273 stran
- Líčeník (1985)  
Líčeník, J.: Geografie v územním plánování. *Skriptum, Brno, VUT Brno, fakulta architektury, 1985, 130 s.*
- Longley, Clarke (1995)

- Longley, P., Clarke, G.: GIS for Business and Service Planning. *Cambridge, Geoinformation International*, 1995, 316 s., ISBN 1-899761-07-1.
- Lukasová, Šarmanová (1985)  
Lukasová, A., Šarmanová, J.: Metody shlukové analýzy. *Praha, SNTL*, 1985, 210 s.
- Lukoszová, Grasseová (2000)  
Lukoszová, X., Grasseová, M.: Jednoduché kvantitativní metody aplikovatelné obchodními organizacemi. *Ostrava, Ekonomická revue*, 2000, 2, VŠB-TU Ostrava, s. 40-47, ISSN 1212-3951.
- Lopez (1996)  
Lopez, A.: Mortality data. In Elliott P., Cuzick J., English D., Stern R.(ed.): *Geographical & Environmental Epidemiology*. New York, Oxford University Press, 1996, s.37-50, ISBN 0-19-262235-8.
- Martin (1991)  
Martin, D.: Geographic information systems and their socioeconomic applications. *New York, Routledge*, 1991, 182 s., ISBN 0-415-05698-5.
- Martin (1995)  
Martin, D.: Censuses and the modelling of population in GIS. In Longley, P., Clarke, G. (ed.): *GIS for Business and Service Planning*. Cambridge, Geoinformation International, 1995, s.48-71, ISBN 1-899761-07-1.
- Mareš (2002)  
Mareš, P.: Nezaměstnanost jako sociální problém. *Praha, Sociologické nakladatelství, řada Studijní texty*, 2002, 172 s., ISBN 80-86429-08-3.
- Maryáš et al. (1995)  
Maryáš, J., Řehák, S., Vystoupil, J., Mládek, J.: Ekonomická geografie I. *Skriptum, Brno, Masarykova univerzita, Ekonomicko-správní fakulta*, 1995, 138 s., ISBN 80-210-1269-2.
- Meloun, Militký (2002)  
Meloun, M., Militký, J.: Kompendium statistického zpracování dat. *Praha, Academia*, 2002, 764 s., ISBN 80-200-1008-4.
- Menegelo, Peckham (1996)  
Menegelo, L., Peckham, R.: A fully integrated tool for site planning using multicriteria evaluation techniques within a GIS. In *Sbor. ref. konference JEC, Volume 1*, 1996, s. 621-630, ISBN 90-5199-268-8.
- Murdych (1969)  
Murdych, Z.: Kartometrická analýza centrality krajských měst Československa. *Geografický časopis, ročník 21, 1969, 4, s. 277-286*.
- Murdych (1973)  
Murdych, Z.: Možnosti přehledného kartografického znázorňování dopravní dostupnosti na příkladech slovenských krajů. *Geografický časopis, ročník 25, 1973, 1, s. 47-53*.
- Molnár (1992)  
Molnár, Z.: Moderní metody řízení informačních systémů. *Praha, Grada*, 1992, 352 s., ISBN 80-85623-07-2.
- Nijkamp, Rietveld (1987)  
Nijkamp, P., Rietveld, P.: Technological development and regional labour markets. In Fischer, M., Nijkamp, P. (ed.): *Regional Labour Markets*. Elsevier Science Publishers, 1987, s. 117-138, ISBN 0-444-70323-3.
- Novotná (2001)  
Novotná, M.: Vimpersko. Geografická analýza příhraničního mikroregionu. *Plzeň, ZČU*, 2001, 121 s.
- Openshaw (1995)

- Openshaw, S.: Marketing spatial analysis: a review of prospects and technologies relevant to marketing. In Longley, P., Clarke, G. (ed.): *GIS for Business and Service Planning*. Cambridge, Geoinformation International, 1995, s. 150-165, ISBN 1-899761-07-1.
- Quinn (1996)  
 Quinn, M.: Confidentiality. In Elliott P., Cuzick J., English D., Stern R.(ed.): *Geographical & Environmental Epidemiology*. New York, Oxford University Press, 1996, s. 132-140, ISBN 0-19-262235-8.
