

**SPACE - TIME MIXTURE
MODELING
FOR DISEASE MAPPING**

*Dankmar Böhning¹
Peter Schlattmann and Heinz Holling*

**Department of Epidemiology
Institute of Social Medicine and Medical
Psychology, Free University Berlin**

Invited Presentation at the 98 IBC in Cape Town

Session on

Space-Time Modeling of Epidemiological Data

Organized by Andrew Lawson

OVERVIEW

- Introduction
- Disease Mapping in Space Using Mixtures
- Disease Mapping in Space *and Time*

FREQUENT OBJECTIVE IN GEOGRAPHIC EPIDEMIOLOGY

to present that part of the **spatial** variation of a disease occurrence distribution, which cannot be explained by the different distribution of **known factors** in the various regions nor is due to **random** variation

HOPE: hints to **unknown** risk factors!

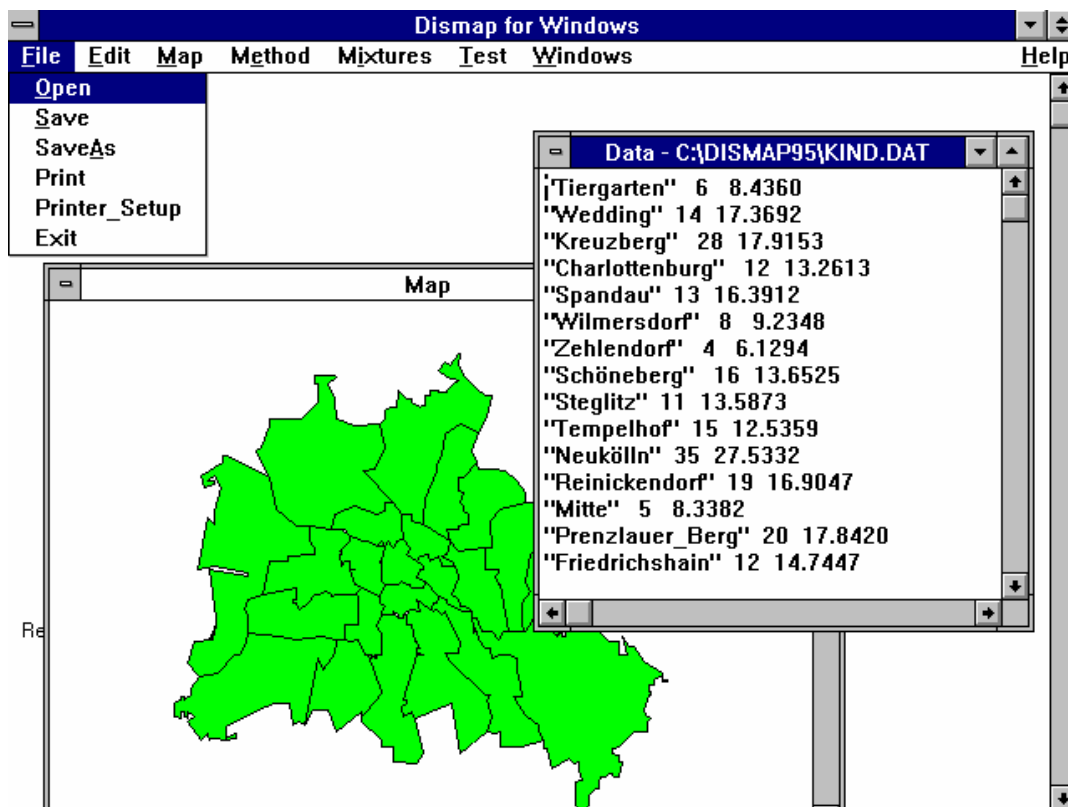
INTRODUCTION

frequently used OCCURRENCE MEASURE in epidemiology and public health institutions

