

Analyzing Categorical Time Series in the Presence of Missing Observations



HELMUT SCHMIDT
UNIVERSITÄT

Universität der Bundeswehr Hamburg

**MATH
STAT**

C.H. Weiß

Department of Mathematics & Statistics,
Helmut Schmidt University, Hamburg

In real applications, time series often exhibit **missing data**, so impossible to apply standard analytical tools.

For real-valued time series with “missingness”, several proposals how to adapt tools, such as sample autocorrelation function (ACF) or spectral estimators, see Scheinok (1965), Bloomfield (1970), Neave (1970), Duns-muir & Robinson (1981), Yajima & Nishino (1999).

Use idea of Parzen (1963) to understand real-valued time series with missingness as resulting from **amplitude modulation**, where amplitude-modulating process binary.

But missingness also happens to **categorical time series**, which consist of qualitative values ordered in time.

Completely different analytical tools for cat. t. s. (Weiß, 2018), so aforementioned solutions for missingness not applicable.

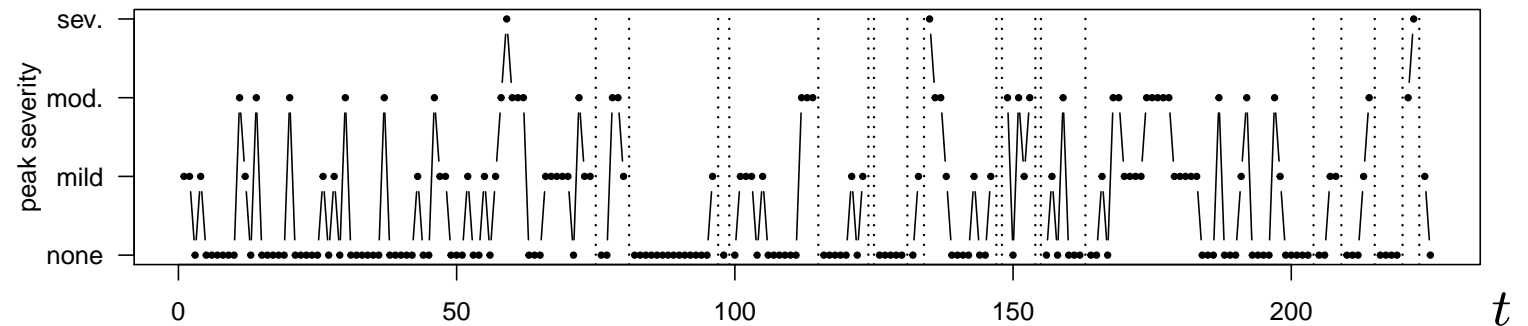
Main practical motivation: collaborative project with Marina Vives-Mestres & Amparo Casanova (Curelator Inc.), on categorical time series from migraine patients.

Daily data on migraine patients from questionnaire in mobile app N1-HeadacheTM. **Missing data** because patients skipped some questions, or stopped before completing questionnaire.

Examples: (dotted lines indicate missing data)

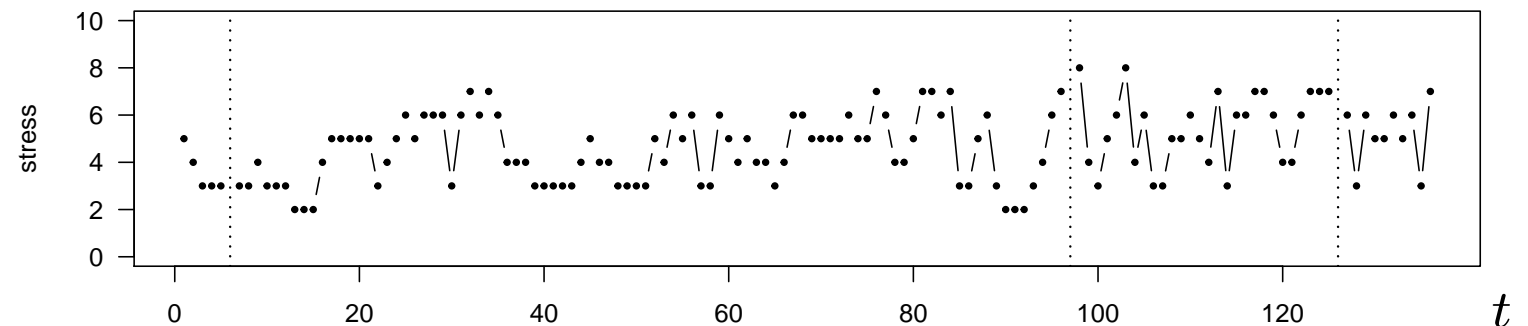
Peak severity (“none”, “mild”, “moderate”, “severe”):

