

Ordinal Compositional Data and Time Series



HELMUT SCHMIDT
UNIVERSITÄT

Universität der Bundeswehr Hamburg

**MATH
STAT**

Christian H. Weiß

Department of Mathematics & Statistics,
Helmut Schmidt University, Hamburg

Funded by DFG – Projektnr. 516522977.



HELMUT SCHMIDT
UNIVERSITÄT
Universität der Bundeswehr Hamburg

**MATH
STAT**

Ordinal Compositional Data and Time Series

■

 ■
Introduction

CoDa = “proportions of some whole” (Aitchison, 1986).

$\mathbf{x} = (x_0, \dots, x_d)^\top$ normalized $(d + 1)$ -part composition
with $d \in \mathbb{N} = \{1, 2, \dots\}$ if \mathbf{x} from $(d + 1)$ -part unit simplex
 $\mathcal{S} := \{\mathbf{x} \in (0; 1)^{d+1} \mid x_0 + \dots + x_d = 1\}$.

Connection to **categorical data** (Agresti, 2002):

CoDa vectors $\mathbf{p} \in \mathcal{S}$ might serve as PMF of
categorical RV Q with qualitative range $\mathcal{S} = \{s_0, \dots, s_d\}$.

In many applications, categories behind CoDa unordered,
so \mathcal{S} is nominal range and Q nominal RV,

notation “ s_0, \dots, s_d ” just uses some lexicographic order.

For such nominal CoDa, reasonable to require that “applying a compositional analysis, the information due to the order of the different classes, plays no role”, see Pawlowsky-Glahn & Buccianti (2011).

But as conceded by Pawlowsky-Glahn & Buccianti (2011), “in some cases the parts can be assumed to be ordered”, so \mathcal{S} exhibits natural order $s_0 < \dots < s_d$ and Q is ordinal RV (Agresti, 2010). Then, $\mathbf{x} \in \mathcal{S}$ called **ordinal composition**.

Example: $\mathbf{x} = (x_0, x_1, x_2)^\top$ expressing proportions of people in ordered age categories < 15 (x_0), $15\text{--}60$ (x_1), and > 60 (x_2).

If concerned with **ordinal CoDa**, certainly still justified to apply well-established CoDa approaches, although ignoring “information due to the order of the different classes”.

However, gain additional insights if existing CoDa approaches complemented by new ones accounting for order within \mathcal{S} .

Such **new ordinal CoDa approaches** derived by adapting well-established concepts from ordinal data analysis.

Benefits illustrated for *descriptive analysis* of ordinal CoDa, *statistical inference* from ordinal CoDa, *monitoring* of ordinal CoDa process, and ordinal compositional *time series* (CoTS).