$$SMR_i = O_i/E_i \quad \text{in the } i\text{-th region}$$

O_i observed death (mortality) or disease (morbidity) cases

E_i expected cases computed from external reference population



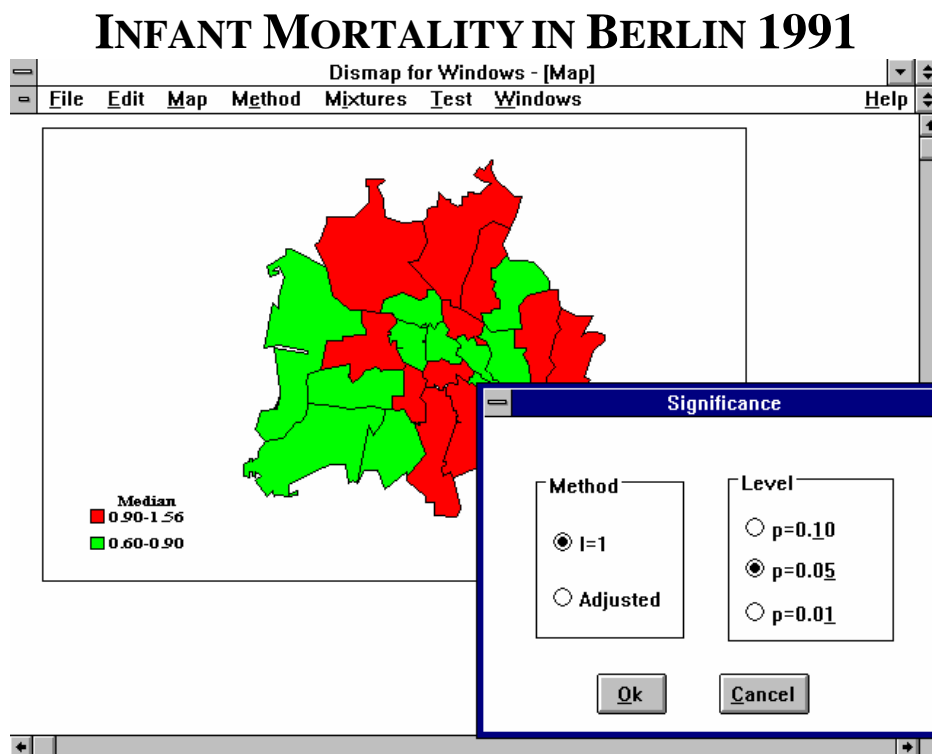
TWO „CLASSIC” BIOMETRIC METHODS

- a) classification based on a certain percentile of the empirical SMR-distribution
- b) Poisson distribution $Po(o_i, \lambda E_i)$

$$= \exp(-\lambda E_i) (\lambda E_i)^{o_i} / o_i!$$

classification based on the P-value under the Poisson distribution

$$P(O_i \geq o_i) = Po(o_i, \lambda E_i) + Po(o_i + 1, \lambda E_i) + \dots$$



Disadvantage of conventional methods:

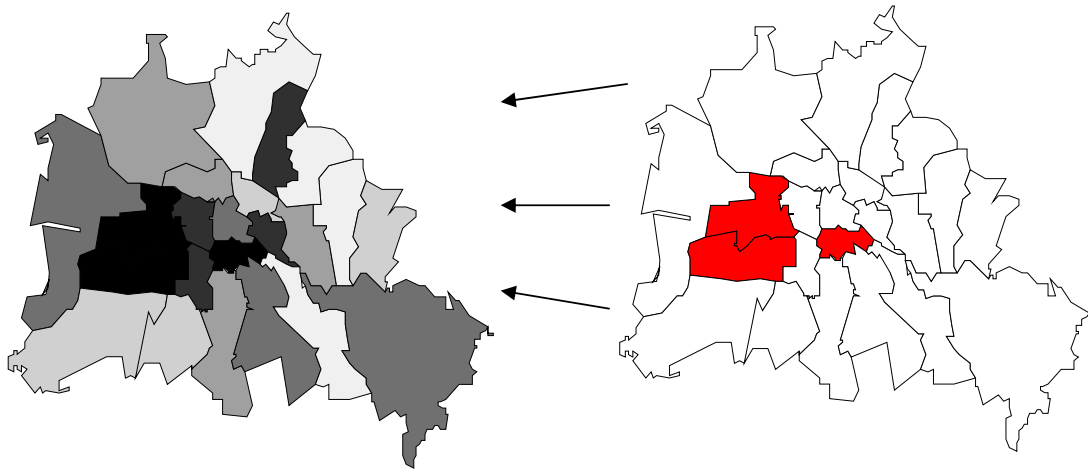
represent random variation on the map

Disease Mapping in Space using Mixtures

observed map of
risk structure

true, but unobserved
map of risk structure
(latent or hidden map)