daily levels
for patient A;
 $n = 225$,
19 missing.



Stress (0–10 Likert scale, “not at all” to “a lot”):

daily levels
for patient B;
 $n = 136$,
3 missing.



- Handle missing data in stationary **categorical** time series, X_1, \dots, X_n with $n \in \mathbb{N} = \{1, 2, \dots\}$.
- Outcomes x_t of X_t have qualitative range $\mathcal{S} = \{s_0, s_1, \dots, s_m\}$; consider both **nominal and ordinal** case.
- Unique approach of incorporating missingness and of deriving asymptotics of proposed statistics.
- Apply novel missing-data approaches to migraine time series.

Full paper: Weiß (2021) Analyzing categorical time series in the presence of missing observations.

Statistics in Medicine **40**(21), 4675–4690. (→ open access)



HELMUT SCHMIDT
UNIVERSITÄT
Universität der Bundeswehr Hamburg

**MATH
STAT**

Basics of Categorical Time Series Analysis

■

 ■
nominal vs. ordinal

Nominal range	Ordinal range
<p>Marginal PMF: (probability mass function)</p> $\mathbf{p} = (p_0, \dots, p_m)^\top \in [0; 1]^{m+1}$ <p>with $p_i = P(X = s_i)$</p>	<p>Marginal CDF: (cumulative distribution fct.)</p> $\mathbf{f} = (f_0, \dots, f_{m-1})^\top \in [0; 1]^m$ <p>with $f_i = P(X \leq s_i)$</p>
<p>Bivariate lag-h PMF:</p> $p_{ij}(h) = P(X_t = s_i, X_{t-h} = s_j)$	<p>Bivariate lag-h CDF:</p> $f_{ij}(h) = P(X_t \leq s_i, X_{t-h} \leq s_j)$
<p>Binarization $(\mathbf{Y}_t)_\mathbb{N}$ with $Y_{t,i} = \mathbb{1}_{\{X_t = s_i\}}$, so $E[\mathbf{Y}_t] = \mathbf{p}$</p>	<p>Binarization $(\mathbf{Z}_t)_\mathbb{N}$ with $Z_{t,i} = \mathbb{1}_{\{X_t \leq s_i\}}$, so $E[\mathbf{Z}_t] = \mathbf{f}$</p>
<p>Sample PMF:</p> $\hat{\mathbf{p}} = \frac{1}{T} \sum_{t=1}^T \mathbf{Y}_t$	<p>Sample CDF:</p> $\hat{\mathbf{f}} = \frac{1}{T} \sum_{t=1}^T \mathbf{Z}_t$
Bivariate (cumulative) relative frequencies	
$\hat{p}_{ij}(h) = \frac{1}{T-h} \sum_{t=h+1}^T Y_{t,i} Y_{t-h,j}$	$\hat{f}_{ij}(h) = \frac{1}{T-h} \sum_{t=h+1}^T Z_{t,i} Z_{t-h,j}$

Note that both **binarizations** lead to different range:

$$\mathbf{Y}_t \in \left\{ \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} \right\} =: \{\mathbf{e}_0, \dots, \mathbf{e}_m\} \subset [0; 1]^{m+1}$$

vs.

$$\mathbf{Z}_t \in \left\{ \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 1 \end{pmatrix}, \dots, \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix} \right\} =: \{\mathbf{c}_0, \dots, \mathbf{c}_m\} \subset [0; 1]^m.$$

