

Towards multivariate modelling of geogenic radon

(parts 1 & 2 together)

P. Bossew¹, W. Meyer¹, M. Bleher²

Bundesamt für Strahlenschutz / German Federal Office for Radiation Protection,
¹Berlin, ²Munich

10th INTERNATIONAL WORKSHOP
on the
GEOLOGICAL ASPECTS OF RADON RISK MAPPING
22 – 25 Sept 2010, Prague, Czech Republic

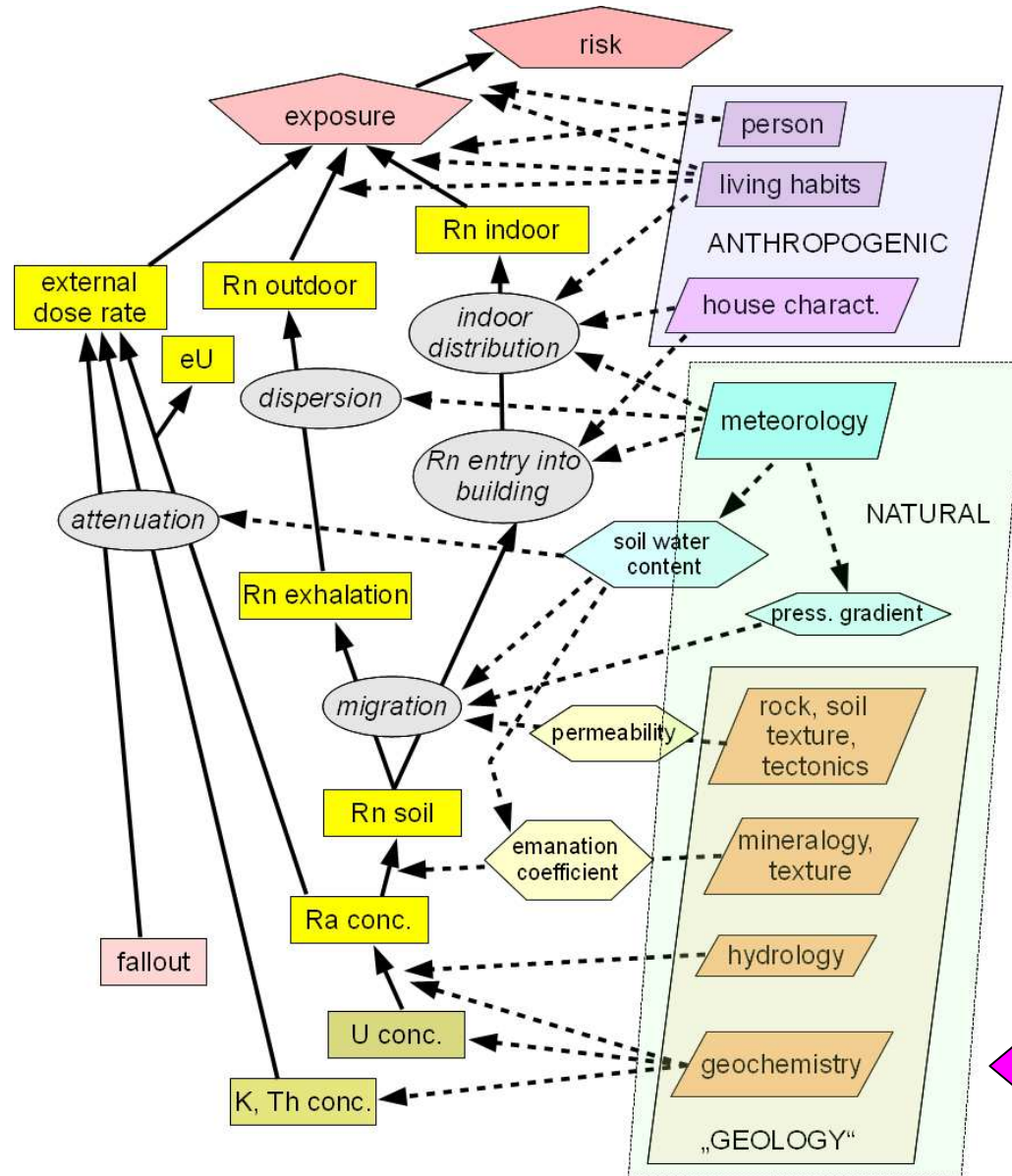
v. 21 Sept 2010

Content

- Reminder: physical complexity
- Concepts
 - types of variables
 - construction of target variable
- Regression modelling
- Correlation between variables
- Geostatistical modelling of residua
- Target variables
- Examples



Physical complexity



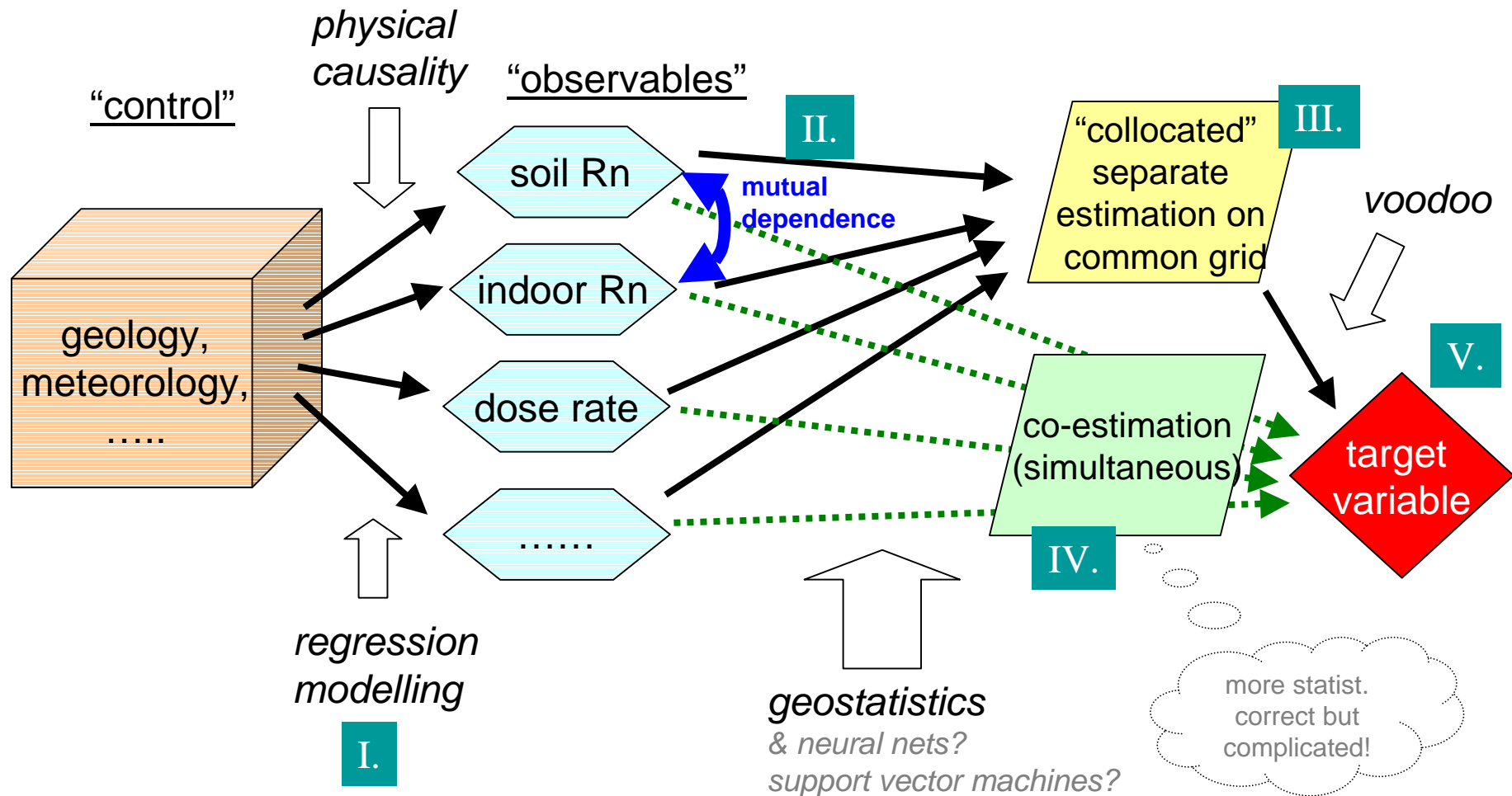
- complex pathway from Rn source – migration – exposure and risk
- some variables not trivial to define accurately
- spatially complicated structure:
 - high nugget
 - hot spots
 - geological predictors not easy to define

- simplified!
 - yellow: some observed Rn variables

Wanted

- spatial prediction of a variable which quantifies Rn hazard
 - its levels
 - its support (“Rn prone areas”)
- hazard:
 - define “hazard” variable Y
 - estimate / predict from *observable quantities*
indoor concentration, conc. in soil air, external dose rate,...
 - and from observable *physical controls*
geology, soil properties, geographic location...
- $E[Y(x)](x \in U)$, $\text{prob}[Y(x) > T](x \in U)$??
 $U =$ prediction support (point, cell, admin. unit,...)