23 city areas of Berlin



n areas given

k (unknown)
components: $\lambda_1, \dots, \lambda_k$

O_1, O_2, \dots, O_n
 E_1, E_2, \dots, E_n } observed data
 Z_1, Z_2, \dots, Z_n unobserved data

$Z_i = (Z_{i1}, \dots, Z_{ik})$ with $Z_{ij} = 1$ meaning:
area i is from component with risk λ_j

MODEL FOR: $SMR_i = O_i/E_i$ is

$$O_i \sim \text{Po}(\lambda_j E_i)$$

conditionally O_i is from area with mortality rate λ_j

let p_j probability of being from component with $\lambda = \lambda_j$

then, **unconditionally**

$$O_i \sim p_1 \text{Po}(\lambda_1 E_i) + p_2 \text{Po}(\lambda_2 E_i) + \dots + p_k \text{Po}(\lambda_k E_i)$$

is a **nonparametric mixture of Poisson** distributions

INTERPRETATION

$k=1$: *homogeneous* risk structure

$k=2$: two risk groups

$k=3$: three risk groups

...

} *heterogeneous*
risk structure

maximum likelihood estimation of the parameters $p_1, \lambda_1, \dots, p_k, \lambda_k$ inclusively number of components k ,

$$P = \begin{pmatrix} \lambda_1 & \lambda_2 & \dots & \lambda_k \\ p_1 & p_2 & \dots & p_k \end{pmatrix}$$

leads to the *nonparametric maximum likelihood estimate* of a mixing distribution

(Laird 1978, Simar 1974, Lindsay 1983, ... Lesperance and Kalbfleisch 1992, Aitkin 1996, ...)

theoretical well-investigated (Lindsay 1995)
algorithmically possible (Böhning 1995 *JSPI*)

Map CONSTRUCTION

each area i is classified into component (color) j such that the *posterior distribution*

$$f(\lambda_j / o_i, E_i, \hat{P}) = Po(o_i, \hat{\lambda}_j E_i) \hat{p}_j / \sum_l Po(o_i, \hat{\lambda}_l E_i) \hat{p}_l$$

is maximized !

TWO APPLICATIONS

Health Region: 219 counties of the former German Democratic Republic (*The 5 New States of Germany*)

- 1) Incidence on Lung Cancer (ICD 162) for Women
- 1980 – 1989

Estimate of Mixing Distribution (NPMLE):

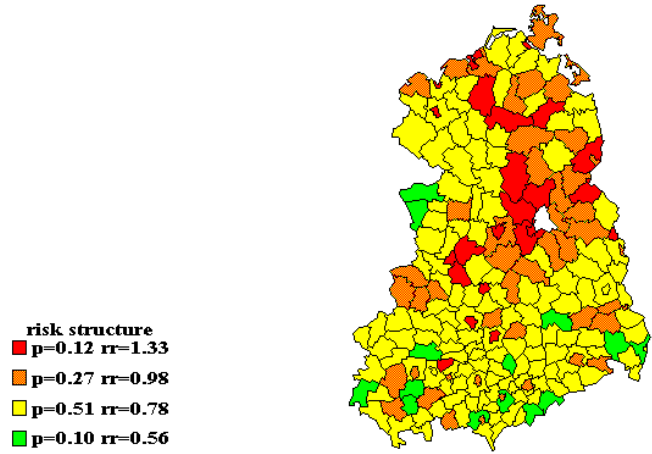
$$\begin{pmatrix} \hat{\lambda}_1 & \hat{\lambda}_2 & \dots & \hat{\lambda}_k \\ \hat{p}_1 & \hat{p}_2 & \dots & \hat{p}_k \end{pmatrix} = \begin{pmatrix} 1.33 & 0.98 & 0.78 & 0.56 \\ 0.12 & 0.27 & 0.51 & 0.10 \end{pmatrix},$$
$$\hat{k} = 4$$

