

# Categorical Time Series: Analysis, Modelling, Monitoring?



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

**Christian H. Weiß**

Department of Mathematics,

Darmstadt University of Technology



*Fachbereich*  
**Mathematik**



# Categorical Time Series

■ 

---

 ■  
Motivation



## Log data of a web server:

Several types of categorical features (IP addresses, web addresses, state codes, etc.).

(→ ENBIS talk in 2006)

Possible applications:

- Measuring success of a web site,
- intrusion detection.



## Medical Diagnoses

of an examiner for a certain type of examination.

(→ ENBIS talk in 2009)

Possible applications:

- Compare diagnosis behavior of different examiners,
- compare examiner with a norm profile.



## **Text:**

Sequence of letters, words, word classes, grammatical entities, etc.

Possible applications:

- Language recognition,
- text recognition,
- part-of-speech tagging, etc.



## Genetic or protein sequences.

Possible applications:

- Structure analysis to understand functionality of segments of the sequence,
- recognition of similar sequences,
- recognition of common evolutionary roots,
- identification of sequences through comparison to consensus model.



## SPC-related examples.

$X_t$  = result of **inspection of item**, with

$X_t = i$  for  $i = 1, \dots, m$  iff item has nonconformity type  $i$ ,

$X_t = 0$  iff conforming.

**Mukhopadhyay (2008)**:  $m = 6$  paint defects of ceiling fan cover ('poor covering', 'bubbles', etc.).

Overall defect category = most predominant defect.

**Ye et al. (2002)** monitor network traffic data (284 different types of audit events) for intrusion detection.



# Categorical Random Variables

---

Properties





## Categorical random variable:

$X$  takes one of **finite** number of **unordered** categories, say  $b_0, \dots, b_m$ . E. g.,

- diagnoses  $b_0, \dots, b_m$ , or
- parts of speech  $b_0, \dots, b_m$ , or
- types of defects  $b_0, \dots, b_m$ , or ...

To simplify notations:

Range of  $X$  is coded as  $\mathcal{V} = \{0, \dots, m\}$ .



$X$  categorical r.v. with range  $\mathcal{V} = \{0, \dots, m\}$

$$\Rightarrow P(X = 0) = 1 - \sum_{j=1}^m P(X = j).$$

Abbreviate marginal probabilities:  $p_i := P(X = i) \in (0; 1)$ ,  
whole distribution determined by  $m$  parameters  $p_1, \dots, p_m$ .

**1<sup>st</sup> problem:** Number  $m$  of parameters might be very large,  
in contrast to real-valued or count-data random variables,  
**no sparse parametric models** available yet!

... except trivial cases:

one-point distribution, uniform distribution.



Typical for **cardinal** random variables:  
quantify basic properties, e. g.,

- location via mean or median,
- dispersion via variance or quartiles.

**But:**

- no arithmetic operations for categorical range  
⇒ no mean, variance, etc.
- unordered range ⇒ no quantiles.



**So how can we quantify  
location and dispersion  
of a categorical random variable?**

**Measures of location:**

only **mode** of  $X$  in use, i. e.,

value  $i \in \mathcal{V}$  such that  $p_i \geq p_j$  for all  $j \in \mathcal{V}$ .

Often not uniquely determined (e. g., uniform distribution).



Intuitive understanding of dispersion:

$X$  shows large dispersion

$\approx$

High uncertainty about the outcome of  $X$

$\Rightarrow$  Uncertainty of categorical random variable?



Two extreme cases:

**Uniform distribution:**

Maximal uncertainty about the outcome of  $X$ .

**One-point distribution:**

Perfect certainty about the outcome of  $X$ .

⇒ Hallmarks for definition of any measure of dispersion!

Contributions in literature (desirable properties, measures),  
e. g., by Uschner (1987), Vogel & Kiesel (1999).



Common standardized measures of dispersion:

**Gini index:**  $\nu_G(X) := \frac{m+1}{m} \left(1 - \sum_j p_j^2\right).$

**Entropy:**  $\nu_E(X) := -\frac{1}{\ln(m+1)} \sum_j p_j \ln p_j.$

**Chebycheff dispersion:**  $\nu_C(X) := \frac{m+1}{m} \left(1 - \max_j p_j\right).$



**Important properties** of these measures:

- continuous and symmetric functions of  $p$ ,
- range  $[0; 1]$ ,
- maximum value 1 in case of uniform distribution,
- minimal value 0 in case of one-point distribution,
- inequality:  $\frac{m+1}{m} (1 - \min_j p_j) \geq \nu_G(X) \geq \nu_C(X)$ .





## Empirical measures of dispersion

based on sample  $X_1, \dots, X_T$ .

**Binarization**  $\mathbf{Y}_t \in \{0, 1\}^{m+1}$  of  $X_t$  via  $Y_{t,i} := \delta_{i, X_t}$ .

Unbiased estimator for  $p_i$ :  $\hat{p}_i(T) := \frac{1}{T} \cdot \sum_{t=1}^T Y_{t,i}$ .

Abbreviate  $\mathbf{p} := (p_0, \dots, p_m)^\top$ ,

$\hat{\mathbf{p}}(T) := (\hat{p}_0(T), \dots, \hat{p}_m(T))^\top = \frac{1}{T} \cdot \sum_t \mathbf{Y}_t$ ,

and  $s_k(\mathbf{p}) := \sum_j p_j^k$  for  $k \in \mathbb{N}$ .