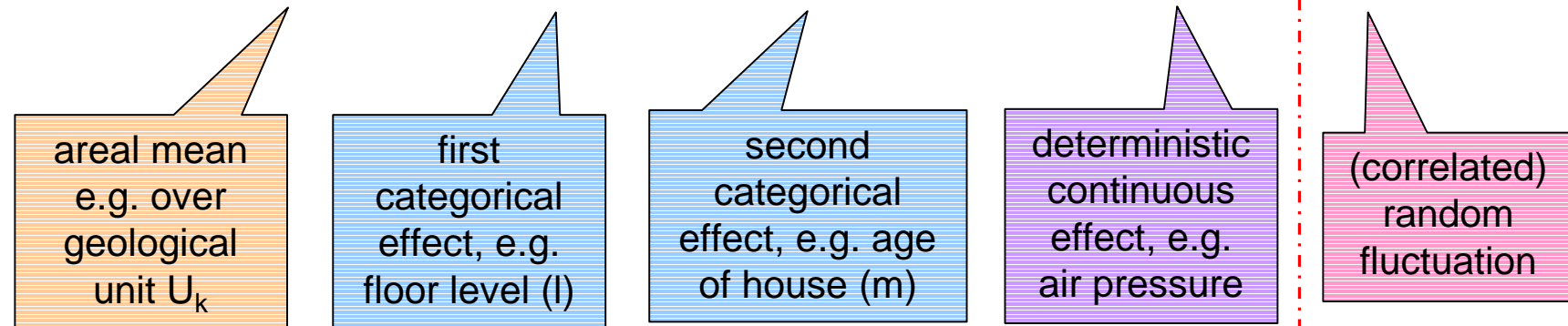
Concept & road map I. to V.



I. regression modelling

General linear model (GLM)

$$Z(\mathbf{U}) = \mu_k(\mathbf{U}) + \alpha_{kl}(\mathbf{U}) + \beta_{klm}(\mathbf{U}) + \dots + f(\mathbf{x}) + \varepsilon(\mathbf{x})$$



factors may or may not be contingent / correlated
ex.: contingent: age of construction ~ building material

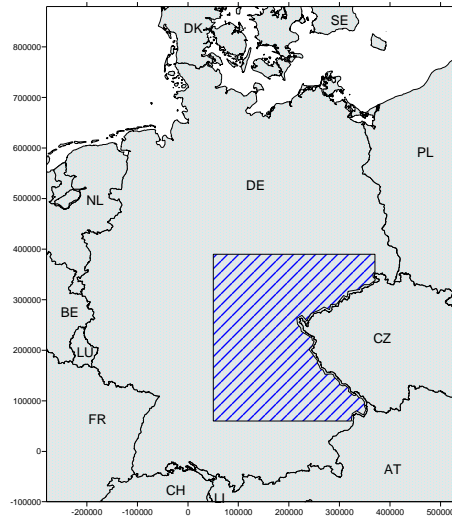
separate estimation: e.g. kriging with external drift
simultaneous estimation: e.g. regression kriging

Z = physical observable, or derived: e.g. log(indoor concentration)
 U = spatial unit, $x \in U$. (e.g. $U = \{x\}$)

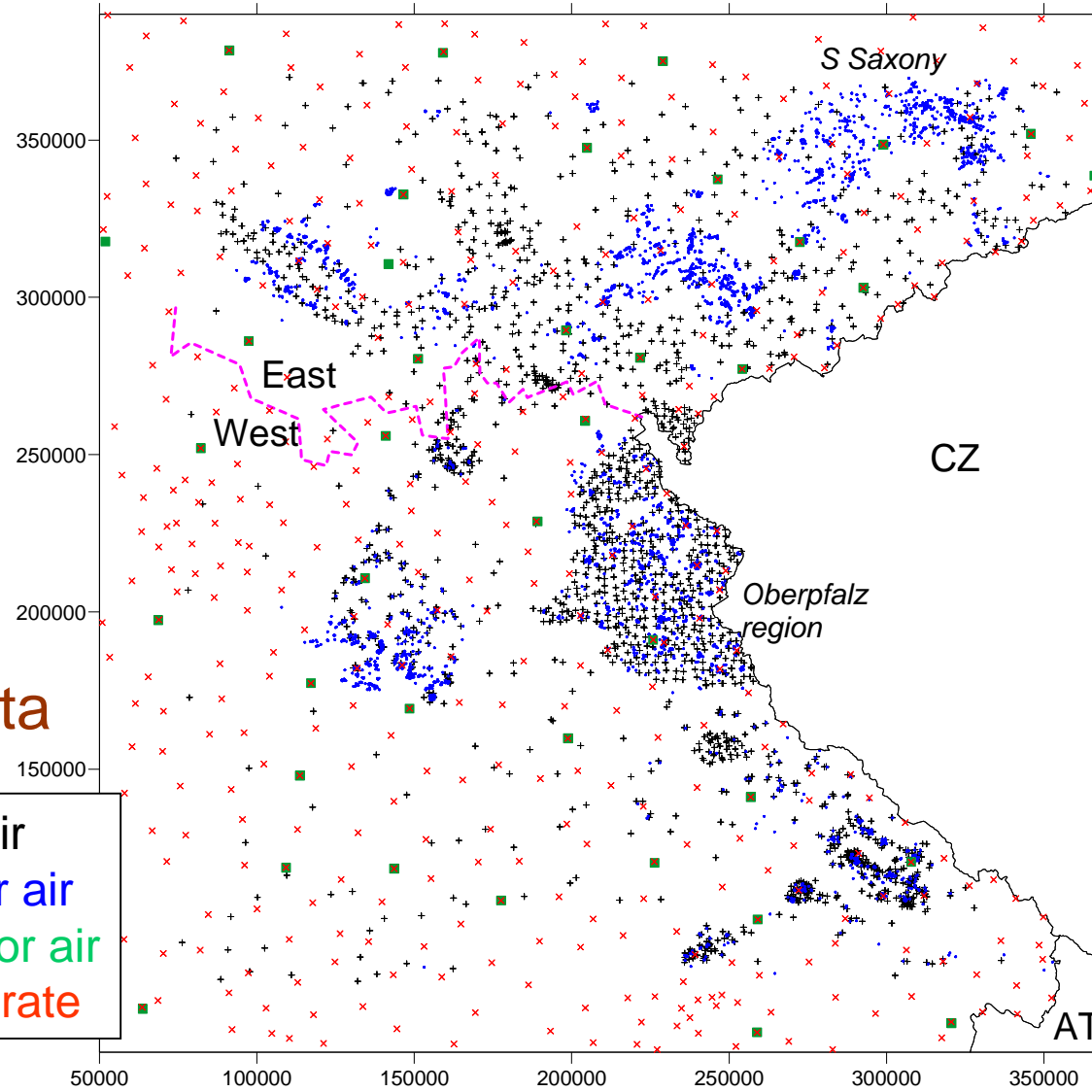
because factors are multiplicative!

division
deterministic //
stochastic part:
non-trivial
question!

example (1/4)



study region

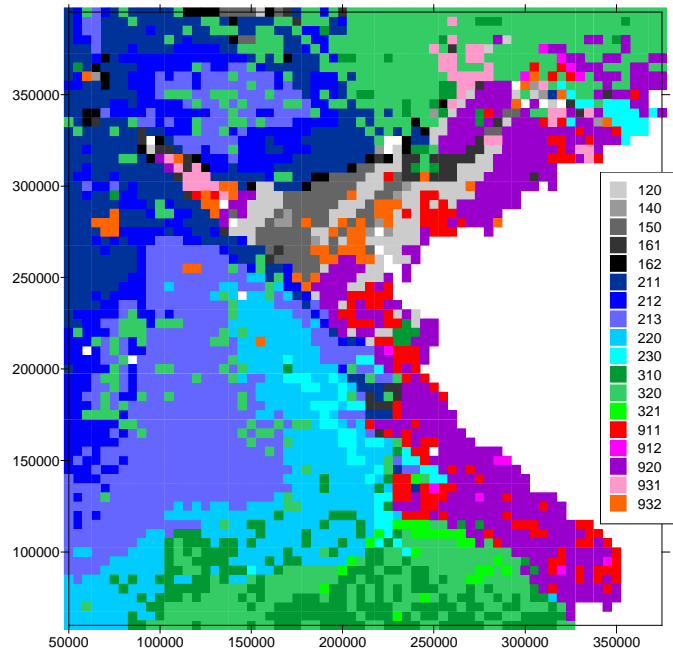


data

- + soil air
- indoor air
- outdoor air
- x dose rate

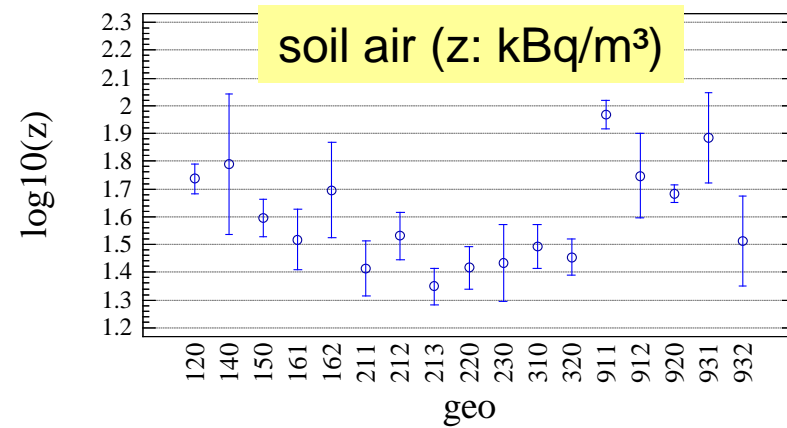
example (2/4)

geo units:

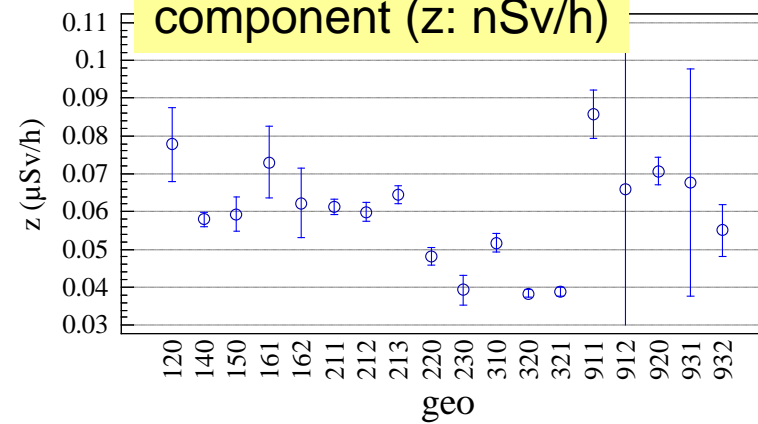


simplified, from German geol. map 1:1M,
coding following Klingel & Kemski

Means and 95.0 Percent Confidence Intervals (internal s)



ext. dose rate, terr. component (z: nSv/h)

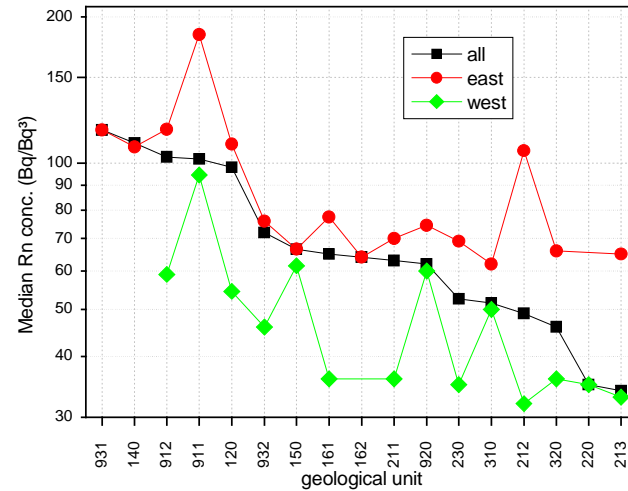


example (3/4)

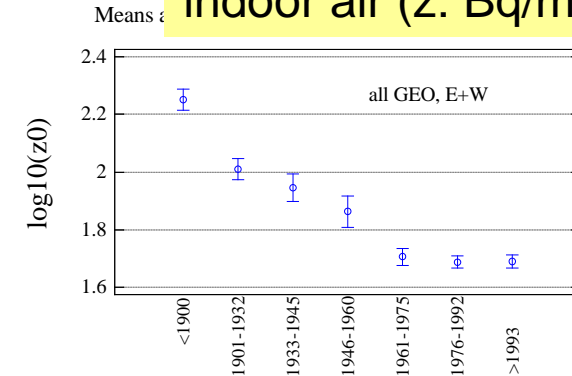
indoor Rn:

dependencies of factors:

- geology
- former East- / West-Germany
- year of construction



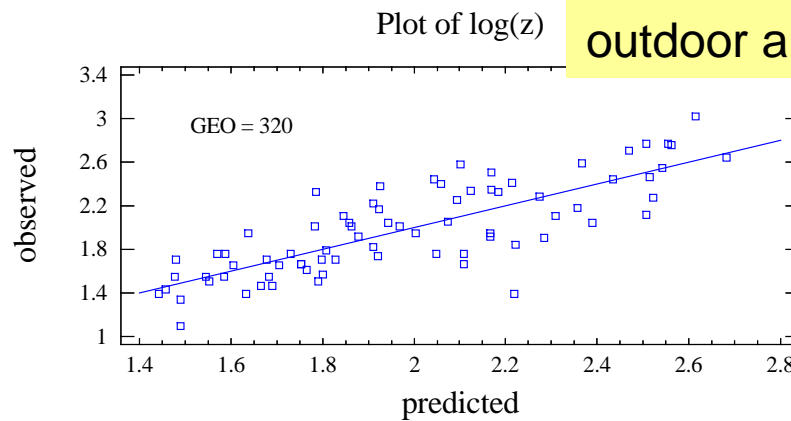
indoor air (z: Bq/m³)



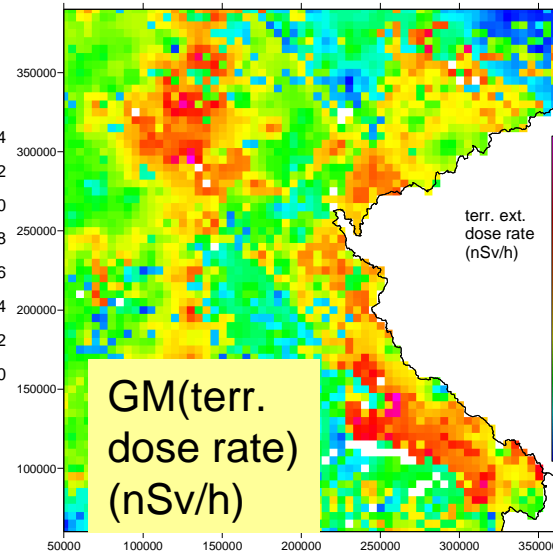
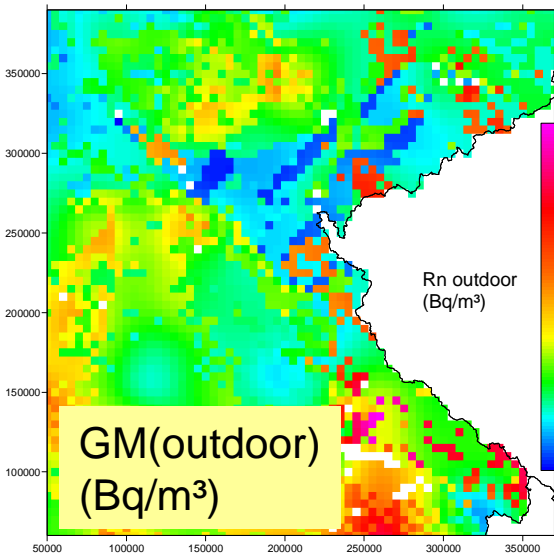
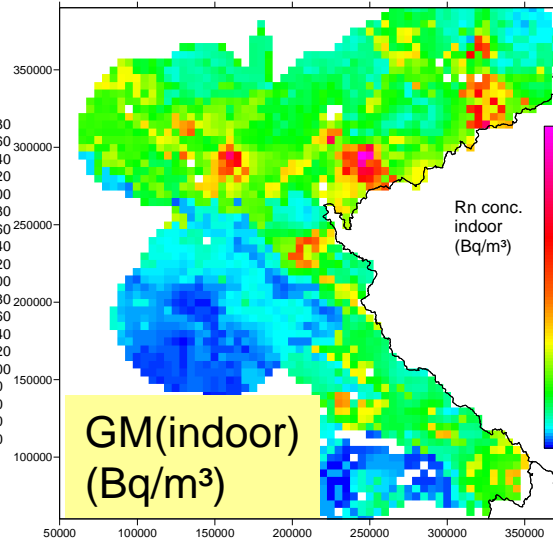
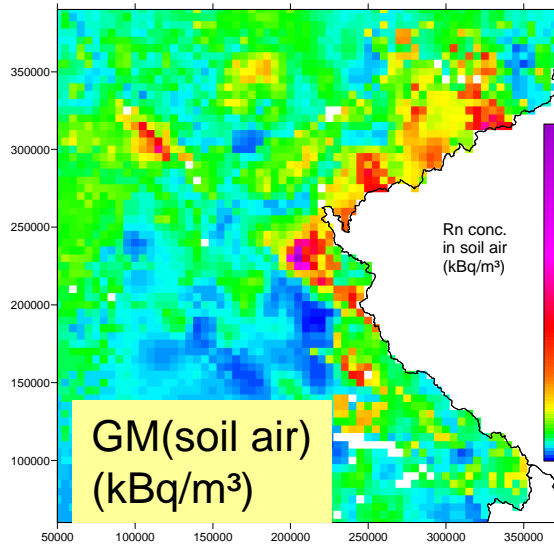
outdoor radon:

factors:

- geology
- distance from N sea



example (4/4)



**result,
separate
estimation**

regression models:

variable	factors
soil Rn, $\ln(Z_s)$	geo-units
indoor Rn, $\ln(Z_{in})$	<ul style="list-style-type: none"> ▪ geo-units ▪ construction year ▪ "East-West"
outdoor Rn, $\ln(Z_{ou})$	<ul style="list-style-type: none"> ▪ geo-units ▪ distance from sea
dose rate, $\ln(Z_d)$	geo-units

sharp edges
along geological
borders !