- 2) Incidence on Mamma-Carcinoma (ICD 174)
- 1980 – 1989

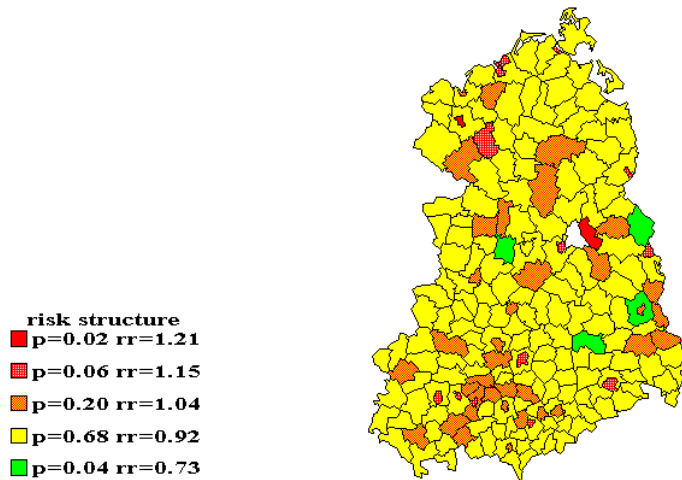
Estimate of Mixing Distribution (NPMLE):

$$\begin{pmatrix} \hat{\lambda}_1 & \hat{\lambda}_2 & \dots & \hat{\lambda}_k \\ \hat{p}_1 & \hat{p}_2 & \dots & \hat{p}_k \end{pmatrix} = \begin{pmatrix} 1.21 & 1.15 & 1.04 & 0.92 & 0.73 \\ 0.02 & 0.06 & 0.20 & 0.68 & 0.04 \end{pmatrix},$$
$$\hat{k} = 5$$

ICD 162 in women GDR 1980-89



ICD 174 in women GDR 1980-89



Disease Mapping in Space *and Time*

n areas given for T periods

$$\begin{aligned} &O_1^{(1)}, O_2^{(1)}, \dots, O_n^{(1)}; O_1^{(2)}, O_2^{(2)}, \dots, O_n^{(2)}; \dots \\ &O_1^{(T)}, O_2^{(T)}, \dots, O_n^{(T)} \\ &E_1^{(1)}, E_2^{(1)}, \dots, E_n^{(1)}; E_1^{(2)}, E_2^{(2)}, \dots, E_n^{(2)}; \dots \\ &E_1^{(T)}, E_2^{(T)}, \dots, E_n^{(T)}; \\ &Z_1^{(1)}, Z_2^{(1)}, \dots, Z_n^{(1)}; \quad Z_1^{(2)}, Z_2^{(2)}, \dots, Z_n^{(2)}; \quad \dots \\ &Z_1^{(T)}, Z_2^{(T)}, \dots, Z_n^{(T)}; \end{aligned}$$

IN THE TWO APPLICATIONS

Incidence on Lung Cancer (ICD 162) for Women

– 1980 – 1989

– 1970 – 1979

– 1960 – 1969

Incidence on Mamma-Carcinoma (ICD 174)

– 1980 – 1989

– 1970 – 1979

– 1960 – 1969

TWO ANALYSIS OPTIONS

a)

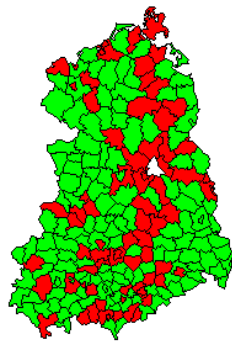
$\mathbf{Z}_i^{(t)} = (Z_{i1}^{(t)}, \dots, Z_{ik}^{(t)})$ with $Z_{ij}^{(t)} = 1$ meaning:

for *each* period: area i is from component with risk $\lambda_j^{(t)}$, $j=1, \dots, k$, where k might depend on t

Result: T mixture models, for each time period one:

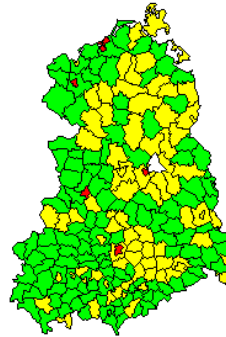
$$O_i^{(t)} \sim p_1^{(t)} \text{Po}(\lambda_1^{(t)} E_i^{(t)}) + p_2^{(t)} \text{Po}(\lambda_2^{(t)} E_i^{(t)}) + \dots + p_k^{(t)} \text{Po}(\lambda_k^{(t)} E_i^{(t)})$$

ICD 162 in women GDR 1960-69



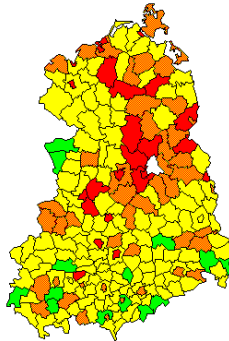
risk structure
■ p=0.38 rr=1.14
■ p=0.62 rr=0.75