HELMUT SCHMIDT
UNIVERSITÄT
Universität der Bundeswehr Hamburg

**MATH
STAT**

Ordinal Compositional Data

■

 ■
Descriptive Analysis

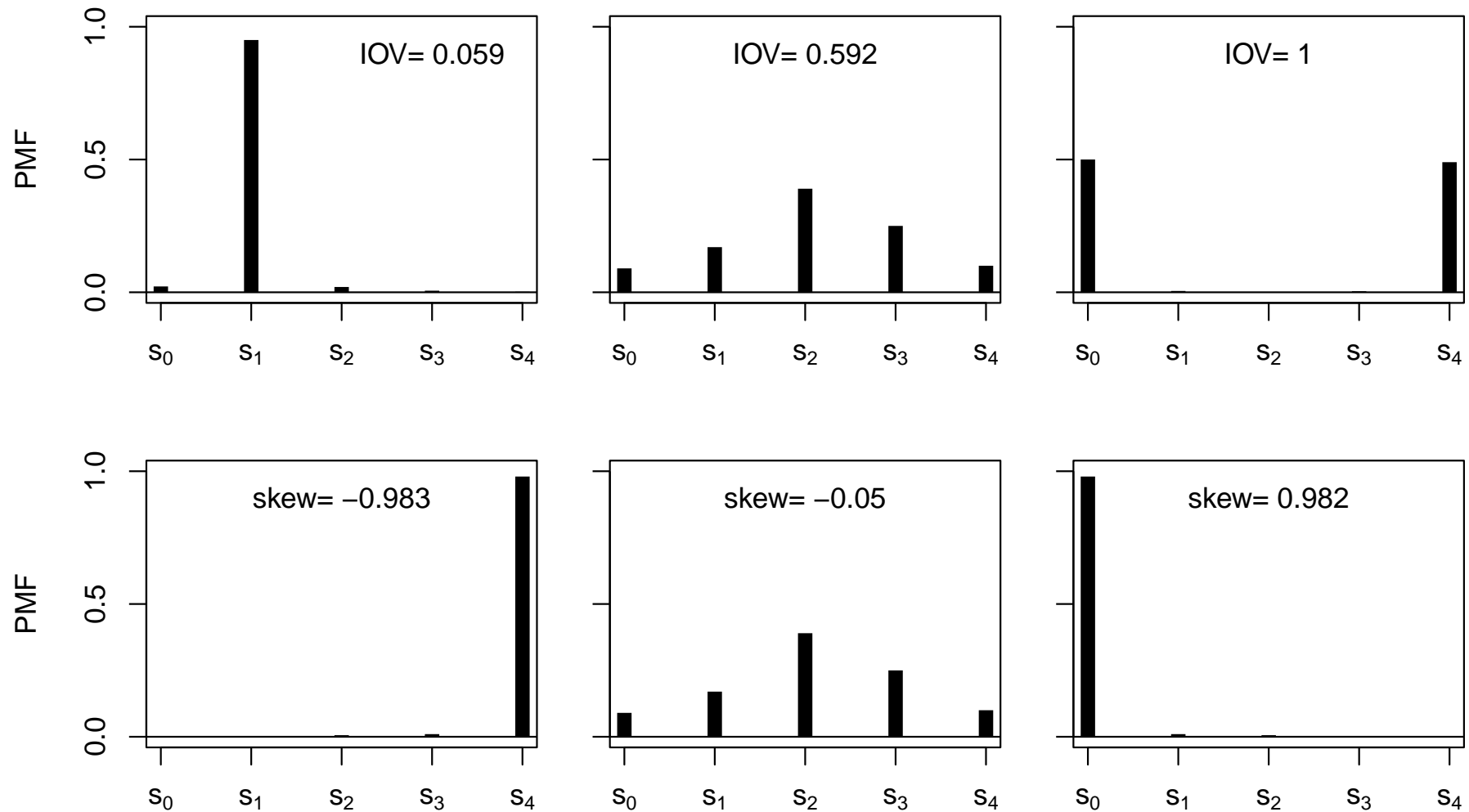
For **ordinal RV** Q , one does not focus on PMF $\mathbf{p} \in \mathbb{S}$, but on d -dimensional CDF vector $\mathbf{f} = (f_0, \dots, f_{d-1})^\top \in [0; 1]^d$, where $f_j = P(Q \leq s_j)$ (note that $f_d = 1$), because accumulation accounts for natural order in range \mathcal{S} .

Formally: $\mathbf{f} = \mathbf{T} \mathbf{p}$, where $\mathbf{T} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 1 & \dots & 1 & 0 \end{pmatrix}$.

Dispersion of Q , e. g., by $\text{IOV}(\mathbf{f}) = \frac{4}{d} \sum_{i=0}^{d-1} f_i(1 - f_i)$.
(one-point vs. extreme two-point distribution)

Skewness of Q , e. g., by $\text{skew}(\mathbf{f}) = \frac{2}{d} \sum_{i=0}^{d-1} f_i - 1$.
(extreme left (right) one-point vs. symmetric distribution)

Examples of IOV and skew:



These ordinal measures can be applied to **ordinal CoDa** $x \in \mathbb{S}$:

First, accumulate CoDa vectors to $c := \mathbf{T} x$, then evaluate ordinal dispersion and skewness via $\text{IOV}(c)$ and $\text{skew}(c)$.

