

# APPROXIMATIONS FOR MULTIDIMENSIONAL DISCRETE SCAN STATISTICS

Doctoral Dissertation

Alexandru Amărioarei

Advisor: Cristian Preda

Laboratoire de Mathématiques Paul Painlevé  
Université de Lille 1, INRIA/Modal Team, France

September 15, 2014

# THE ONE DIMENSIONAL SCAN STATISTICS

Let  $m_1 \leq T_1$  be a positive integers and  $X_1, X_2, \dots, X_{T_1}$  a sequence of r.v.'s.  
 If we consider the moving sums

$$Y_{i_1} = \sum_{j=i_1}^{i_1+m_1-1} X_j$$

then the discrete one dimensional scan statistics is defined as

$$S_{m_1}(T_1) = \max_{1 \leq i_1 \leq T_1 - m_1 + 1} Y_{i_1}.$$

EXAMPLE ( $T_1 = 20$ ,  $m_1 = 3$  AND  $X_{i_1} \sim \mathcal{B}(p)$ ,  $1 \leq i_1 \leq 20$ )

# RELATED STATISTICS

Let  $X_1, \dots, X_{T_1}$  be a sequence of i.i.d.  $0 - 1$  Bernoulli of parameter  $p$

- $W_{m_1, k}$  - the waiting time until we first observe at least  $k$  successes in a window of size  $m_1$

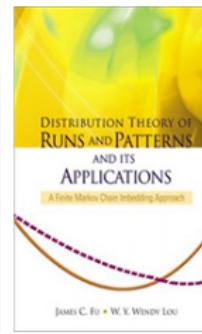
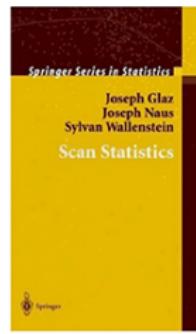
$$\mathbb{P}(W_{m_1, k} \leq T_1) = \mathbb{P}(S_{m_1}(T_1) \geq k)$$

- $D_{T_1}(k)$  - the length of the smallest window that contains at least  $k$  successes

$$\mathbb{P}(D_{T_1}(k) \leq m_1) = \mathbb{P}(S_{m_1}(T_1) \geq k)$$

- $L_{T_1}$  - the length of the longest success run

$$\mathbb{P}(L_{T_1} \geq m_1) = \mathbb{P}(S_{m_1}(T_1) \geq m_1) = \mathbb{P}(S_{m_1}(T_1) = m_1)$$



# PROBLEM AND APPROACHES

## PROBLEM

Find a good estimate for the distribution of the discrete scan statistic

$$\mathbb{P}(S_{m_1}(T_1) \leq \tau).$$

## Previous work:

- One dimensional scan statistics
  - Exact results ([Naus, 1974], [Fu, 2001], [Gao et al., 2005])
  - Approximations: product-type, Poisson ([Naus, 1982], [Chen and Glaz, 1997], [Glaz et al., 2001])
  - Bounds ([Glaz, 1990], [Glaz and Naus, 1991])
- Two dimensional scan statistics
  - Approximations: product-type, Poisson ([Chen and Glaz, 1996], [Boutsikas and Koutras, 2000])
  - Bounds ([Chen and Glaz, 1996], [Boutsikas and Koutras, 2003])
- Three dimensional scan statistics
  - Approximations: product-type, Poisson ([Guerriero et al., 2010])

### ► Product-Type Approximations

# THE FOCUS OF THIS THESIS

We consider the  $d$  dimensional discrete scan statistics over a random field generated by:

- i.i.d. observations
- dependent (block-factor type) observations

We present

- accurate approximations
- error bounds
- simulation aspects

# OUTLINE

## 1 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (I.I.D. MODEL)

- Framework
- Extremes of 1-dependent stationary sequences
- Scan statistics and 1-dependent sequences
- Simulation methods and computational aspects
- Numerical examples