II. correlation between variables

symbolically:

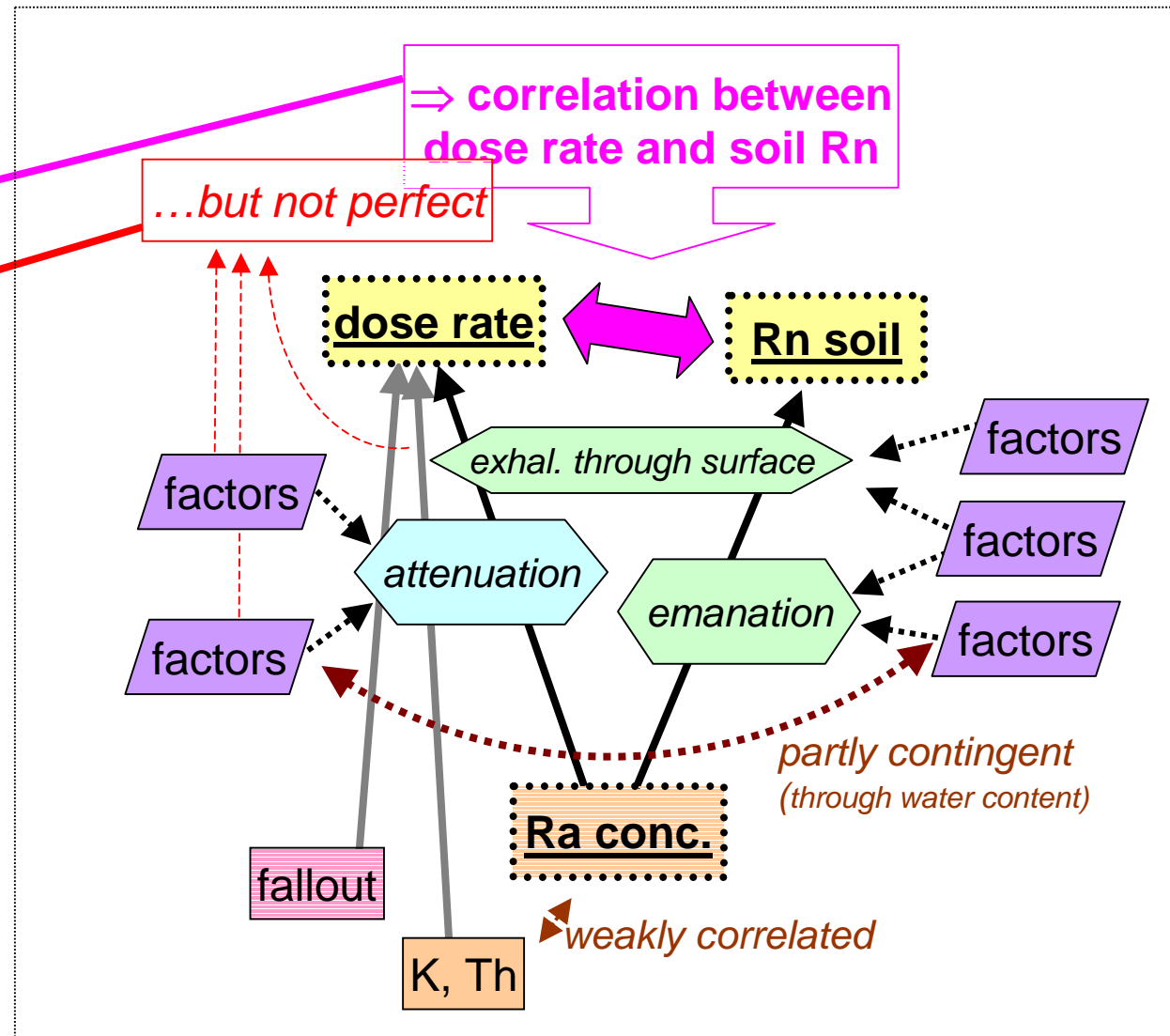
$$(Rn)(x) = (dose\ rate)(x) + \epsilon(x)$$

statistically:

$$\begin{aligned} & cov(Z^1, Z^2), \\ & r^2(Z^1, Z^2), \\ & \tau(Z^1, Z^2) \text{ (rank corr.)} \end{aligned}$$

Z may be derived from original variables, like log, some $g(Z)$,...

variables indexed by upper indices.
Lower indices: sample index



examples for mutual dependencies

- “ $Z^1(U) = f_{\theta}(Z^2(U))$ ”

- Example 1:

$$C_{\text{indoor}} = B + T * C_{\text{soil}}$$

= physical model, valid in steady state

B = building material, T = transfer “factor”

- Example 2:

$$D(\text{ext. dose rate}) = D_0 + \delta C_{\text{soil}}$$

D_0 ... influence of K, T; δ ~ emanation, attenuation,...

- Example 3:

$$C_{\text{outdoor}} = C_0 + \beta C_{\text{soil}} + \gamma * (\text{dist. from sea})$$

C_0 : contribution from distant locations,

β : exhalation (~pressure diff., snow, humidity,...)

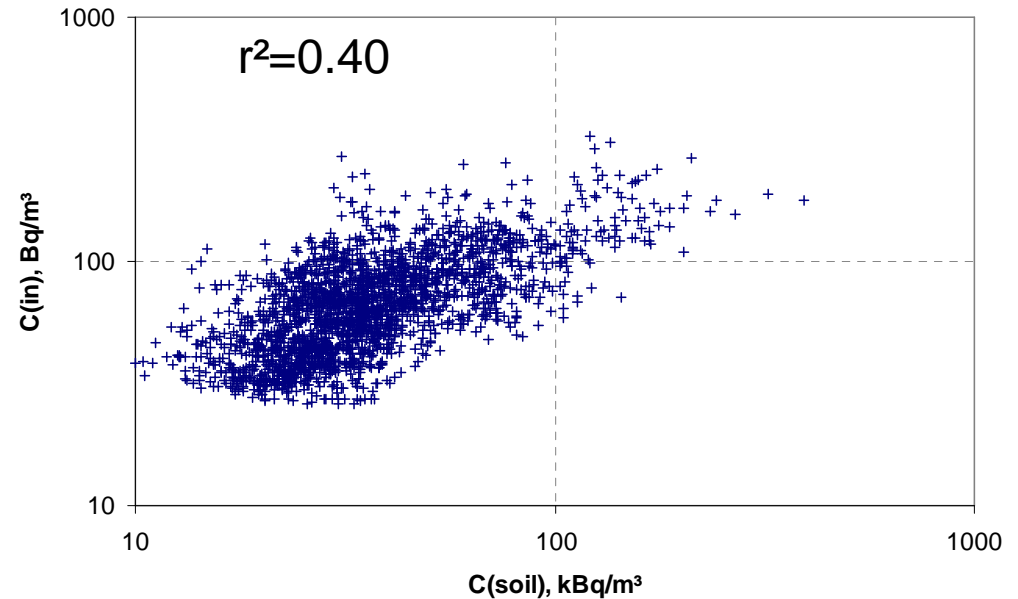
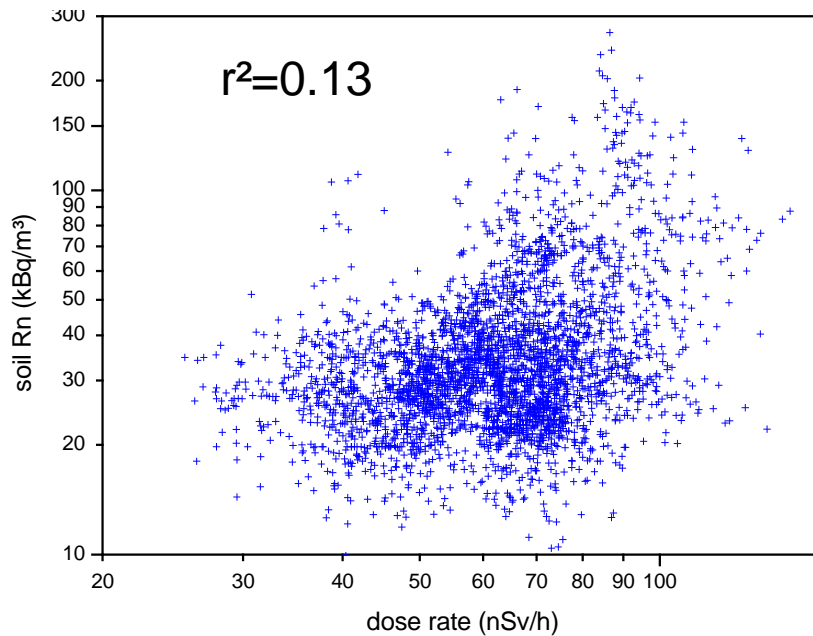
γ : proxy for dilution by “clean” sea air

all model parameters θ are themselves random variables $\theta(x)$

hopefully:
spatial variabilities
of $\theta(x)$ are low.

not yet
examined !