Dispersion concepts for qualitative random variables:

- Minimal dispersion iff one-point distribution,
i. e., $\mathbf{p} \in \{\mathbf{e}_0, \dots, \mathbf{e}_m\}$ or $\mathbf{f} \in \{\mathbf{c}_0, \dots, \mathbf{c}_m\}$, respectively.
- Max. nominal dispersion iff uniform, $\mathbf{p} = (\frac{1}{m+1}, \dots, \frac{1}{m+1})^\top$.
- Max. ordinal disp. iff extreme two-point, $\mathbf{f} = (\frac{1}{2}, \dots, \frac{1}{2})^\top$.

Nominal range	Ordinal range
Index of qualitative variation: $IQV = \frac{m+1}{m} \left(1 - \sum_{i=0}^m p_i^2\right)$	Index of ordinal variation: $IOV = \frac{4}{m} \sum_{i=0}^{m-1} f_i(1 - f_i)$
	Ordinal skewness: $\text{skew} = \frac{2}{m} \sum_{i=0}^{m-1} f_i - 1$
Nominal Cohen's κ : $\kappa_{\text{nom}}(h) = \frac{\sum_{j=0}^m (p_{jj}(h) - p_j^2)}{1 - \sum_{i=0}^m p_i^2}$	Ordinal Cohen's κ : $\kappa_{\text{ord}}(h) = \frac{\sum_{j=0}^{m-1} (f_{jj}(h) - f_j^2)}{\sum_{i=0}^{m-1} f_i(1 - f_i)}$

The signed κ -measures serve as substitutes of ACF, where positive (negative) values express extend of (dis)agreement between X_t and X_{t-h} , see Weiß (2018, 2020).

Aforementioned approaches assume time series to be fully observed!

Now, categorical time series **with missing observations**, unified framework for handling missingness in *both ordinal and nominal* time series.

Time restrictions & clarity of talk:

focus on ordinal case, but nominal case in ...

Full paper: Weiß (2021) Analyzing categorical time series in the presence of missing observations.

Statistics in Medicine **40**(21), 4675–4690. (→ open access)

Ordinal Time Series with Missing Data

■

 ■
Amplitude Modulation

Let (X_t) be ordinal process, corresponding binarization (Z_t) .

Idea: Adapt amplitude modulation of Parzen (1963) to binarization (Z_t) .

Define **amplitude-modulating process** (O_t) as

$O_t = 1$ if X_t observed, and $O_t = 0$ otherwise.

Then, **amplitude modulation** of (Z_t) is $(O_t \cdot Z_t)$

(O_t) might be deterministic or i. i. d. with $E[O_0] = \pi$ or stationary with some dependence structure, $E[O_h O_0] = \pi(h)$.

But we assume that (O_t) is independent of (X_t) , (Z_t) .

Estimation of marginal CDF f :

Let $\hat{f} := \frac{1}{n} \sum_{t=1}^n O_t \mathbf{Z}_t$, then $E[\hat{f}] = \left(\frac{1}{n} \sum_{t=1}^n E[O_t] \right) f$.

Thus, estimate f by

$$\hat{f}^* := \frac{\frac{1}{n} \sum_{t=1}^n O_t \mathbf{Z}_t}{\frac{1}{n} \sum_{t=1}^n O_t} =: \hat{f} / \bar{O}.$$

Estimators for IOV and skew by using \hat{f}^* :

$$\widehat{\text{IOV}} = \frac{4}{m} \sum_{i=0}^{m-1} \hat{f}_i^* (1 - \hat{f}_i^*), \quad \widehat{\text{skew}} = \frac{2}{m} \sum_{i=0}^{m-1} \hat{f}_i^* - 1.$$

Estimation of bivariate CDF $f_{ij}(h)$:

Let $\hat{f}_{ij}(h) = \frac{1}{n-h} \sum_{t=h+1}^n O_t O_{t-h} Z_{t,i} Z_{t-h,j}$,

then $E[\hat{f}_{ij}(h)] = \left(\frac{1}{n-h} \sum_{t=h+1}^n E[O_t O_{t-h}] \right) f_{ij}(h)$.