ICD 162 in women GDR 1970-79



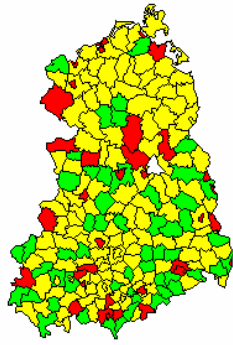
risk structure
■ p=0.06 rr=1.35
■ p=0.39 rr=1.02
■ p=0.55 rr=0.71

ICD 162 in women GDR 1980-89

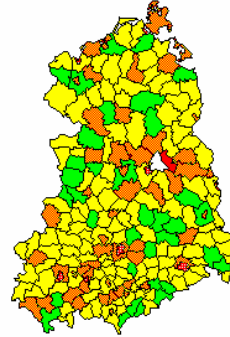


risk structure
■ p=0.12 rr=1.33
■ p=0.27 rr=0.98
■ p=0.51 rr=0.78
■ p=0.10 rr=0.56

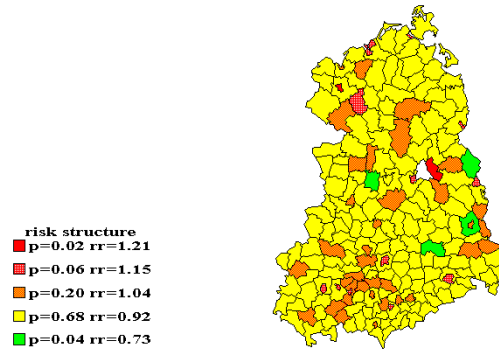
ICD 174 in women GDR 1960-69



ICD 174 in women GDR 1970-79



ICD 174 in women GDR 1980-89



b)

$\mathbf{Z}_i^{(t)} = (Z_{i1}^{(t)}, \dots, Z_{ik}^{(t)})$ with $Z_{ij}^{(t)} = 1$ meaning:

for *all* periods : area i is from component with risk λ_j , $j=1, \dots, k$, where λ_j and k *does not* depend on t

Result: one mixture model

$$O_i^{(t)} \sim p_1 \text{Po}(\lambda_1 E_i^{(t)}) + p_2 \text{Po}(\lambda_2 E_i^{(t)}) + \dots + p_k \text{Po}(\lambda_k E_i^{(t)})$$

Note: in *both* cases T maps are drawn, however in

a) there are $k^{(1)} + \dots + k^{(T)}$ colors

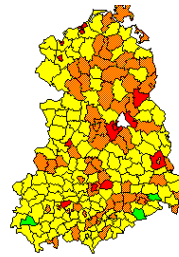
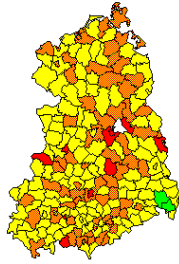
b) there are k colors

INTERPRETATION:

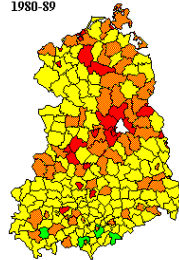
option b) is attractive since it is looking for joint *space-time* components (clusters)

ICD 162 GDR 1960-69

risk structure
■ p=0.10 rr=1.31
■ p=0.32 rr=1.01
■ p=0.52 rr=0.76
■ p=0.06 rr=0.53

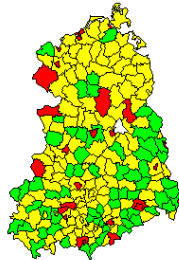


1980-89

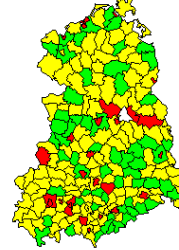


ICD 174 GDR 1960-69

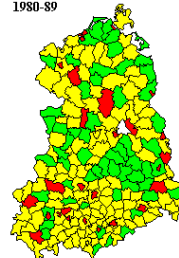
risk structure
■ p=0.16 rr=1.13
■ p=0.54 rr=0.96
■ p=0.30 rr=0.84



ICD 174 GDR 1970-79



1980-89



Option c)