## 2 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (BLOCK-FACTOR MODEL)

- Model and discussion
- Applications

## 3 CONCLUSIONS AND PERSPECTIVES

- Conclusions
- Perspectives

## 4 REFERENCES

# OUTLINE

## 1 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (I.I.D. MODEL)

- Framework

- Extremes of 1-dependent stationary sequences
- Scan statistics and 1-dependent sequences
- Simulation methods and computational aspects
- Numerical examples

## 2 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (BLOCK-FACTOR MODEL)

- Model and discussion
- Applications

## 3 CONCLUSIONS AND PERSPECTIVES

- Conclusions
- Perspectives

## 4 REFERENCES

# Framework

# THE $d$ -DIMENSIONAL DISCRETE SCAN STATISTICS

Let  $T_1, T_2, \dots, T_d$  be positive integers, with  $d \geq 1$

- The rectangular region,  $\mathcal{R}_d = [0, T_1] \times [0, T_2] \times \cdots \times [0, T_d]$
- The r.v.'s  $X_{s_1, s_2, \dots, s_d}$ ,  $1 \leq s_j \leq T_j$ ,  $j \in \{1, 2, \dots, d\}$

Let  $2 \leq m_j \leq T_j$ ,  $1 \leq j \leq d$ , be positive integers

- Define for  $1 \leq i_l \leq T_l - m_l + 1$ ,  $1 \leq l \leq d$ ,

$$Y_{i_1, i_2, \dots, i_d} = \sum_{s_1=i_1}^{i_1+m_1-1} \sum_{s_2=i_2}^{i_2+m_2-1} \cdots \sum_{s_d=i_d}^{i_d+m_d-1} X_{s_1, s_2, \dots, s_d}$$

- The  $d$ -dimensional discrete scan statistic,

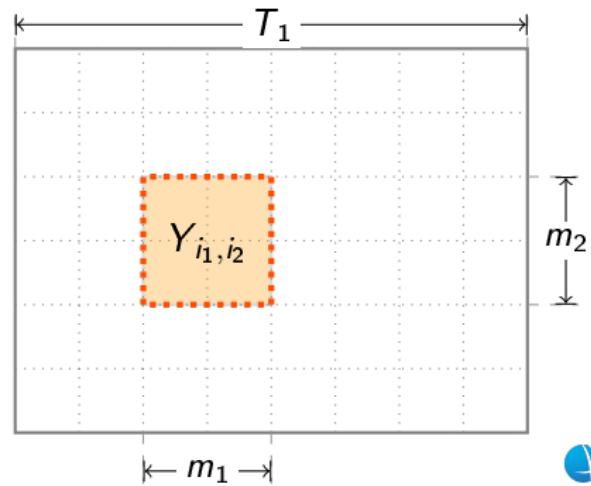
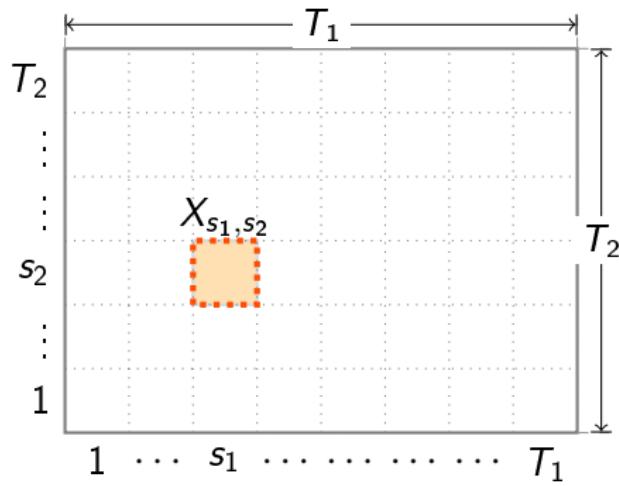
$$S_{\mathbf{m}}(\mathbf{T}) = \max_{\substack{1 \leq i_j \leq T_j - m_j + 1 \\ j \in \{1, 2, \dots, d\}}} Y_{i_1, i_2, \dots, i_d}$$

with  $\mathbf{m} = (m_1, m_2, \dots, m_d)$  and  $\mathbf{T} = (T_1, T_2, \dots, T_d)$