Wei (2011): **Empirical Gini index** defined by

$$\hat{\nu}_G(X) := \frac{m+1}{m} \cdot \frac{T}{T-1} \cdot \left(1 - s_2(\hat{\mathbf{p}}(T))\right).$$

If computed from i.i.d. data, then

$\hat{\nu}_G(X)$  is exactly unbiased and asymptotically normally distributed with variance determined by

$$V\left[1 - s_2(\hat{\mathbf{p}}(T))\right] = \frac{4}{T} \cdot \left(s_3(\mathbf{p}) - s_2^2(\mathbf{p})\right) + O(T^{-2}).$$

**Current research:**

$\hat{\nu}_G(X)$  and  $\hat{\nu}_E(X)$  for NDARMA processes (see below).



# Categorical Time Series

---

Terms & Notations



## Categorical process:

$(X_t)_{\mathbb{N}}$  with  $\mathbb{N} = \{1, 2, \dots\}$ , where each  $X_t$  takes one of **finite** number of **unordered** categories.

## Categorical time series:

Realizations  $(x_t)_{t=1, \dots, T}$  from  $(X_t)_{\mathbb{N}}$ .

$(X_t)_{\mathbb{N}}$  said to be **stationary**

if joint distribution of  $(X_t, \dots, X_{t+k})$

independent of  $t$  for all  $k \in \mathbb{N}_0$ .



**Notations** for time-invariant probabilities:

If  $(X_t)_{\mathbb{N}}$  stationary:

- marginal probabilities  $p_i := P(X_t = i) \in (0; 1)$ .  
 $\mathbf{p} := (p_0, \dots, p_m)^\top$ , and  
 $s_k(\mathbf{p}) := \sum_j p_j^k$  for  $k \in \mathbb{N}$ ; obviously  $s_1(\mathbf{p}) = 1$ .
- bivariate probabilities  $p_{ij}(k) := P(X_t = i, X_{t-k} = j)$ ,  
conditional probabilities  $p_{i|j}(k) := P(X_t = i \mid X_{t-k} = j)$ .



# Categorical Time Series



Visual Analysis



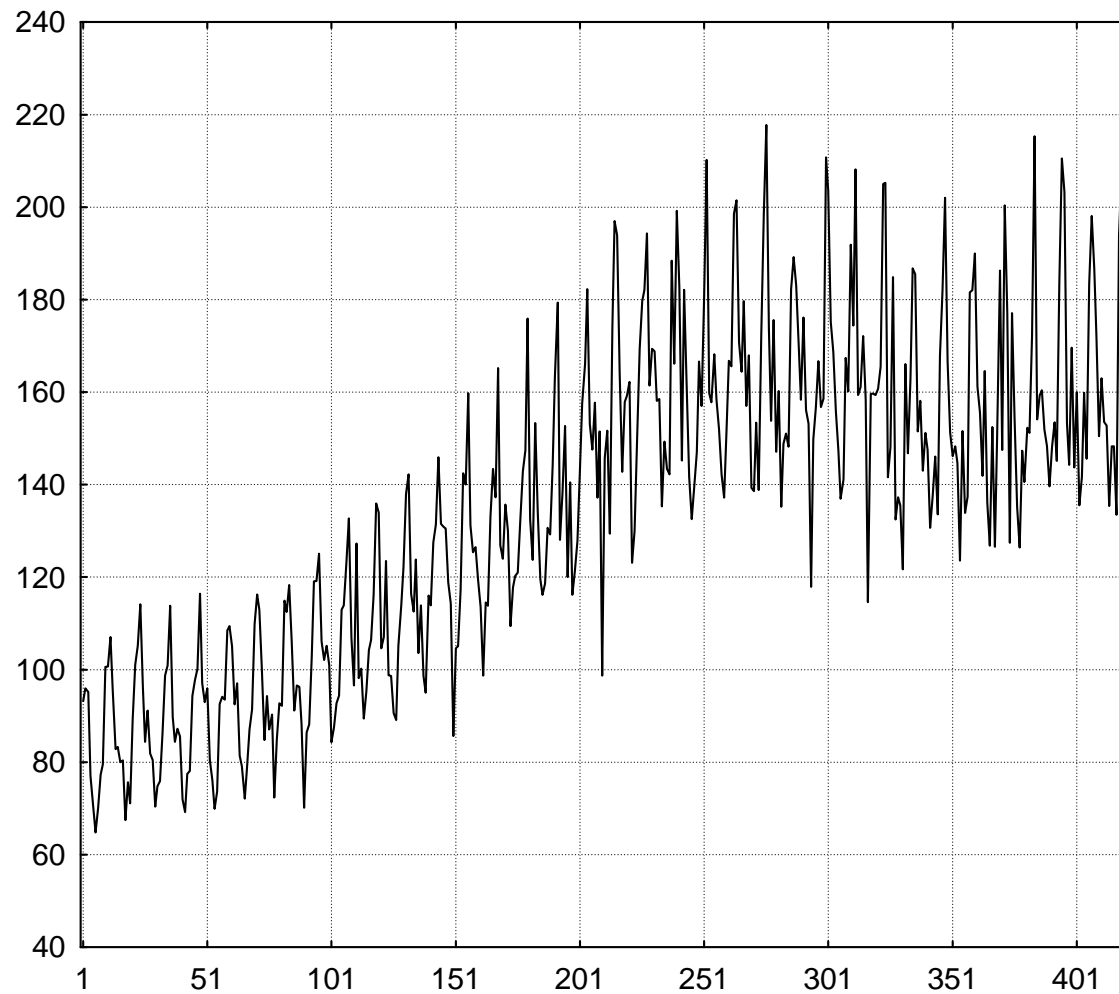
*“Little existing work deals directly with categorical time series analysis, and much less deals with the visualization of categorical time series.”* (Ribler, 1997, p. 11)

Several visual tools for real-valued time series:  
line plot, periodogram, correlogram, etc.