Thus, estimate $f_{ij}(h)$ by

$$\hat{f}_{ij}^*(h) = \frac{\frac{1}{n-h} \sum_{t=h+1}^n O_t O_{t-h} Z_{t,i} Z_{t-h,j}}{\frac{1}{n-h} \sum_{t=h+1}^n O_t O_{t-h}} =: \hat{f}_{ij}(h) / \overline{O_t O_{t-h}}.$$

Estimator for ordinal κ by using $\hat{f}^*, \hat{f}_{ij}^*(h)$:

$$\hat{\kappa}_{\text{ord}}(h) = \frac{\sum_{j=0}^{m-1} (\hat{f}_{jj}^*(h) - (\hat{f}_j^*)^2)}{\sum_{i=0}^{m-1} \hat{f}_i^* (1 - \hat{f}_i^*)} \quad \text{for } h \in \mathbb{N}.$$



HELMUT SCHMIDT
UNIVERSITÄT
Universität der Bundeswehr Hamburg

**MATH
STAT**

Ordinal Time Series with Missing Data

■

 ■
Asymptotics

Theorem 1: Let (X_t) and (O_t) be α -mixing with exponentially decreasing weights.

Then, the corrected CDF-estimator \hat{f}^* satisfies

$$\sqrt{n} \left(\hat{f}^* - f \right) \xrightarrow{d} N(\mathbf{0}, \Sigma^*), \quad \text{with } \Sigma^* = (\sigma_{ij}^*)_{i,j=0,\dots,m-1}, \quad \text{where}$$
$$\sigma_{ij}^* = \frac{1}{\pi} \left(f_{\min\{i,j\}} - f_i f_j \right) + \frac{1}{\pi^2} \sum_{h=1}^{\infty} \pi(h) \left(f_{ij}(h) + f_{ji}(h) - 2 f_i f_j \right).$$

Furthermore, the bias $E[\hat{f}_j^*] - f_j$ is of order $o(n^{-1})$.

Theorem 2: Let assumptions of Theorem 1 hold. Then, corrected sample IOV and skew asymptotically normal with

$$E[\widehat{\text{IOV}}] \approx \text{IOV} \left(1 - \frac{1}{n} \left(\frac{1}{\pi} + \frac{2}{\pi^2} \sum_{h=1}^{\infty} \pi(h) \kappa_{\text{ord}}(h) \right) \right),$$

$$V[\widehat{\text{IOV}}] \approx \frac{1}{n} \frac{16}{m^2} \sum_{i,j=0}^{m-1} (1 - 2f_i)(1 - 2f_j) \sigma_{ij}^*,$$

and

$$E[\widehat{\text{skew}}] \approx \text{skew} + o(n^{-1}), \quad V[\widehat{\text{skew}}] \approx \frac{1}{n} \frac{4}{m^2} \sum_{i,j=0}^{m-1} \sigma_{ij}^*.$$

Theorem 3: Let assumptions of Theorem 1 hold. Denote

$$\pi(h_1, \dots, h_r) := E[O_0 \cdot O_{h_1} \cdots O_{h_r}] \quad \text{with } 0 < h_1 < \dots < h_r.$$

Under null hypothesis of (X_t) being i. i. d.,
 distribution of $\hat{\kappa}_{\text{ord}}(h)$ at lag $h \in \mathbb{N}$ approximately normal

with mean $-\frac{1}{n\pi}$

and variance

$$\frac{1}{n} \frac{\sum_{i,j=0}^{m-1} \left(f_{\min\{i,j\}} - f_i f_j \right) \left(f_{\min\{i,j\}} - f_i f_j + 2 \left(1 + \frac{\pi(h, 2h)}{\pi(h)} - 2 \frac{\pi(h)}{\pi} \right) f_i f_j \right)}{\pi(h) \left(\sum_{k=0}^{m-1} f_k (1 - f_k) \right)^2}.$$

Possible applications: confidence intervals, hypothesis tests.