Example: ageCatWorld data from R-package robCompositions.

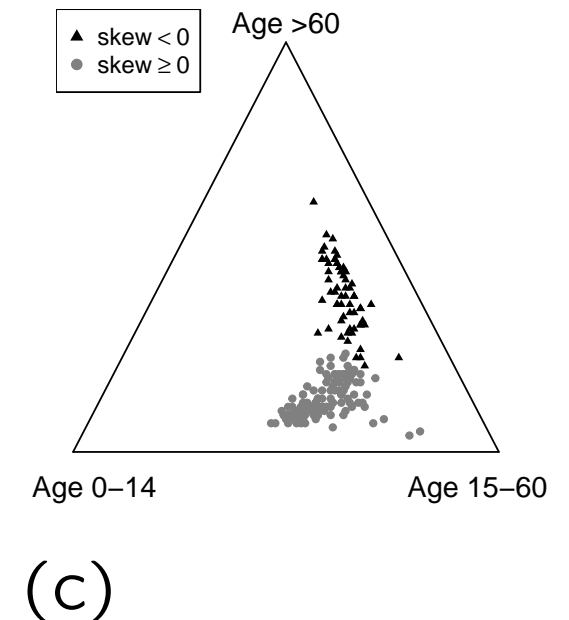
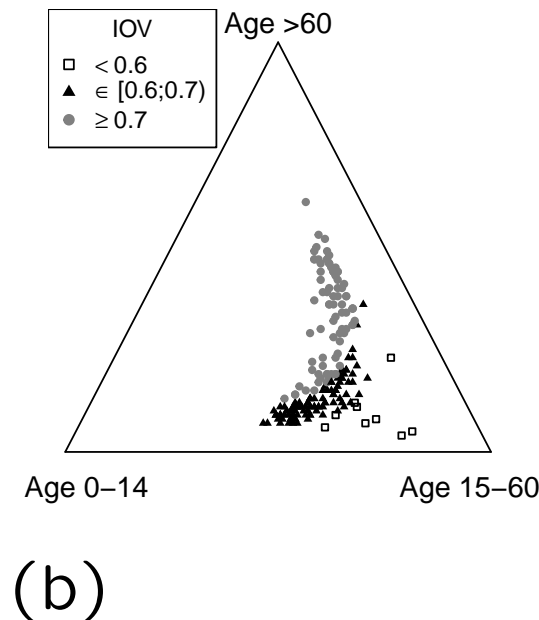
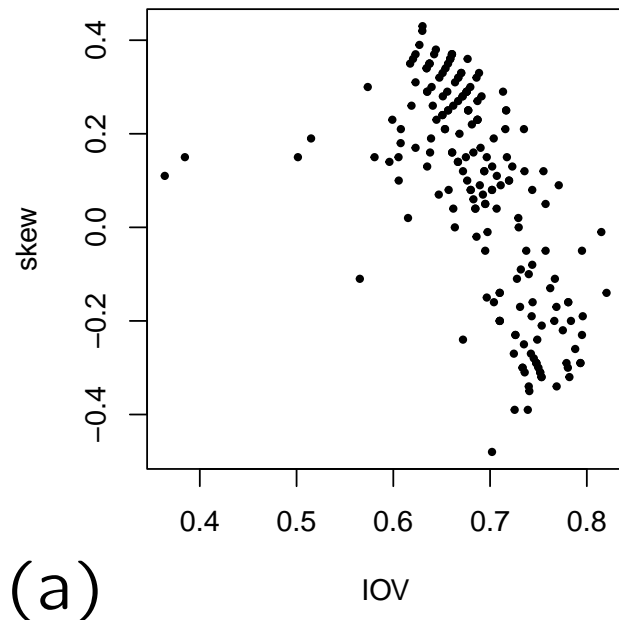
Age proportions $x_i = (x_{i,0}, x_{i,1}, x_{i,2})^\top$ in $n = 195$ countries, people with age < 15 ($x_{i,0}$), 15–60 ($x_{i,1}$), and > 60 ($x_{i,2}$).