At least line plot simple and universal tool!



# Categorical Time Series — Visual Analysis







## **Visual tools for categorical time series:**

only few proposals,

often from computer science and biology,

see survey by Weiß (2008).

In fact problematic: Analogue of line plot?

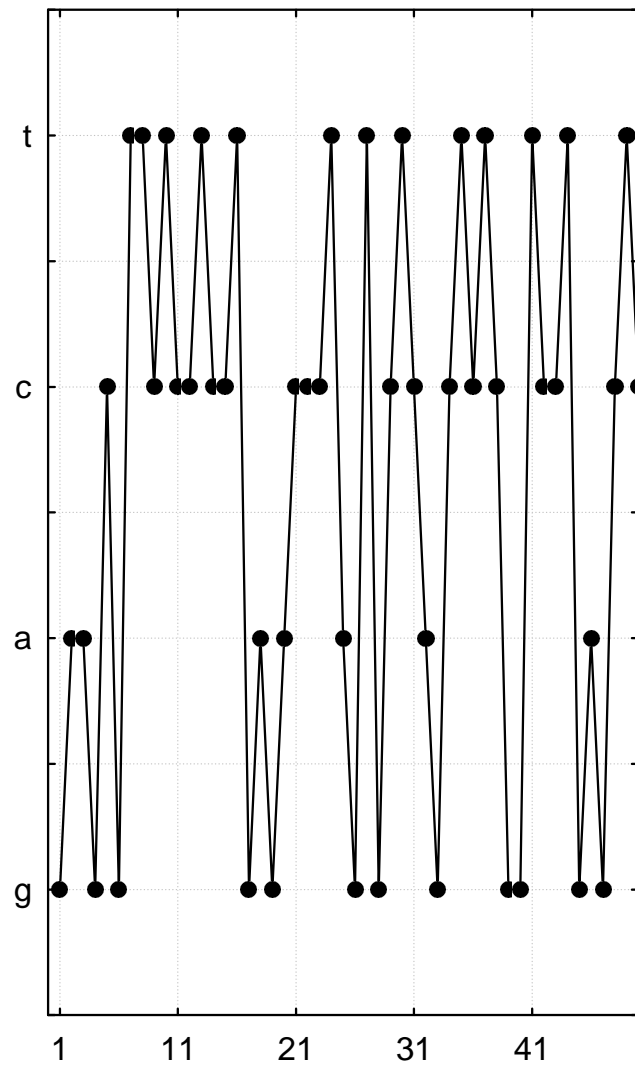
Lack of a natural order within purely categorical range

⇒ arrangement of range along ordinate

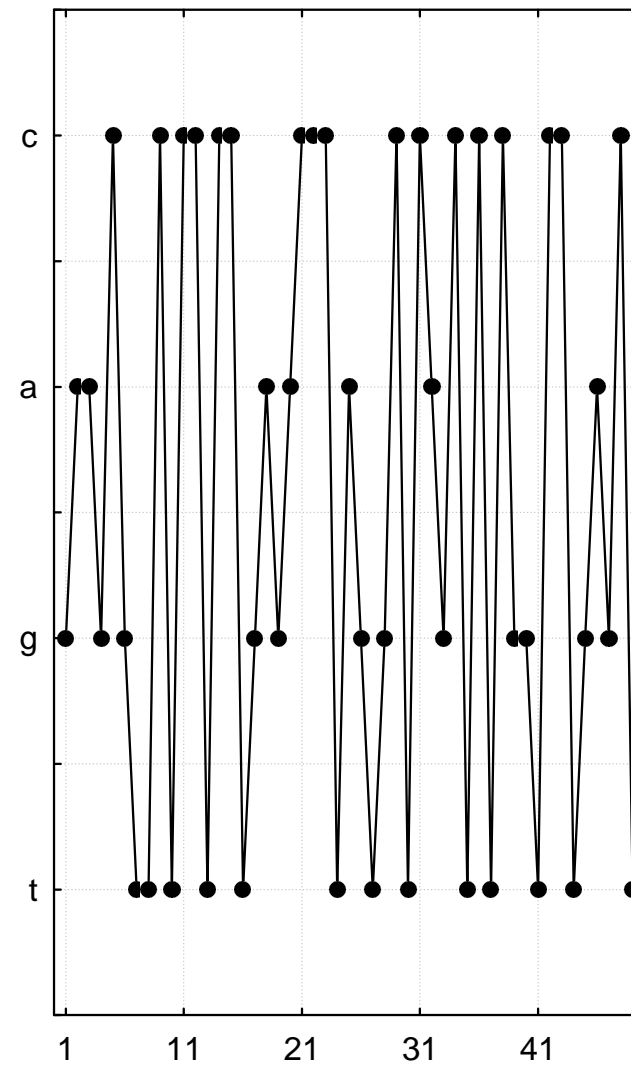
arbitrary and misleading.