Example: Theorem 3 to test for serial independence at lag h .
 Let (O_t) be i. i. d. (“missing at random”), then simplification

$$V[\hat{\kappa}_{\text{ord}}(h)] \approx \frac{1}{n} \frac{1}{\pi^2} \frac{\sum_{i,j=0}^{m-1} (f_{\min\{i,j\}} - f_i f_j)^2}{\left(\sum_{k=0}^{m-1} f_k(1 - f_k)\right)^2} + \frac{2}{n} \frac{1 - \pi}{\pi^2} \frac{\sum_{i,j=0}^{m-1} f_i f_j (f_{\min\{i,j\}} - f_i f_j)}{\left(\sum_{k=0}^{m-1} f_k(1 - f_k)\right)^2}.$$

Plug-in estimated probabilities $\hat{\pi}, \hat{\mathbf{f}}^*$.

Then **critical values** $-1/(n\hat{\pi}) \mp z_{1-\alpha/2} \hat{\sigma}_{\kappa}$,

where $z_{1-\alpha/2}$ denotes $(1 - \alpha/2)$ -quantile of $N(0, 1)$.

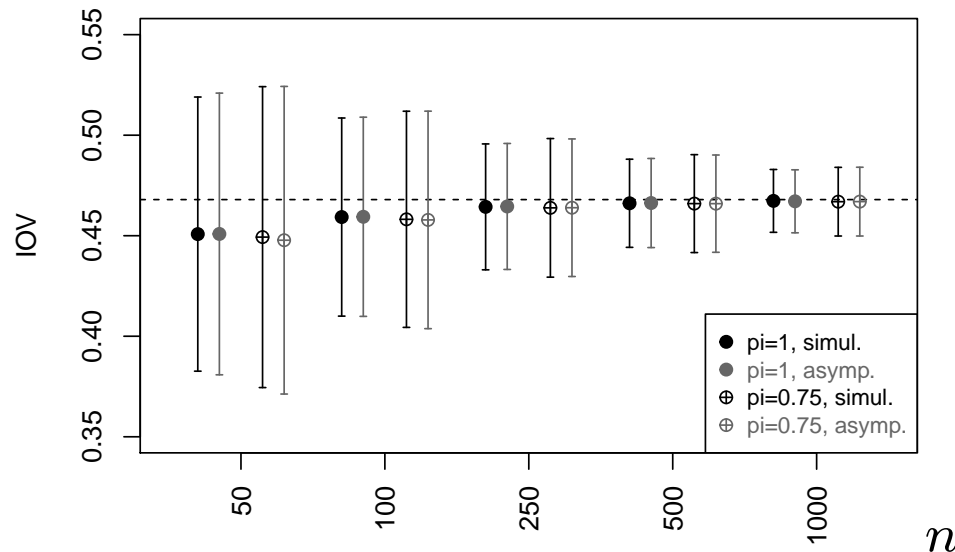
Empirical Illustrations

- Simulations & Data Example ■

Simulation study in full paper, Weiß (2021), confirms good finite-sample performance of asymptotic approximations.

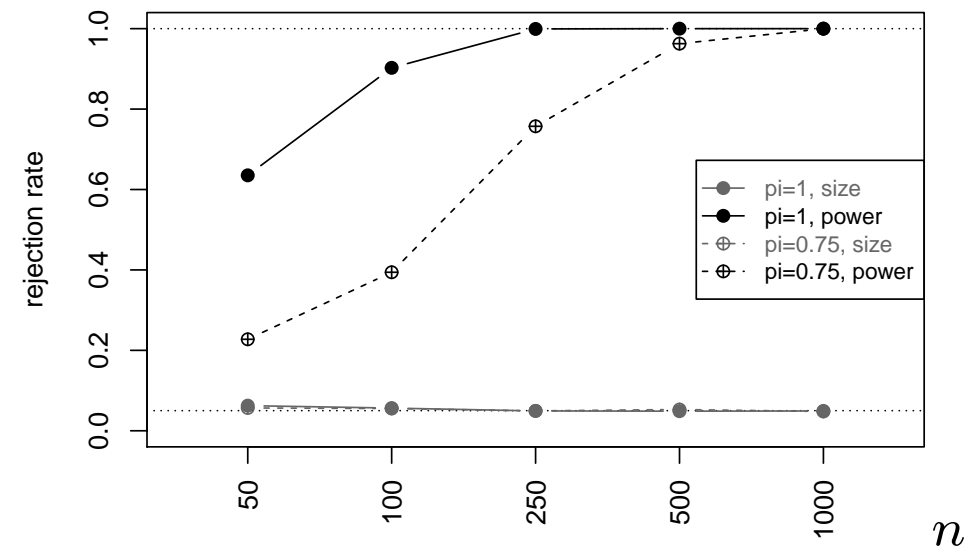
Illustrative examples for models . . .