- Pavlík, Kühnl (1981)  
 Pavlík, Z., Kühnl, K.: Úvod do kvantitativních metod pro geografii. Praha, SPN, 1981, 267 s.
- Pebesma (1999)  
 Pebesma E.: Manuál ke GSTAT 2.1.2. 1992-1999. Dostupný na WWW <<http://www.geog.uu.nl/gstat>>
- Peňáz et al. (2000)  
 Peňáz, T., Horák, J., Horáková, B.: Analýza územní dostupnosti významných firem na území okresu Nový Jičín. In Sbor. ref. konference GIS Seč 2000, Seč 7.-9.6.2000, 2000, s. 280-289, ISSN 1211-7439, ISBN 80-86143-17-1.
- Píšek, Hanuš (1996)  
 Píšek, M., Hanuš, F.: Rozhodovací analýza. Vybrané modely a metody řešení na PC. Skriptum, Praha, ČVUT, fakulta strojní, 1996, 78 s., ISBN 80-01-01534-3.
- Richardson (1996)  
 Richardson, S.: Statistical methods for geographical correlation studies. In Elliott P., Cuzick J., English D., Stern R.(ed.): *Geographical & Environmental Epidemiology*. New York, Oxford University Press, 1996, s. 181-204, ISBN 0-19-262235-8.
- Robinson et al. (1995)  
 Robinson, A., Morrison, J., Muehrcke, P., Kimerling, A., Guptill, S.: Elements of cartography. New York, John Wiley & Sons Ltd., 6.vydání, 1995, 674 s., ISBN 0-471-55579-7.
- Rogalewitz (1993)  
 Rogalewitz, V.: Stochastické procesy (analýza časových řad). Skriptum, Praha, ČVUT, fakulta elektrotechnická, 1993, 106 s., ISBN 80-01-00905-X.
- Rochovská, Horňák (2002)  
 Rochovská, A., Horňák, M.: Sociálna polarizácia spoločnosti a jej regionálny priemet na území Slovenska a Bratislavy. Plzeň, *Miscellanea geographica* 9, 2002, s. 131-142, ISBN 80-7082-805-6.
- Rouwendal, Nijkamp (1987)  
 Rouwendal J., Nijkamp P.: Regional economic research on labour markets. In Fischer, M., Nijkamp, P. (ed.): *Regional Labour Markets*. Elsevier Science Publishers, 1987, s. 95-115, ISBN 0-444-70323-3.
- Rölc (2001)  
 Rölc, R.: Dopravní dostupnost a regionální význam krajských měst. *Geografie - sborník České geografické společnosti, ročník 106, 2001, 4, s. 222-233.*
- Řehák (1992)  
 Řehák, S.: Sídelně dopravní model ČSFR a jeho územní souvislosti. *Geografický časopis, ročník 44, 1992, 1, s. 59-72.*
- Saaty (1977)  
 Saaty, T.: A scaling Method for Priorities in Hierarchical Structures. *Journal of Math. Psychology, 15, 1977, s.234-281.*
- Sabel et al. (2005)  
 Sabel C.E., Kingham S., Nicholson A., Bartie P.: Road Traffic Accident Simulation Modelling - A Kernel Estimation Approach. In proc. Of The 17th Annual Colloquium of the Spatial

- Senjuk (2001)  
Senjuk, I.: Základy dopravního inženýrství. Logistika a marketing. *Skriptum, Praha, ČVUT, fakulta dopravní, 2001, 192 s., ISBN 80-01-02338-9.*
- Schubert et al.(1987)  
Schubert U., Gerking S., Isserman A., Taylor C.: Regional labour market moelling: A state of the art review. In *Fischer, M., Nijkamp, P. (ed.): Regional Labour Markets. Elsevier Science Publishers, 1987, s.53-94, ISBN 0-444-70323-3.*
- Smans, Esteve (1996)  
Smans, M., Esteve, J.: Practical approaches to disease mapping. In *Elliott P., Cuzick J., English D., Stern R.(ed.): Geographical & Environmental Epidemiology. New York, Oxford University Press, 1996, s. 141-150, ISBN 0-19-262235-8.*
- Spurná P. (2006)  
Spurná P.: Současné trendy v kvantitativní analýze geografických dat se zaměřením na využití metody geograficky vážené regrese DP Katedra sociální geografie a regionálního rozvoje PřF UK, Praha, 150 s.