# Categorical Time Series — Visual Analysis



or



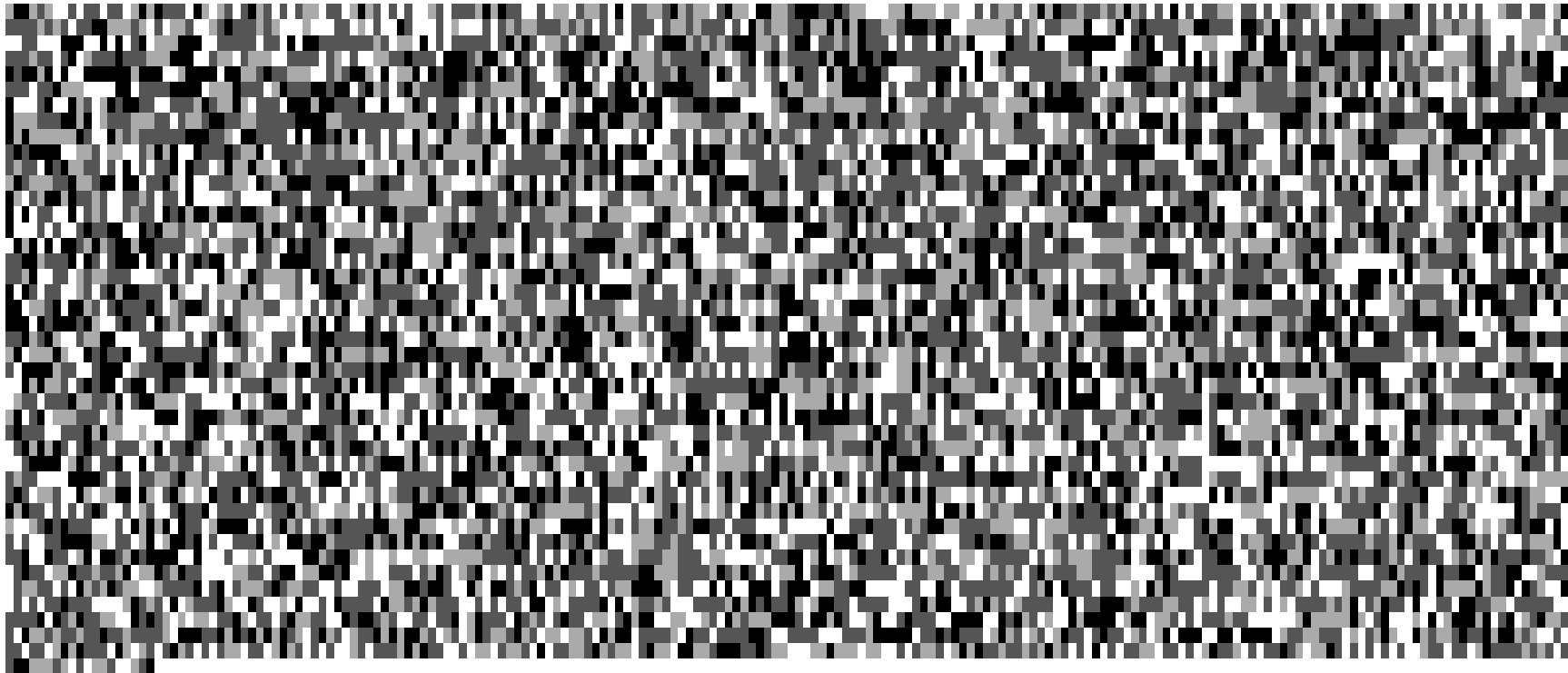


Alternative (Keim & Kriegel, 1996): Map range onto set of colors or symbols, plot  $x_1, x_2, \dots$  successively on space-filling curve (line-by-line, column-by-column, Peano-Hilbert curves, spiral, etc.).

However: These graphs perform poorly, characteristic features of time series difficult to recognize, interpretation of resulting plot is problematic, appearance depends heavily on choice of underlying curve.



Genome of Bovine leukemia virus:



■ a ■ c ■ g □ t



So how to analyze categorical time series visually?

Proposal by Ribler (1997): **Rate evolution graph.**

Categorical process  $(X_t)_{\mathbb{N}}$ , range coded as  $\mathcal{V} = \{0, \dots, m\}$ .

Binarization  $\mathbf{Y}_t \in \{0, 1\}^{m+1}$  with  $Y_{t,i} = \delta_{i,X_t}$ ,  $i = 0, \dots, m$ .

Define the cumulated sums  $\mathbf{C}_t := \sum_{s=1}^t \mathbf{Y}_s$ , i. e.,

$C_{t,i}$  = number of  $X_s$ ,  $s = 1, \dots, t$ , equal to  $i$ .



**Rate evolution graph** of  $(X_t)_{\mathbb{N}}$ : (Ribler, 1997)

Multiple line plot of all component series  $C_{t,i}$ ,  $i = 0, \dots, m$ ,  
i. e., all  $C_{t,i}$  are plotted simultaneously into one chart.