Location measures not much discriminative power:

median category s_1 in 191 out of 195 cases,

mode category s_1 in 163 out of 195 cases.

IOV-skew diagram and ternary diagram:



IOV-skew diagram easy to interpret. In the present example, similar insights into data as ternary diagram.

But IOV-skew diagram also for higher-dimensional CoDa!



HELMUT SCHMIDT
UNIVERSITÄT
Universität der Bundeswehr Hamburg

**MATH
STAT**

Ordinal Compositional Data

■

 ■
Statistical Inference

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be i. i. d. ordinal CoDa with mean \mathbf{p} ,
then standard CLT implies

$$\sqrt{n} (\bar{\mathbf{X}} - \mathbf{p}) \sim N(\mathbf{0}, \Sigma), \quad \Sigma = (\sigma_{ij})_{i,j=0,\dots,d}, \quad \sigma_{ij} = \text{Cov}[X_i, X_j].$$

Accumulate $\mathbf{C}_k = \mathbf{T} \mathbf{X}_k$ with $E[\mathbf{C}_k] = \mathbf{f} = \mathbf{T} \mathbf{p}$, then

$$\sqrt{n} (\bar{\mathbf{C}} - \mathbf{f}) \sim N(\mathbf{0}, \Sigma') \quad \text{with } \Sigma' = \mathbf{T} \Sigma \mathbf{T}^\top, \quad \sigma'_{ij} = \sum_{r=0}^i \sum_{s=0}^j \sigma_{rs}.$$

Asymptotics of $\text{IOV}(\bar{\mathbf{C}})$ and $\text{skew}(\bar{\mathbf{C}})$

via Taylor expansions (“Delta method”) ...

Theorem 1: $\text{IOV}(\bar{\mathbf{C}})$ and $\text{skew}(\bar{\mathbf{C}})$ asymptotically normally distributed with

$$E[\text{IOV}(\bar{\mathbf{C}})] \approx \text{IOV}(\mathbf{f}) - \frac{1}{n} \frac{4}{d} \sum_{i=0}^{d-1} \sigma'_{ii},$$

$$V[\text{IOV}(\bar{\mathbf{C}})] \approx \frac{1}{n} \frac{16}{d^2} \sum_{i,j=0}^{d-1} (1 - 2f_i)(1 - 2f_j) \sigma'_{ij},$$

and

$$E[\text{skew}(\bar{\mathbf{C}})] = \text{skew}(\mathbf{f}), \quad V[\text{skew}(\bar{\mathbf{C}})] \approx \frac{1}{n} \frac{4}{d^2} \sum_{i,j=0}^{d-1} \sigma'_{ij}.$$

Example: closed-form formulae for Dirichlet distribution.

Simulation study: good finite-sample performance.

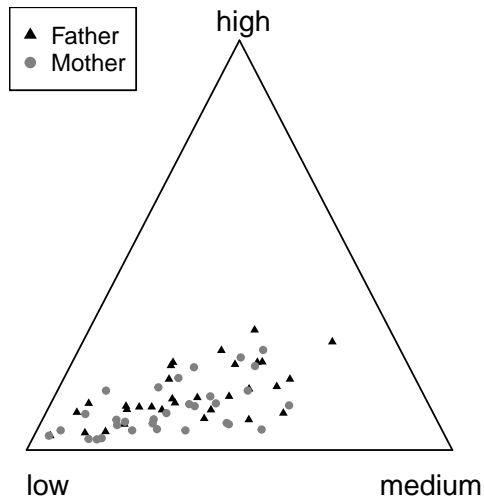
Example: educFM data from R-package `robCompositions`.
Proportions of low (s_0), medium (s_1), and high (s_2) education levels of father ($\mathbf{x}_{f;i}$) and mother ($\mathbf{x}_{m;i}$) in $n = 31$ European countries.

Category s_0 (“low”) is median (mode) of $\mathbf{x}_{f;i}$ in 19 (22) cases, and of $\mathbf{x}_{m;i}$ in 22 (25) cases; otherwise, it is s_1 .