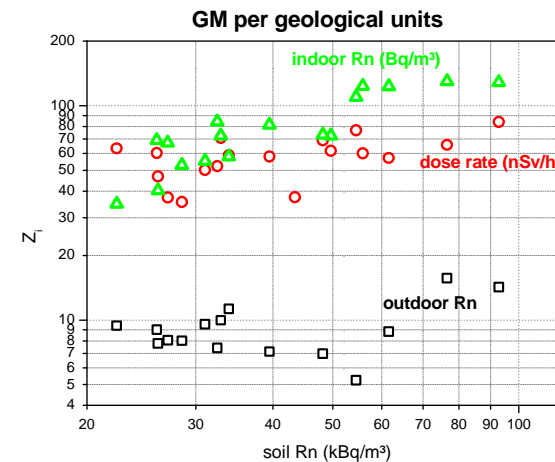
examples: soil Rn ~ dose rate, indoor Rn



so far bad to moderate correlations!

likely reason:

regression models are insufficient,
in particular geological classification
still questionable



III. collocating

- Problem: how to
 - (a) construct target variable of different Z^i
 - (b) estimate correlation parameters θ , if the Z^i not sampled at the same locations ?
- (1) first separate estimation on same set of locations (“sampling set”), e.g. grid or locations of one chosen variable. (“collocated estimation”)
- (2) model afterwards:
 - $Z^i = f_{ij}^{\theta_{ij}}(Z^j)$ transfer models
 - $Y = Y(Z^i)$ target variable

} *later slides!*

excursion: sampling sets

- variable Z^i : samples $\{z_j^i\} \equiv \{z^i(x_j)\}$, $j=1 \dots n^i$
 $\xi^i := \{x_j\} = \text{sampling set}$ of variable Z^i .
- in general: $\xi^i \neq \xi^k$!
(locations of indoor R_n and soil R_n samples are different)
- **joint sampling set**: $\xi := \cup \xi^i$
(=all sampling locations of all variables)
- **grid**: $\Xi = \text{set of grid nodes or cells}$
 $\xi \rightarrow \Xi$ requires interpolation

IV. co-estimation

- estimating variable Z^1 , exploiting information contained in Z^2, \dots, Z^m
- based on cross-covariances
 $C^{ij}(h) = \text{cov}(Z^i(x), Z^j(x+h))$
 - a) co-kriging
 - b) co-simulation
- problems:
 - estimation of co-variograms
 - technical implementation for $m > 2$
- different approach: $\zeta \zeta \zeta$ neural networks and support vector machines ???
- problem: implementation? theoretically more complicated, but seems to have big potential!
- in my view: most statistically satisfying methods!
- not done --- for now!



target variable



*now starts
the more
complicated
part !*

- method 1: “transfer”

first estimate one common Z^1 out of available Z^i using transfer models from collocated estimates $\rightarrow Z^{1(i)*}$, then $Y=f(Z^{1(i)*})$

inspired by H. Friedmann's proposal (JRC geogenic expert group)

- method 2: “multi-variate proper”

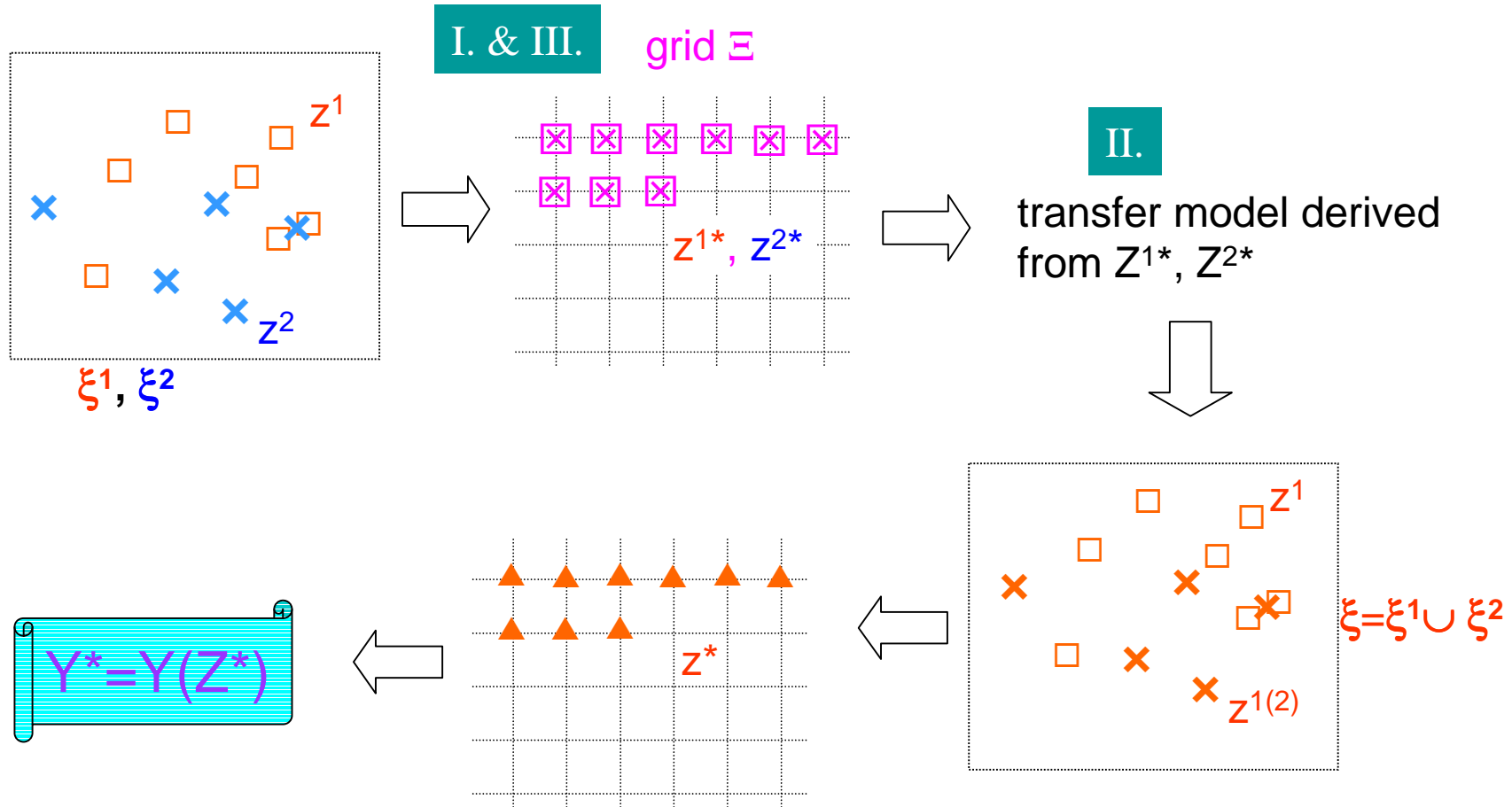
$Y=f(Z^1, \dots, Z^m)$. (Some Z^j may be missing.)

“global” and “local” versions

method 1: cooking recipe

- I. estimate all m variables Z^i **separately** on **common grid** $\Xi \Rightarrow Z^{1*}, \dots, Z^{m*}$ III.
- estimate transfer models between variables, $Z^i = f_{\theta}(Z^j)$, data = $(Z^{i*}, Z^{j*})(\Xi)$ II.
- apply the models to the **original** data $\{z^i\}$ on ξ^i .
- \Rightarrow e.g., $Z^{1(2)} = f^{12, \theta_{12}}(Z^2)$, $Z^{1(3)} = f^{13, \theta_{13}}(Z^3)$
new dataset $\{z^{1, \text{new}}\} := \{z^1\} \cup \{z^{1(2)}\} \cup \{z^{1(3)}\}$ on $\xi = \xi^1 \cup \xi^2 \cup \xi^3$
(Z^1 =soil Rn, Z^2 =indoor Rn, Z^3 =dose rate; $Z^{1, \text{new}}$ = information of all)
- use $\{z^{1, \text{new}}\}$ for modelling $Z^{1, \text{new}}$ on grid Ξ .
(Contains information of original Z^1, Z^2, Z^3 .)
- target variable $Y = Y(Z^{1, \text{new}})$ (could be $Z^{1, \text{new}}$ itself)

method 1



method 2

1. estimate all available or wanted variables Z^i **separately** on **common grid** $\Xi \Rightarrow Z^{1*}, \dots, Z^{m*}$ (= method 1)
wanted but not available: use regression, **I.**
2. multivariate target variable $Y=Y(Z^1, \dots, Z^m)$
3. How to?....
 - multivariate **classification** (CZ, DE, USA, others)
 - **continuous**
 - (a) estimate *local* $F_Y(x)$; (\rightarrow 3rd next slide)
 - (b) *global* distribution G_{Z^1, \dots, Z^m} for a region B . (\rightarrow 4th next slide)

in both cases, spatial association enters through the input variables Z^i .
If $|B|$ large $\Rightarrow G$ independent of B .

proven and viable concept, could be implemented relatively easily

missing input

in general, some Z^j missing!

consistency:

strong version:

at each point x , $Y(Z^1, \dots, \blacksquare^i, \dots, \blacksquare^j, \dots, Z^m)$ must be independent of which Z^i are missing (up to statistics).