## **Interpretation:**

Slope of graphs is estimate for corresp. marginal probability.

If  $(X_t)_{\mathbb{N}}$  stationary and at most moderately serially dependent, then graphs approximately linear in  $t$

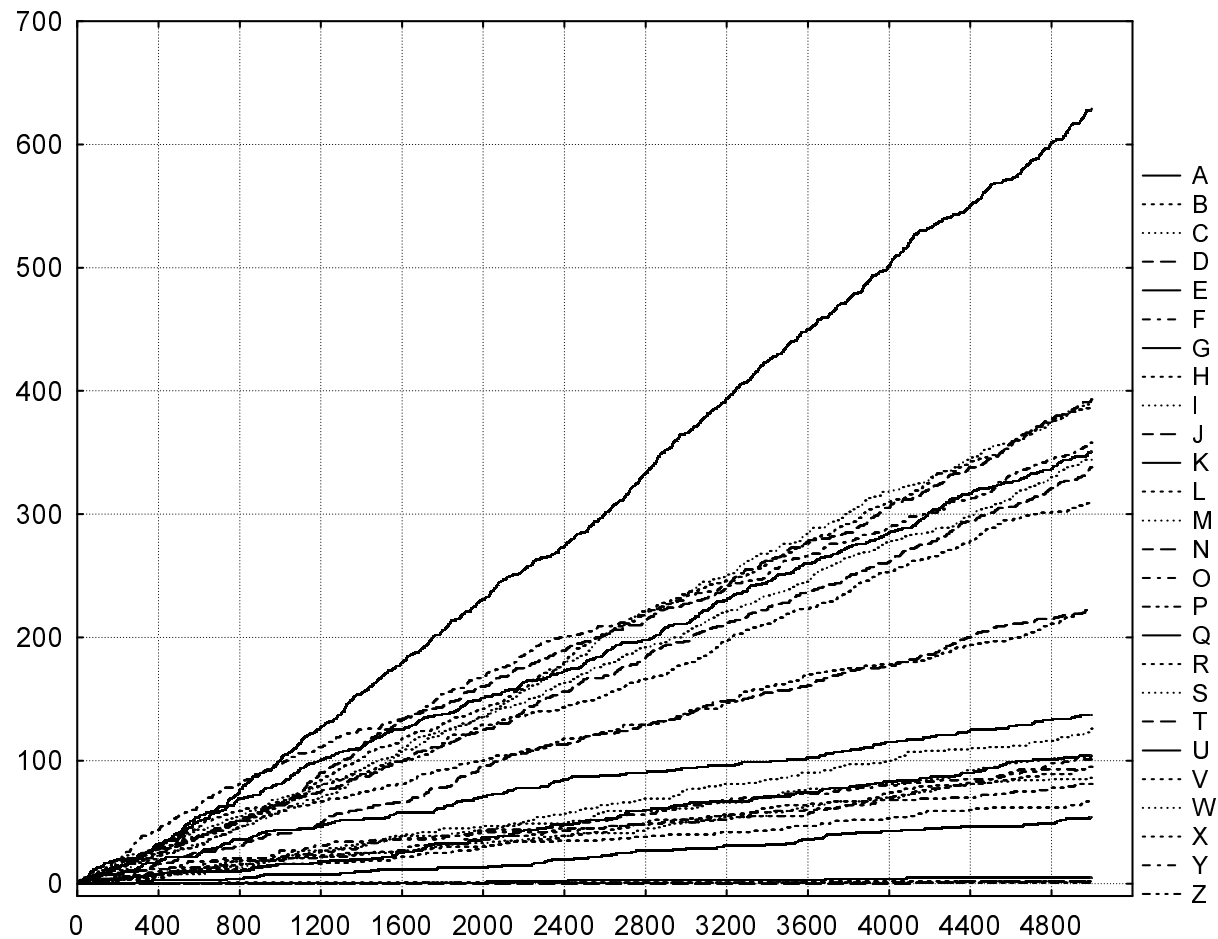
⇒ Simple visual tool for checking stationarity.



# Categorical Time Series — Visual Analysis



Shakespeare's (1593) poem "Venus and Adonis":