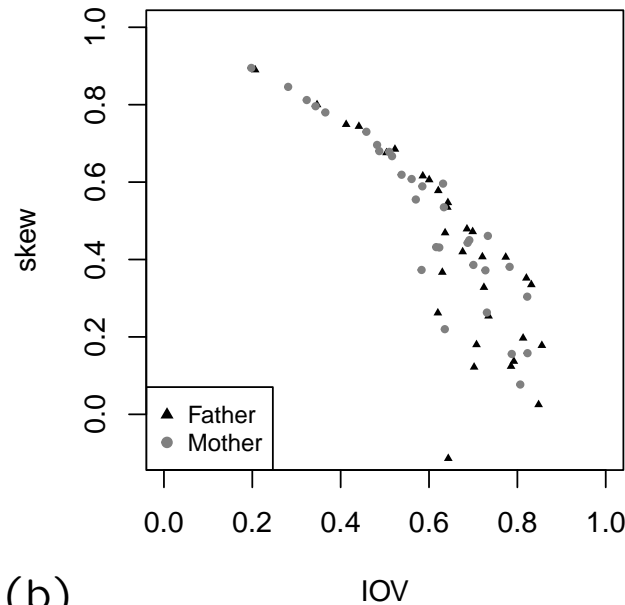
Ternary plot: concentration in low-to-medium corner.

IOV-skew diagram: low dispersion if strong positive skewness.

Hardly any differences between father and mother visible.



(a)



(b)

Father	Estim.	Bias	SE
IOV(\bar{C})	0.731	-0.003	0.022
skew(\bar{C})	0.414	0.000	0.043
Mother			
IOV(\bar{C})	0.652	-0.002	0.033
skew(\bar{C})	0.516	0.000	0.041

(c)

Fit logistic normal distributions (Aitchison, 1986) to CoDa, approximate bias and standard error (SE):

Significantly (5%-level) less dispersion and stronger positive skewness for mothers \Rightarrow unequal opportunities for education.



HELMUT SCHMIDT
UNIVERSITÄT
Universität der Bundeswehr Hamburg

**MATH
STAT**

Monitoring Ordinal Compositional Data

■

 ■
Control Charts

Default control chart for i. i. d. CoDa process: χ^2 -chart

$$X_t^2 = (\mathbf{Y}_t - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_0^{-1} (\mathbf{Y}_t - \boldsymbol{\mu}_0) \quad \text{for } t = 1, 2, \dots$$

Proposal: Complement by **IOV-** and **skew-chart**, preferably with additional EWMA smoothing, i. e.,

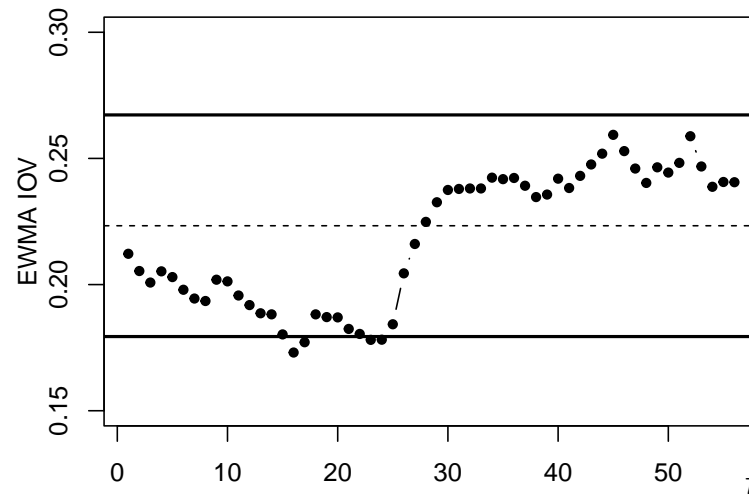
$$C_{t,\lambda} = \lambda C_t + (1 - \lambda) C_{t-1,\lambda} \quad \text{with } C_{0,\lambda} = f_0 \quad \text{and } \lambda \in (0; 1).$$

Plot $\text{IOV}(C_{t,\lambda})$

or $\text{skew}(C_{t,\lambda})$

for $t = 1, 2, \dots$

against control limits:





HELMUT SCHMIDT
UNIVERSITÄT
Universität der Bundeswehr Hamburg

**MATH
STAT**

Ordinal Compositional Time Series

■

 ■
Stochastic Modelling

Existing models for CoTS do not account for ordering.