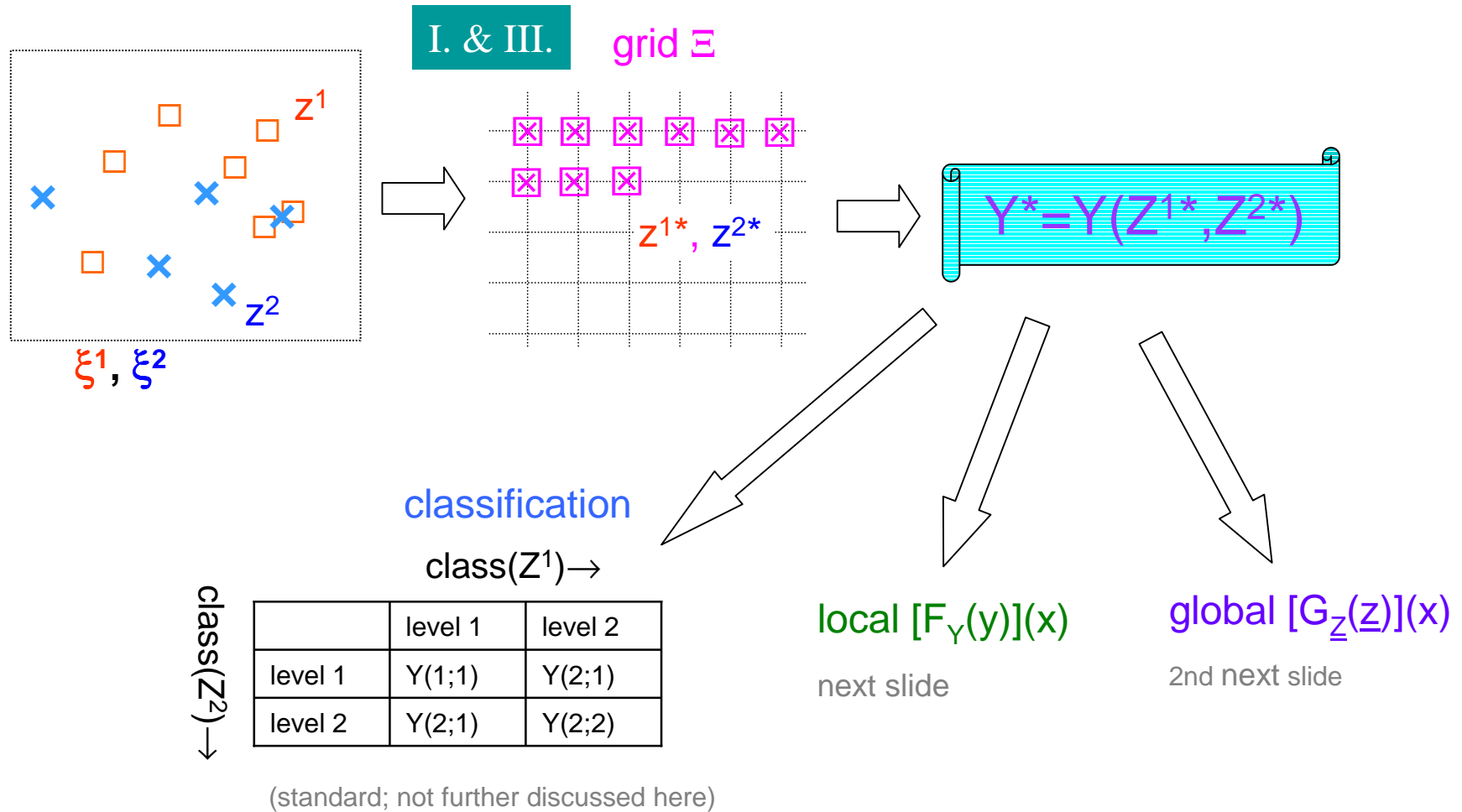
E.g.: $Y(Z^1, Z^2, Z^3) \cong Y(Z^1, Z^2, \blacksquare)$... maybe not realistic

weak version: “conservative”

$Y(Z^1, Z^2, \blacksquare) \geq Y(Z^1, Z^2, Z^3)$... \geq means “higher hazard”

Serious constraint on admissible Y !

method 2

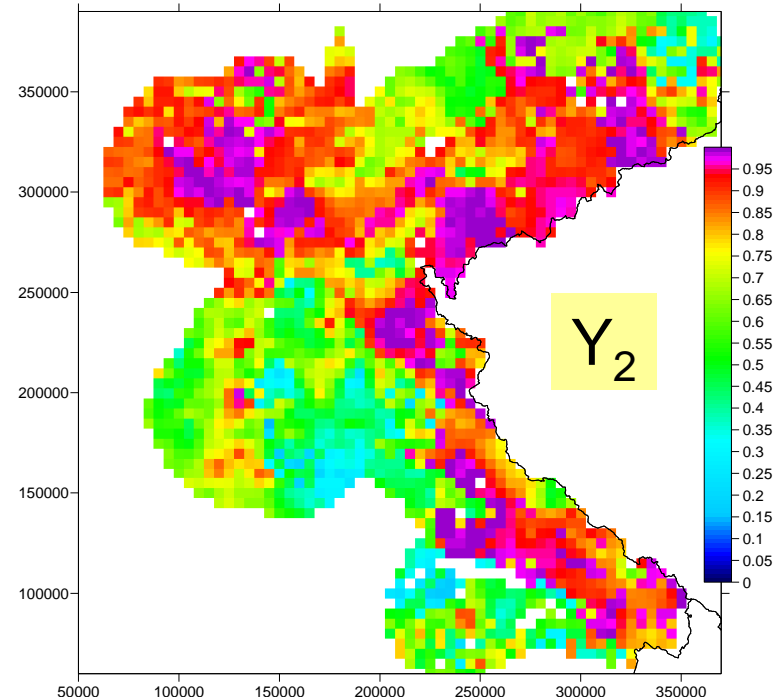
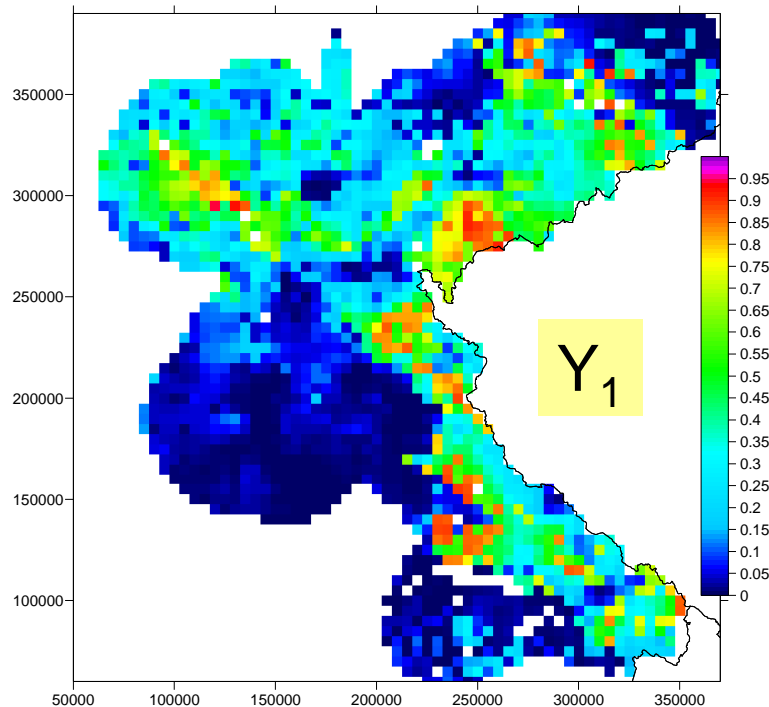


example for Y from local $F_Y(x)$ (a)

- at x , generate many realizations
 $(z^1, \dots, z^m)^{(k)} =: \underline{z}^{(k)}$
using global $\text{cov}(Z^1, \dots, Z^m)$ from II.
- for each $\underline{z}^{(k)} \rightarrow Y^{(k)}$
- statistics of Y over (k) (e.g. $\text{AM}\{Y^{(k)}\}$) $\rightarrow Y$ } $\forall x$
- technically: e.g. Cholesky method: A , such that $AA^T = \text{cov}$.;
generate $u^i \sim N(0, 1)$ indep.;
 $\Rightarrow \underline{z} = \underline{\mu} + A \underline{u} \sim N(\underline{\mu}, \text{cov})$
- needs to be demonstrated yet !