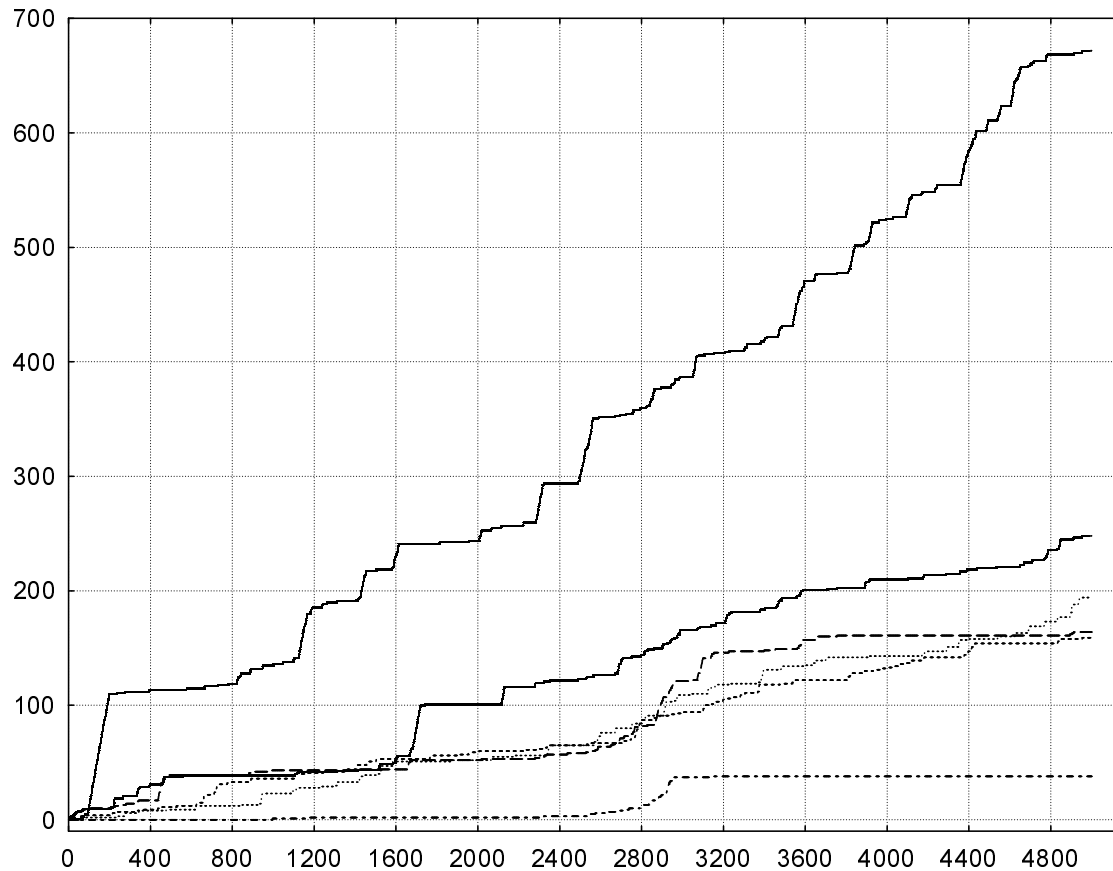




# Categorical Time Series — Visual Analysis



Log data (2005) of Statistics server at Univ. of Würzburg:  
Access to home directory of five members.





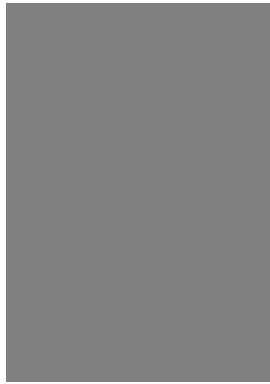


**Some further tools:** (see Weiß (2008) for references)

- **VizTree** or **IFS circle transformation** for visualizing string frequencies; (→ ENBIS talk in 2005)
- **pattern histograms** (e. g., based on runs or cycles);
- **categorical control charts** (also see below).

**However:** These tools are very specialized.

Universal instrument (like line plot), providing multiple types of information at once, still missing.



# Categorical Time Series

- ---

 Serial Dependence ▪



## **Cardinal time series:**

Convenient measure of serial dependence:  
(Partial) Autocorrelation.

## **Categorical time series:**

Measures of serial dependence?



Weiß & Göb (2008): same strategy as for dispersion measures, i. e., we start with

## Extreme cases:

- $X_t, X_{t-k}$  (stochastically) independent iff  $p_{ij}(k) = p_i \cdot p_j$ .
- $X_t$  perfectly depends on  $X_{t-k}$  iff for every  $j = 0, \dots, m$ : conditional distribution of  $X_t$ , conditioned on  $X_{t-k} = j$ , is a one-point distribution.



Weiß & Göb (2008) proposed, among others,

Goodman and Kruskal's  $\tau$ : 
$$\tau(k) = \sqrt{\sum_{i,j} \frac{(p_{ij}(k) - p_i p_j)^2}{p_j (1 - s_2(\mathbf{p}))}}$$

Cramer's  $v$ : 
$$v(k) = \sqrt{\frac{1}{m} \cdot \sum_{i,j} \frac{(p_{ij}(k) - p_i p_j)^2}{p_i p_j}}$$

**Some properties:**  $\tau(k), v(k)$  have range  $[0; 1]$  with

- $X_t, X_{t-k}$  independent  $\Leftrightarrow \tau(k) = v(k) = 0$ .
- $X_t$  depends perfectly on  $X_{t-k}$   $\Leftrightarrow \tau(k) = v(k) = 1$ .



**Cardinal case:** Positive and negative autocorrelation.

⇒

Concept of **signed dependence:** (Weiß & Göb, 2008)

- $X_t, X_{t-k}$  **perfectly positively dependent**,  
if perfectly dependent and  
if  $p_{i|i}(k) = 1$  for all  $i$ .
- $X_t, X_{t-k}$  **perfectly negatively dependent**  
if perfectly dependent and  
if  $p_{i|i}(k) = 0$  for all  $i$ .