MORE COMPLEX MIXED MODELING

$$O_i^{(t)} \sim p_1 \text{Po}(\lambda_{i1}^{(t)} E_i^{(t)}) + p_2 \text{Po}(\lambda_{i2}^{(t)} E_i^{(t)}) + \dots + p_k \text{Po}(\lambda_{ik} E_i^{(t)})$$

however,

$$\begin{aligned} \log(\lambda_{ij}^{(t)} E_i^{(t)}) &= \log(E_i^{(t)}) + \log(\lambda_{ij}^{(t)}) \\ &= \log(E_i^{(t)}) + \alpha_j + \beta_j t + \text{further covariates} \end{aligned}$$

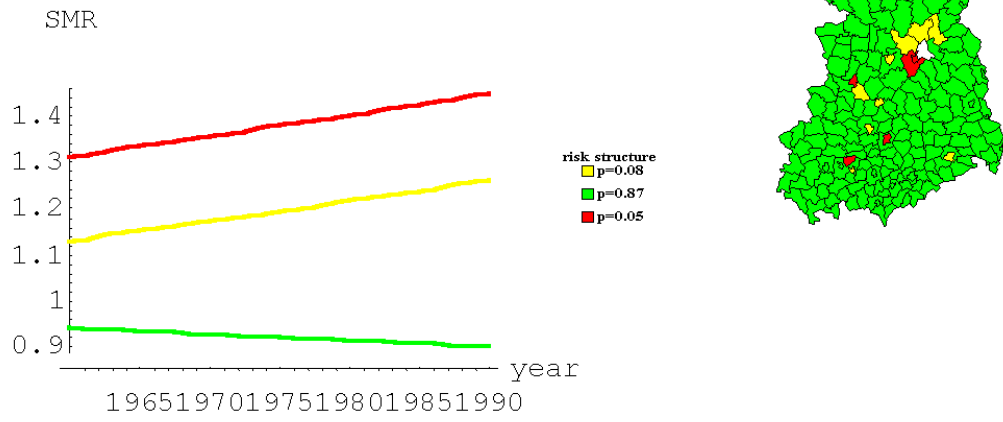
leading to

$$O_i^{(t)} \sim \sum_{j=1}^k p_j \text{Po}\{E_i^{(t)} \exp(\alpha_j + \beta_j t)\}$$

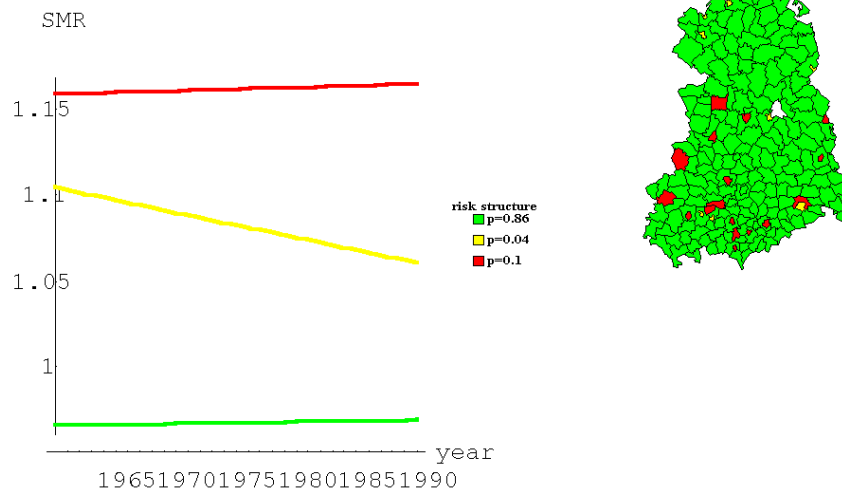
NOTE: in this case the map consists of the
k log-linear models (each model one color)

DISADVANTAGE: increased complexity: mixing
over intercept or effect parameter? or both?

ICD 162 GDR 1960-89



ICD 174 GDR 1960-89



In Conclusion

- ⇒ investigated possibilities of consistently estimating heterogeneity via nonparametric mixture models
- ⇒ availability of these procedures in C.A.MAN and DISMAP

SOME HINTS TO RECENT REFERENCES

Lawson, Böhning, Biggeri, Lesaffre, Viel, and Bertolini (Eds.), 1998. *Disease Mapping and Risk Assessment for Public Health*. Wiley & Sons.

Böhning, 1998. *Computer Assisted Analysis of Mixtures and Applications: Disease Mapping, Meta-Analysis and others*. Chapman & Hall