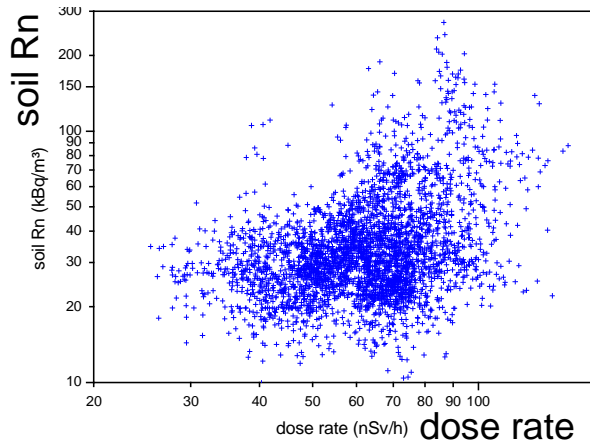
example for Y from global G (b); 1/4

- $Y_1(x) := \text{prob}(z^1 < \zeta^1, \dots, z^m < \zeta^m) \equiv F_{\underline{z}}(\underline{\zeta})$ ($\underline{z} := (z^1, \dots, z^m)$)
- $Y_2(x) := 1 - \text{prob}(z^1 > \zeta^1, \dots, z^m > \zeta^m)$
- **3-variate scoring** ($m=3$): indoor conc., soil conc., dose rate
- use (for now) empirical distributions F_{emp}



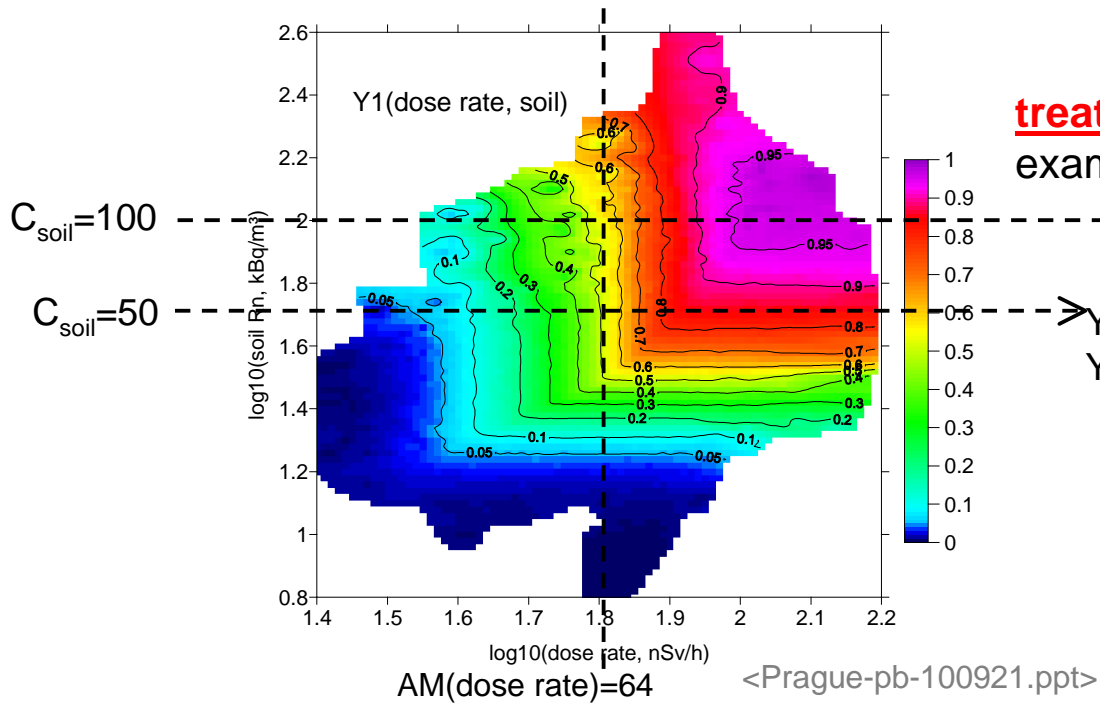
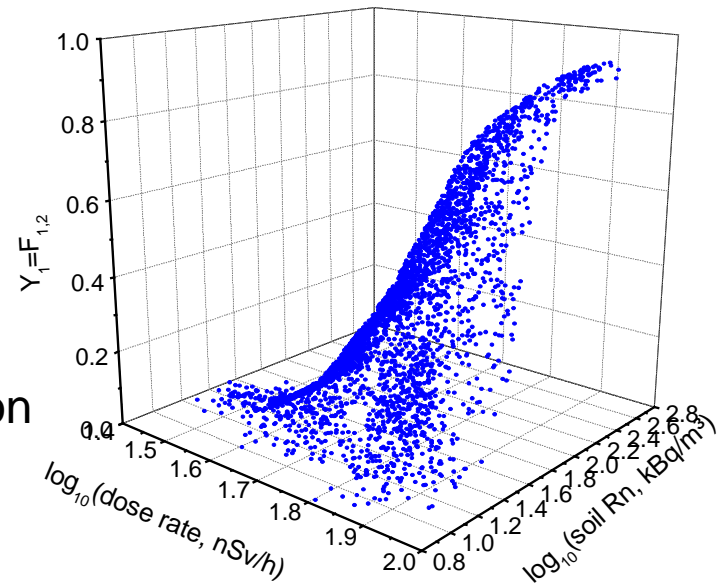
missing: e.g., $Y_1(\zeta^1, \zeta^2, \blacksquare)(x) = F_{\underline{z}}(\zeta^1, \zeta^2, \infty)$ (not yet implemented)

continued, 2/4: soil Rn & dose rate



soil Rn ~ dose rate:
bad correlation !
 $r^2=0.13$

bivariate distribution
 $Y_1 = F_{Z^1, Z^2}$



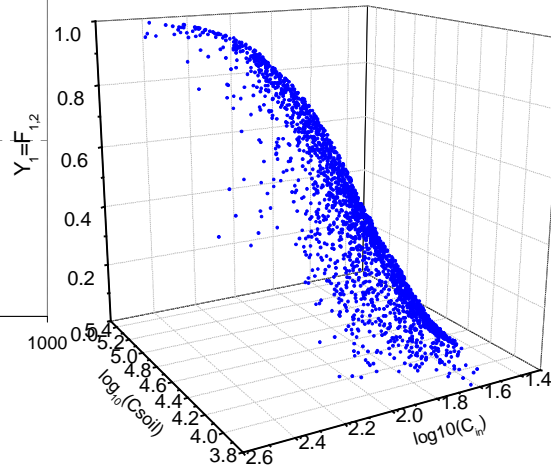
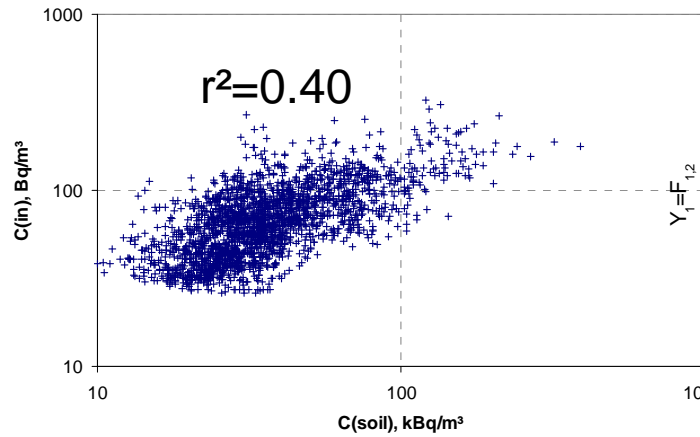
treatment of missing variable,
example: dose rate = missing

$$\begin{aligned} &\rightarrow Y_1(100, \infty) = 1 \\ &Y_1(100, AM(\text{dose rate})) = 0.55 \\ &\rightarrow Y_1(50, \infty) = 0.83 \\ &Y_1(50, AM(\text{dose rate})) = 0.55 \end{aligned}$$