## Measures of signed dependence:

(Weiß, 2011; Weiß & Göb, 2008)

Cohen's  $\kappa$ :

$$\kappa(k) = \frac{\sum_j p_{jj}(k) - s_2(\mathbf{p})}{1 - s_2(\mathbf{p})} \text{ with range } \left[-\frac{s_2(\mathbf{p})}{1 - s_2(\mathbf{p})}; 1\right],$$

Modified  $\kappa$ :

$$\kappa^*(k) = \frac{1}{m} \cdot \left(\sum_j p_{j|j}(k) - 1\right) \text{ with range } \left[-\frac{1}{m}; 1\right].$$

## Some properties:

- $X_t, X_{t-k}$  independent  $\Rightarrow \kappa(k) = \kappa^*(k) = 0$ .
- $X_t, X_{t-k}$  perf. positively dep.  $\Leftrightarrow \kappa(k) = \kappa^*(k) = 1$ .
- $X_t, X_{t-k}$  perf. negatively dep.  $\Rightarrow \kappa(k), \kappa^*(k)$  minimal.



## Empirical measures of (signed) serial dependence

based on time series  $X_1, \dots, X_T$ .

Weiß (2011): **empirical Cohen's  $\kappa$**  and **modified  $\kappa$** :

$$\hat{\kappa}(k) := 1 + \frac{1}{T} - \frac{1 - \sum_j \hat{p}_{jj}(k, T)}{1 - s_2(\hat{\mathbf{p}}(T))},$$

$$\hat{\kappa}^*(k) := \frac{1}{T} + \frac{1}{m} \cdot \left( \sum_j \frac{\hat{p}_{jj}(k, T)}{\hat{p}_j(T)} - 1 \right).$$





## Empirical measures of (signed) serial dependence

If computed from i.i.d. data, then

$\hat{\kappa}(k)$ ,  $\hat{\kappa}^*(k)$  are asymptotically normally distributed with

$$E[\hat{\kappa}(k)] = E[\hat{\kappa}^*(k)] = 0 + O(T^{-2}),$$

and variances given by

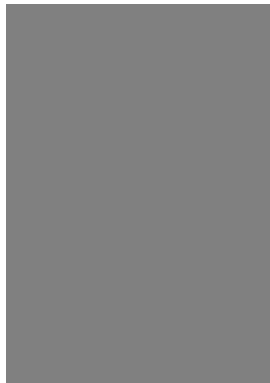
$$V[\hat{\kappa}(k)] = \frac{1}{T} \cdot \left( 1 - \frac{1 + 2s_3(\mathbf{p}) - 3s_2(\mathbf{p})}{(1 - s_2(\mathbf{p}))^2} \right) + O(T^{-2});$$

$$V[\hat{\kappa}^*(k)] = \frac{1}{m \cdot T} + O(T^{-2}).$$



## Current and future research:

- $\hat{\kappa}(k)$ ,  $\hat{\kappa}^*(k)$  and also  $\hat{\tau}(k)$ ,  $\hat{v}(k)$  for NDARMA processes;
- goodness-of-fit tests for NDARMA processes.



# Categorical Time Series



Modelling



## Models for stationary categorical processes

with range  $\mathcal{V} = \{0, \dots, m\}$ :

- **$p^{\text{th}}$  order Markov model**:  $(m + 1)^p \cdot m$  parameters;
- **variable length M. m.** of Bühlmann & Wyner (1999):  
more parsimonious, but model choice difficult;
- **MTD( $p$ ) model** of Raftery (1985):  
still  $m(m + 1) + p - 1$  parameters;
- **NDARMA( $p, q$ ) models** of Jacobs & Lewis (1983):  
 $m + p + q$  parameters, also non-Markovian dependence.



$(X_t)_{\mathbb{Z}}$ ,  $(\epsilon_t)_{\mathbb{Z}}$ : categorical processes with range  $\mathcal{V} = \{0, \dots, m\}$ ;  
 $(\epsilon_t)_{\mathbb{Z}}$ : i.i.d. with marginal distribution  $P(\epsilon_t = j) = p_j > 0$ ,  
 $\epsilon_t$  independent of  $(X_s)_{s < t}$ .

For  $\varphi_q > 0$ , with  $\phi_p > 0$  if  $p \geq 1$ , let

$$D_t = (\alpha_{t,1}, \dots, \alpha_{t,p}, \beta_{t,0}, \dots, \beta_{t,q}) \sim \text{MULT}(1; \phi_1, \dots, \phi_p, \varphi_0, \dots, \varphi_q)$$

be i.i.d. and independent of  $(\epsilon_t)_{\mathbb{Z}}$ ,  $(X_s)_{s < t}$ .

$(X_t)_{\mathbb{Z}}$  is **NDARMA(p, q) process** if

$$X_t = \alpha_{t,1} \cdot X_{t-1} + \dots + \alpha_{t,p} \cdot X_{t-p} + \beta_{t,0} \cdot \epsilon_t + \dots + \beta_{t,q} \cdot \epsilon_{t-q}.$$



## Some properties:

- marginal distribution  $P(X_t = j) = p_j$ ;
- only shows positive serial dependence, and
$$\kappa(k) = \kappa^*(k) = v(k) = \tau(k);$$

- Yule-Walker-type equations

$$\kappa(k) = \sum_{j=1}^p \phi_j \cdot \kappa(|k-j|) + \sum_{i=0}^{q-k} \varphi_{i+k} \cdot r(i) \quad \text{for } k \geq 1,$$

where  $r(i) = 0$  for  $i < 0$ ,  $r(0) = \varphi_0$ , and

$$r(i) = \sum_{j=\max\{0, i-p\}}^{i-1} \phi_{i-j} \cdot r(j) + \sum_{j=0}^q \delta_{ij} \cdot \varphi_j \quad \text{for } i > 0.$$



## **Disadvantage of NDARMA models:**

Only describe positive dependence, i. e.,  
categorical time series with long runs of symbols.  
(e. g., genetic sequence vs. letters of text)

## **Issues for future research:**

- Find simple and sparsely parametrized models that also allow for negative dependence!
- Definition of trend or seasonality?
- And: How to remove such trend or seasonality?