Example: conditional Dirichlet model (Zheng & Chen, 2017),

$\mathbf{X}_t | \mathcal{F}_{t-1} \sim \text{Dir}(\mathbf{p}_t, \nu)$ with $E[\mathbf{X}_t | \mathcal{F}_{t-1}] = \mathbf{p}_t$ and $\nu > 0$.

“ARMA recursion” for $\mathbf{p}_t = \phi(\mathbf{X}_{t-1}, \dots, \mathbf{X}_{t-p}, \mathbf{p}_{t-1}, \dots, \mathbf{p}_{t-q})$,
with AR-order p and feedback terms $\mathbf{p}_{t-1}, \dots, \mathbf{p}_{t-q}$.

Idea: To account for natural order, combine with cumulative logit approach in Pruscha (1993), Fokianos & Kedem (2003).

Proposal: Let $F_{0,1}(x) = [1 + \exp(-x)]^{-1}$ be inverse logit.

Define $\tilde{\mathbf{X}}_t = (X_{t,0}, \dots, X_{t,d-1})^\top$ and $\tilde{\mathbf{p}}_t = (p_{t,0}, \dots, p_{t,d-1})^\top$,

let $\mathbf{f}_t = (f_{t,0}, \dots, f_{t,d-1})^\top$ be conditional CDF vector.

Full ARMA-type model for ordinal CoTS by

$$f_{t,i} = F_{0,1}\left(\eta_i + \sum_{k=1}^p \alpha_k^\top \tilde{\mathbf{X}}_{t-k} + \sum_{l=1}^q \beta_l^\top \tilde{\mathbf{p}}_{t-l}\right) \quad \text{for } i = 0, \dots, d-1,$$

with $d(p + q + 1)$ parameters and $-\infty < \eta_0 < \dots < \eta_{d-1} < +\infty$.

More parsimonious model with $d + p + q$ parameters:

$$f_{t,i} = F_{0,1}\left(\eta_i + \sum_{k=1}^p \alpha_k \mathbf{1}^\top \tilde{\mathbf{X}}_{t-k} + \sum_{l=1}^q \beta_l \mathbf{1}^\top \tilde{\mathbf{p}}_{t-l}\right).$$

Include possible covariate z_t by summand “ $+\gamma^\top z_t$ ”.

Dirichlet model with additional scale parameter ν :

$\mathbf{X}_t | \mathcal{F}_{t-1} \sim \text{Dir}(\mathbf{p}_t, \nu)$, where \mathbf{p}_t from \mathbf{f}_t by discrete differences.

Example: Yearly proportions of categories “not overweight” (BMI < 25), “overweight” (BMI in [25; 30)), and “obese” (BMI ≥ 30) in Germany for period 1975–2016 (data from WHO).

Partition: $\mathbf{x}_1, \dots, \mathbf{x}_{36}$ (1975–2010) for model fitting, $\mathbf{x}_{37}, \dots, \mathbf{x}_{42}$ (2011–2016) for out-of-sample forecasting.

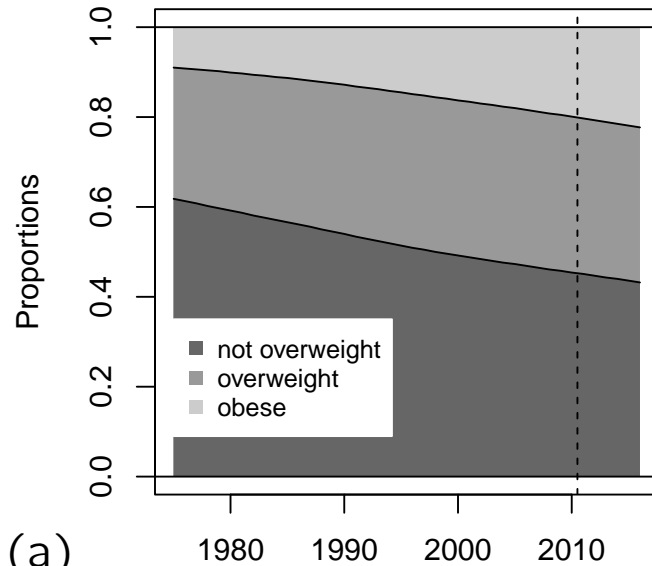
For similar data set, Mills (2010) proposed **linear alr-model**:

$$(1) \quad E[\text{alr}(\mathbf{x}_t)_i] = a_i + b_i t \quad \text{for } i = 1, \dots, d.$$

$$(2) \quad f_{t,i} = F_{0,1}(\eta_i + \gamma_i t) \quad \text{for } i = 0, \dots, d - 1.$$

$$(3) \quad f_{t,i} = F_{0,1}(\eta_i + \gamma_i t + \boldsymbol{\alpha}_1^\top \tilde{\mathbf{X}}_{t-1}) \quad \text{for } i = 0, \dots, d - 1.$$

Plot of proportions and table with model fits and CSS values:



Model	i	η_{i-1} or a_i	γ_{i-1} or b_i	$\alpha_{1,i-1}$	CSS _{in}	CSS _{out}
(1)	1	1.957	-0.032		1.135	2.119
	2	1.230	-0.018			
(2)	1	0.480	-0.019		0.698	1.269
	2	2.350	-0.027			
(3)	1	0.469	-0.016	0.322	0.410	0.379
	2	2.354	-0.025	-0.651		

(b)

Novel (logistic-)linear model, especially with AR(1)-component, clearly superior in conditional sum of squares (CSS, times 10^3):

$$\text{CSS}(\boldsymbol{\theta}) := \sum_t \|\mathbf{X}_t - \mathbf{p}_t\|^2.$$

- Beneficial to consider natural order of ordinal CoDa:
 - (visual) descriptive analysis of ordinal CoDa,
 - statistical inference from i. i. d. CoDa samples,
 - EWMA control charts for i. i. d. CoDa processes,
 - conditional regression models for ordinal CoTS.
- **Future research:** Among others, ...
 - performance of ordinal CoDa control charts for *serially dependent* CoDa (i. e., monitoring of CoTS);
 - *analysis* of ordinal CoTS by adapting methods from ordinal time series analysis (Weiß, 2020).

Thank You for Your Interest!



HELMUT SCHMIDT
UNIVERSITÄT

Universität der Bundeswehr Hamburg

**MATH
STAT**

Christian H. Weiß

Department of Mathematics & Statistics

Helmut Schmidt University, Hamburg

weissc@hsu-hh.de

*This research was funded by the
Deutsche Forschungsgemeinschaft
(DFG, German Research Foundation),
Projektnummer 516522977.*

Weiß (2023) Ordinal compositional data and time series.
Statistical Modelling, in press.

Agresti (2002) *Categorical Data Analysis*. Wiley.

Agresti (2010) *Analysis of Ordinal Categorical Data*. Wiley.

Aitchison (1986) *Statistical Analysis of Compositional Data*. Chapman.

Filzmoser et al. (2018) *Applied Compositional Data Analysis*. Springer.

Fokianos & Kedem (2003) Regression theory for ... *Stat Sci* **18**, 357–376.

Mills (2010) Forecasting compositional ... *Qual & Quant* **44**, 673–690.

Pawlowsky-Glahn & Buccianti (2011) *Compositional Data Analysis*. Wiley.

Pruscha (1993) Categorical time series ... *Statistics* **24**, 43–57.

Vives-Mestres et al. (2014) Individual T^2 control ... *JQT* **46**, 127–139.

Weiß (2020) Distance-based analysis of ... *JASA* **115**, 1189–1200.

Zheng & Chen (2017) Dirichlet ARMA models ... *JMA* **158**, 31–46.