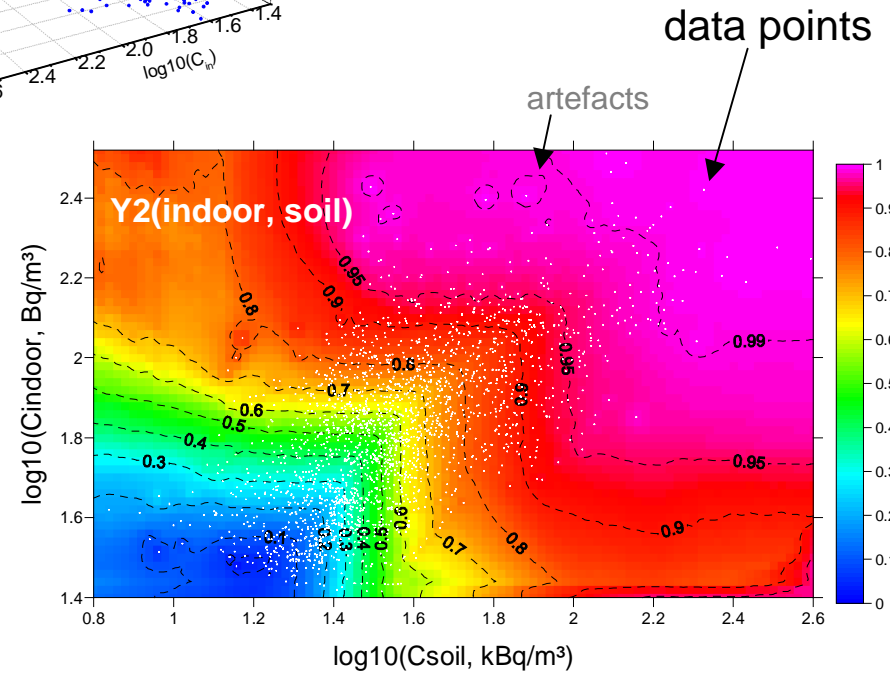
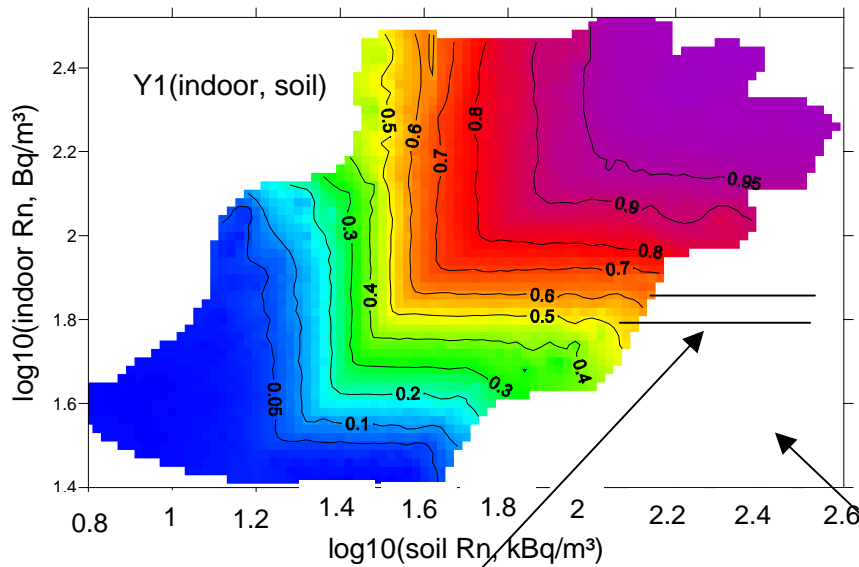
does not require a priori knowledge of $F_1(Z^1, Z^2)$!

requires a priori knowledge of $F_1(Z^1, Z^2)$!

continued, 3/4: soil Rn & indoor Rn



Y_2 more conservative than Y_1 !



theoretically: straight lines

area without data support blanked

improvement / generalization

on $[0,1]^n$, define a p-norm $\|\cdot\|_p$, (\rightarrow metric space)

$$Y^p(\mathbf{x}) := \|\mathbf{F}_{\underline{z}}(\mathbf{x})\|_p := (m^{-1/p})(\sum_i (F_{z_i}(\mathbf{x}))^p)^{1/p}$$

n= number of possible covariates,

m= number of available covariates

projection onto main diagonal of sub-cube $[0,1]^m$.

i.e. Y^p is functional f^p : $f^p[\underline{Z}(\mathbf{x})]=\|\mathbf{F}_{\underline{z}}(\mathbf{x})\|_p$

$p=0$... corresponds score 1 (Y_1)

$p=\infty$... corresponds score 2 (Y_2 , „maximum norm“)

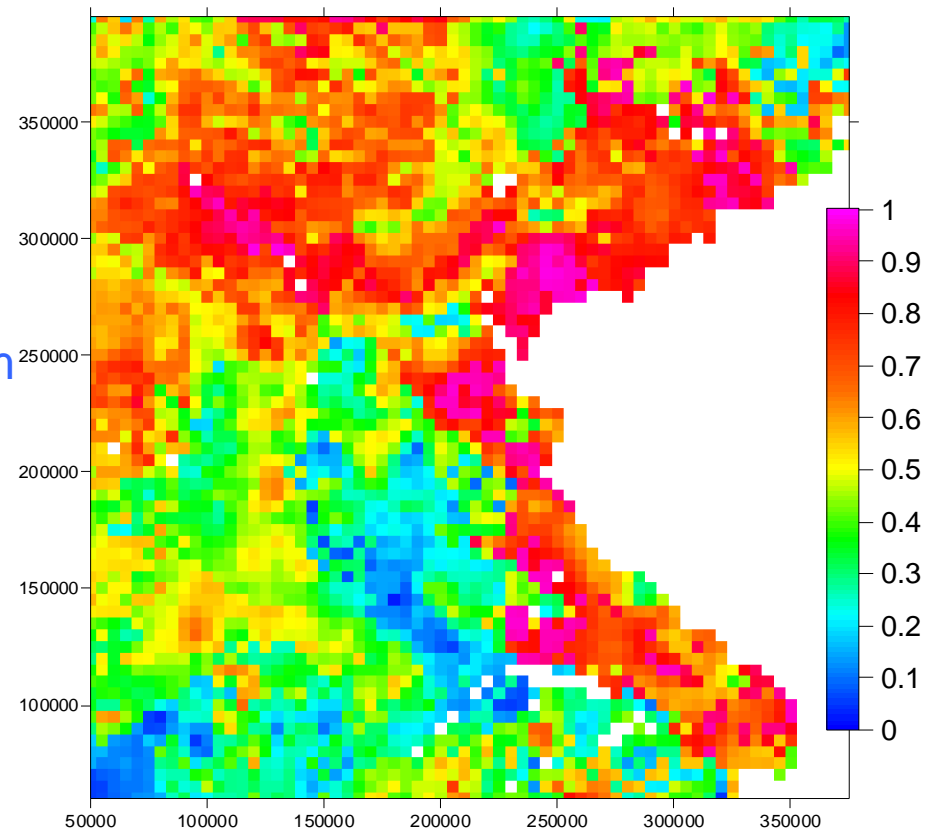
$p=1$... inner product ($\mathbf{F}_{\underline{z}}(\mathbf{x}) \bullet \underline{1}$)

Y^1 = joint F_{z_i} as if the z_i were perfectly correlated; Y^1 =average of F_{z_i}

advantage of concept: more adaptive and flexible !

natural treatment of missing data if $m < n$

example: $p=2$;
soil Rn, indoor Rn, dose rate



4/4: in practice

- how to estimate $Y(z^1, \dots, z^m)$?
 1. given F_{z^1}, \dots, F_{z^m} from previous experience
 2. nscore $z \rightarrow w \Rightarrow F_{w^1}, \dots, F_{w^m} \sim N(\mu^i, \sigma^i); \sigma^{ij}$
 3. values $\{z_k^i\} \rightarrow \{w_k^i\}$ from transform table \Rightarrow values = m-tuples $(w^1, \dots, w^m)_k$
 4. generate N (many!) (u^1, \dots, u^m) random $\sim F_{\underline{w}}$
 5. $Y_1(\underline{z}) = Y_1(\underline{w}) = \text{prob}(w'^1 < w^1, \dots, w'^m < w^m) \approx \min\{\#(u^1 < w^1)/N, \dots, \#(u^m < w^m)/N\}$
 Y_2 and Y^p in analogy
- if sufficient data for F_{z^i} : easy to implement automatically !

Σ method 2: cooking recipe

1. select your input variables
(= the ones you have)
2. regression modelling on available factors **I.**
(geology,...) $\rightarrow Z^i = \text{factor}_{1,j} + \text{factor}_{2,jk} + \dots + f(\text{factor}_n) + \dots$
3. estimate on common grid $\rightarrow (z^{1*}, \dots, z^{m*})(\Xi)$ **III.**
4. correlation analysis $\rightarrow \text{cov}(Z^i, Z^j)$ **II.**
5. select target variable Y
.... Z^1 // classification // F_Y // G_Z **V.**
6. estimate Y

1 - 4: \pm straight forward;
5 - 6: still problems

problems (as usual)



1. **Geological control: proper definition of geological classes remains to be done.**
⇒ hopefully also better correlation between variables!
2. **Definition of input variables, in particular soil-Rn, permeability**
3. **Regression modelling: data! data! data!**
4. **target quantity: more discussion still needed.**
5. **estimation methods to be improved.**

literature review would be helpful!

*here !
round table Friday!*

*dig deeper into geo-
statistical methodology!*