# Categorical Processes



Monitoring





Only very few approaches for monitoring processes of **unordered** and **mutually exclusive** categories.

Monitoring **samples** from an **i.i.d.** categorical process:

Duncan (1950), Marcucci (1985), Nelson (1987), Mukhopadhyay (2008) plot Pearson's  $\chi^2$ -statistic for goodness of fit on a control chart.

**Continuously** monitoring cat. process (100 % inspect.):

Weiß (2010) proposes

- two moving-average-type charts,  
based either on Pearson statistic or Gini index,
- a  $(k, r)$ -runs chart.





If  $(X_t)_{\mathbb{N}}$  Markov chain, then

$(Y_n^{(k,r)})_{\mathbb{N}}$  i.i.d. process, range  $\mathbb{N}_k := \{k, k + 1, \dots\}$ .

**Properties:** (Chryssaphinou et al., 1994)

Let  $c_{k,r}(z) := \sum_{i=1}^r \frac{(1 - p_i z)(p_i z)^k}{1 - (p_i z)^k}$ , then

$$E[Y^{(k,r)}] = \frac{1}{c_{k,r}(1)}, \quad V[Y^{(k,r)}] = \frac{1 + c_{k,r}(1) - 2c'_{k,r}(1)}{c_{k,r}^2(1)}.$$



**$(k, r)$ -runs chart:**

(Weiß, 2010)

$(Y_n^{(k,r)})_{\mathbb{N}}$  plotted on chart with  $k \leq LCL < UCL$ .

**Exact ANE computation** with Markov chain approach.

**Issues for future research:**

- charts based on different patterns,  
e. g., cycles instead of runs;
- CUSUM or EWMA methods for categorical processes,  
e. g., related to patterns or certain statistics; ...



**In a nutshell,**

the field of categorical time series . . .

- is relevant for practice, and
- offers a lot of topics for future research,  
in any of the disciplines

analysis, modelling and monitoring!

**Thank You**  
**for Your Interest!**



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT



*Fachbereich*  
**Mathematik**

Christian H. Weiß

Department of Mathematics

Darmstadt University of Technology



# Literature



- Box et al. (1994): *Time series analysis – forecasting and control*. 3<sup>rd</sup> edition, Prentice Hall, Inc., New Jersey.
- Bühlmann & Wyner (1999): Variable length Markov chains. *Ann. Stat.*, **27**, 480–513.
- Chryssaphinou et al. (1994): On the waiting time of appearance of given patterns. In: Godbole & Papastavridis (Eds.): *Runs and Patterns in Probability*, Kluwer Academic Publishers, 231–241.
- Duncan (1950): A chi-square chart for controlling a set of percentages. *Industrial Quality Control*, **7**, 11–15.
- Jacobs & Lewis (1983): Stationary discrete autoregressive-moving average time series generated by mixtures. *J. Time Series Anal.*, **4**(1), 19–36.
- Keim & Kriegel (1996): Visualization techniques for mining large databases: A comparison. *IEEE Trans. Knowl. Data Eng.*, **8**(6), 923–938.
- Marcucci (1985): Monitoring multinomial processes. *J. Qual. Tech.*, **17**(2), 86–91.
- Mukhopadhyay (2008): Multivariate attribute control chart using Mahalanobis  $D^2$  statistic. *J. Appl. Statist.*, **35**(4), 421–429.
- Nelson (1987): A chi-square control chart for several proportions. *J. Qual. Tech.*, **19**(4), 229–231.
- Raftery (1985): A model for high-order Markov chains. *J. Royal Stat. Soc., B*, **47**(3), 528–539.

(...)



# Literature

---



(...)

Ribler (1997): *Visualizing categorical time series data with applications to computer and communications network traces*. PhD thesis, Virginia State University.

Uschner (1987): *Streuungsmessung nominaler Merkmale mit Hilfe von Paarvergleichen*. Doctoral dissertation, University Erlangen-Nürnberg.

Vogel & Kiesl (1999): Deskriptive und induktive Eigenschaften zweier Streuungsmaße für nominale Merkmale. In: Vogel (Edt.): *Arbeiten aus der Statistik*, Univ. Bamberg.

Weiß (2008): Visual analysis of categorical time series. *Statist. Meth.*, **5**(1), 56–71.

Weiß (2010): Continuously monitoring categorical processes. *Qual. Tech. Quant. Manag.*, to appear.

Weiß (2011): Empirical measures of signed serial dependence in categorical time series. *J. Statist. Comp. Simulation*, **81**(4), 411–429.

Weiß & Göb (2008): Measuring serial dependence in categorical time series. *Adv. Statist. Anal.*, **92**(1), 71–89.

Ye et al. (2002): Multivariate statistical analysis of audit trails for host-based intrusion detection. *IEEE Trans. Computer*, **51**(7), 810–820.