$$(m, p, r) = (3, 0.20, 0.35)$$



Mean \mp std. dev. for IOV

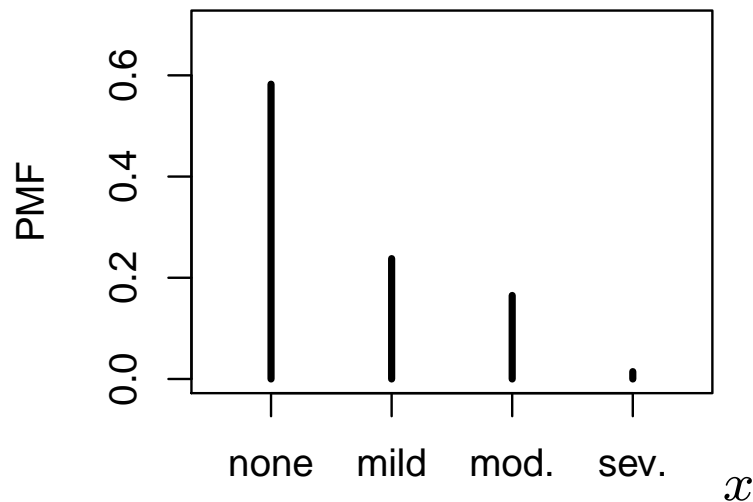
$$(m, p, r) = (3, 0.20, 0 \text{ vs. } 0.35)$$



Rej. rate based on $\hat{\kappa}_{\text{ord}}(1)$

Data application: migraine patients

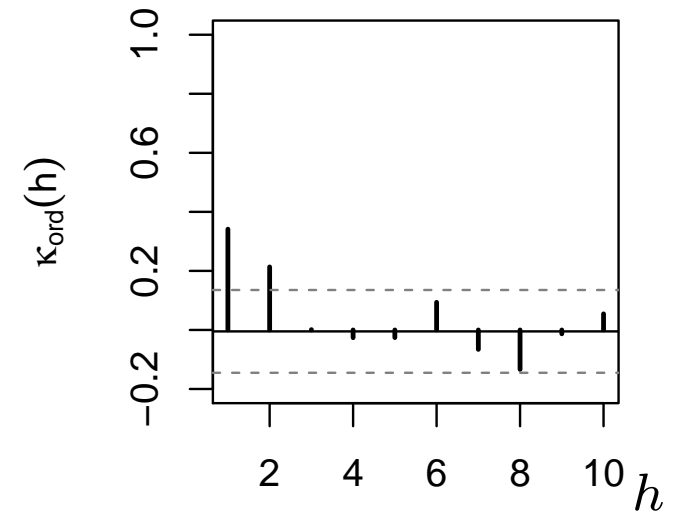
Estimates and tests *accounting for missingness* for **peak-severity** time series of patient A ($n = 225$, 19 missing):



Median:
"none"