- Spurná P. (2008)  
Spurná P.: Geograficky vážená regrese: metoda analýzy prostorové nestacionarity geografických jevů. *Geografie- Sborník ČGS, 2008, 113, 2, 125-139.*
- Székely (1999)  
Székely, V.: Časovo-priestorová diferenciácia nezamestnanosti a jej tokov na Slovensku v rokoch 1997-1999. *Geografický časopis, ročník 53, 2001, 2, s. 147-170.*
- Šotkovský (1996)  
Šotkovský, I.: Úvod do studia demografie. *Skriptum, Ostrava, VŠB, fakulta ekonomická, 1996, 159 s., ISBN 80-7078-327-3.*
- Tiefelsdorf (2000)  
Tiefelsdorf, M.: „Modelling spatial processes: the identification and analysis of spatial relationships in regression residuals by means of Moran’s. *Berlín, Springer-Verlag, řada Lecture notes in earth sciences; 87, 2000, 167 s., ISBN 3-540-66208-1.*
- Thill, Horowitz (1997)  
Thill, J., Horowitz, J: Modelling non-work destination choices with choice sets defined by travel-time constraints. In *Fischer M., Getis A. (ed.): Recent Developments in Spatial Analysis. Spatial Statistics, Behavioural Modelling and Computational Intelligence. Berlín, Springer-Verlag, 1997, s. 186-208, ISBN 3-540-63180-1.*
- Tuček (1998):  
Tuček, J.: Geografické informační systémy. Principy a praxe. *Praha, ComputerPress, 1998, 424 s., ISBN 80-7226-091-X.*
- Turk (1994)  
Turk, A.: Cogent GIS visualizations. In *Hearnshaw H., Unwin D. (ed.): Visualization in Geographical Information Systems. UK, John Wiley & Sons Ltd., 1994, s.26-33, ISBN 0-471-94435-1.*
- Turčan (2002)  
Turčan, M.: Statistické metody. *Skriptum pro PGS, Ostrava, 2002, 166 s.*
- You-Hong (1996)  
You-Hong, Ch.: Exploring Spatial Analysis in Geographic Information Systems. *Santa Fe, Onword Press, 1996, 473 s., ISBN 1-56690-118-9.*
- Unwin (1981)  
Unwin, D.: Introductory Spatial Analysis. *Londýn, Methuen, 1981.*
- Vencálek (1991)

- Vencálek, J.: Socioekonomická geografie I. *Ostrava, Pedagogická fakulta, 1991, 191 s., ISBN 80-7042-056-1.*
- Veverka (1997)  
Veverka, B.: Topografická a tematická kartografie. *Skriptum, Praha, ČVUT, fakulta stavební, 1997, 203 s., ISBN 80-01-0124-X.*
- Veverka (1989)  
Veverka, B.: Teorie systémů a kybernetika. *Skriptum, Praha, ČVUT, fakulta stavební, 1989. 154 s.*
- Virus (1999)  
Virus, M.: Aplikace matematické statistiky. Metoda Monte Carlo. *Skriptum, Praha, ČVUT, fakulta jaderná a fyzikálně inženýrská, 1999, 168 s.*
- Víšek (1998)  
Víšek, J.: Statistická analýza dat. *Skriptum, Praha, ČVUT, fakulta jaderná a fyzikálně inženýrská, 1998, 187 s.*
- Vojta (2009):  
Vojta M. Analýza dojíždění pro podniky Lanex a MSA. DP, Institut geoinformatiky, HGF VŠB-TU Ostrava, 2009.
- Voženílek (2001)  
Voženílek, V.: Aplikovaná kartografie I. Tematické mapy. *Olomouc, Univerzita Palackého, přírodovědecká fakulta, 2001, 1987 s., ISBN 80-244-0270-X.*
- W3C (2001)  
W3C: Extensible Markup Language (XML) [on-line]. 2001. Dostupný na WWW: <<http://www.w3.org/TR/REC-xml>>
- Walford (1995)  
Walford, N.: Geographical Data Analysis. *UK, John Wiley & Sons Ltd., 1995, 446 s., ISBN 0-471-94162-X.*
- Weis (2002)  
Weis, V.: Registry veřejné správy, zákon o registrech veřejné správy. [CD-ROM] In *Sbor. ref. konference 3.konference Městské informační systémy, Praha, 2002.*