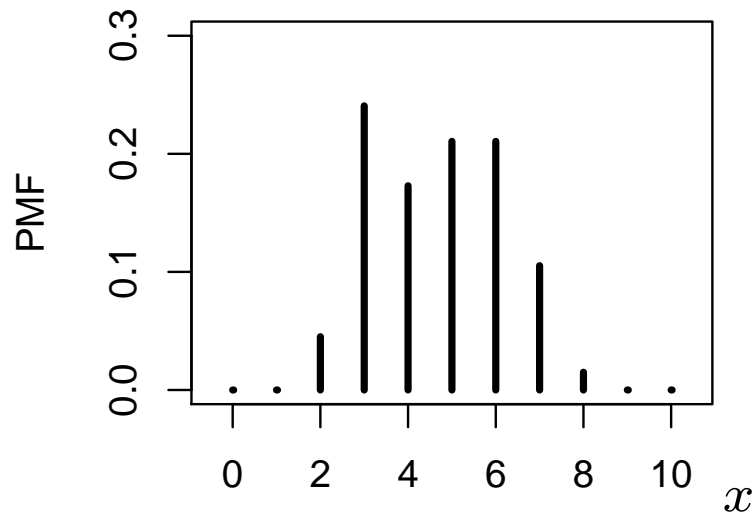
\widehat{IOV}
 ≈ 0.540

\widehat{skew}
 ≈ 0.592



Data application: migraine patients

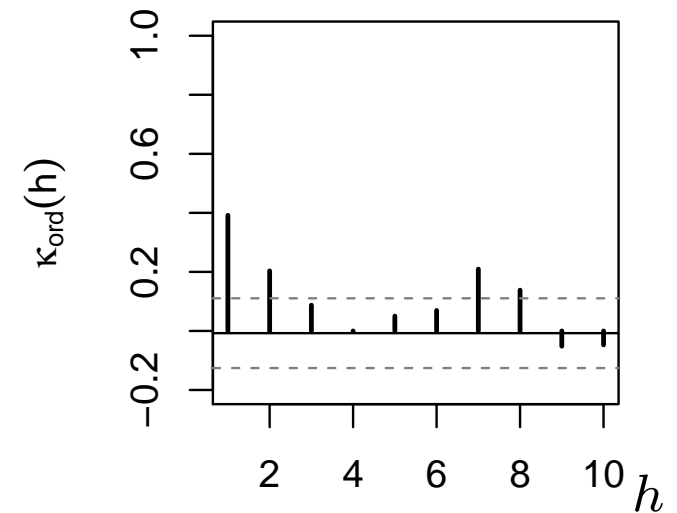
Estimates and tests *accounting for missingness* for **stress** time series of patient B ($n = 136$, 3 missing):



Median:
5

\widehat{IOV}
 ≈ 0.335

$\widehat{\text{skew}}$
 ≈ 0.065



What would happen if just ignore (skip) missing data?

Then complete time series of reduced lengths

$\tilde{n} = 206$ (peak severity) and $\tilde{n} = 133$ (stress).

$\hat{\kappa}_{\text{ord}}(h)$ changes from 0.341, 0.213, ... to 0.258, 0.138, ...
(approximate SE from 0.072 to 0.055)

for peak-severity series (19 out of 225),

and from 0.392, 0.203, ... to 0.370, 0.195, ...
(approximate SE from 0.060 to 0.057)

for stress series (3 out of 136).

⇒ **carefully consider missingness** for time series analysis!

**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`

Weiß (2021) Analyzing categorical time series in the presence of missing observations. *Statistics in Medicine* 40(21), 4675–4690.

(→ open access)

Bloomfield (1970) Spectral analysis with ... *JRSS B* **32**, 369–380.

Dunsmuir & Robinson (1981) Asymptotic theo... *Sankhyā A* **43**, 260–281.

Neave (1970) Spectral analysis of ... *Biometrika* **57**, 111–122.

Parzen (1963) On spectral analysis with ... *Sankhyā A* **25**, 383–392.

Scheinok (1965) Spectral analysis with ... *AMS* **36**, 971–977.

Vives-Mestres & Casanova (2021) Modelling ... *Stat Med* **40**, 213–225.

Weiß (2018) *An Introduction to Discrete-Valued Time Series*. Wiley.

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

Yajima & Nishino (1999) Estimation of the ... *Sankhyā A* **61**, 189–207.