# EXAMPLE: TWO DIMENSIONAL SCAN STATISTICS

We have for  $d = 2$

$$Y_{i_1, i_2} = \sum_{s_1=i_1}^{i_1+m_1-1} \sum_{s_2=i_2}^{i_2+m_2-1} X_{s_1, s_2}, \quad S_{m_1, m_2}(T_1, T_2) = \max_{\substack{1 \leq i_1 \leq T_1 - m_1 + 1 \\ 1 \leq i_2 \leq T_2 - m_2 + 1}} Y_{i_1, i_2}$$



# ANIMATION FOR 3 DIMENSIONAL SCAN STATISTICS

# OBJECTIVE

Find a good estimate for the distribution of  $d$ -dimensional discrete scan statistic

$$Q_{\mathbf{m}}(\mathbf{T}) = \mathbb{P}(S_{\mathbf{m}}(\mathbf{T}) \leq \tau)$$

The distribution of  $S_{\mathbf{m}}(\mathbf{T})$  is used for testing the null hypotheses of randomness against the alternative hypothesis of clustering.

## EXAMPLE

Bernoulli model

$H_0$ : The r.v.'s  $X_{s_1, s_2, \dots, s_d}$  are i.i.d.  $\mathcal{B}(p)$

$H_1$ : There exists

$$\mathcal{R}(i_1, i_2, \dots, i_d) = [i_1 - 1, i_1 + m_1 - 1] \times \dots \times [i_d - 1, i_d + m_d - 1] \subset \mathcal{R}_d$$

where the r.v.'s  $X_{s_1, s_2, \dots, s_d} \sim \mathcal{B}(p')$ ,  $p' > p$  and  $X_{s_1, s_2, \dots, s_d} \sim \mathcal{B}(p)$

outside  $\mathcal{R}(i_1, i_2, \dots, i_d)$

# OUTLINE

## 1 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (I.I.D. MODEL)

- Framework
- Extremes of 1-dependent stationary sequences
- Scan statistics and 1-dependent sequences
- Simulation methods and computational aspects
- Numerical examples

## 2 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (BLOCK-FACTOR MODEL)

- Model and discussion
- Applications

## 3 CONCLUSIONS AND PERSPECTIVES

- Conclusions
- Perspectives

## 4 REFERENCES

# Extremes of 1-dependent stationary sequences

# DEFINITIONS AND NOTATIONS

Let  $(Z_n)_{n \geq 1}$  be a sequence of random variables

## *m*-DEPENDENCE

The sequence  $(Z_n)_{n \geq 1}$  is *m*-dependent,  $m \geq 1$ , if for any  $h \geq 1$  the  $\sigma$ -fields generated by  $\{Z_1, \dots, Z_h\}$  and  $\{Z_{h+m+1}, \dots\}$  are independent.

## STATIONARITY (IN THE STRONG SENSE)

The sequence  $(Z_n)_{n \geq 1}$  is stationary if for all  $k \geq 1$ , for all  $h \geq 0$  and for all  $t_1, \dots, t_k$  the families  $\{Z_{t_1}, \dots, Z_{t_k}\}$  and  $\{Z_{t_1+h}, \dots, Z_{t_k+h}\}$  have the same joint distribution.

## NOTATION

For  $x < \sup\{u | \mathbb{P}(Z_1 \leq u) < 1\}$ ,

$$q_n = q_n(x) = \mathbb{P}(\max(Z_1, \dots, Z_n) \leq x)$$

# THE MAIN RESULT

## THEOREM [HAIMAN, 1999]

For  $x$  such that  $\mathbb{P}(Z_1 > x) = 1 - q_1 < 0.025$  and  $n > 3$  we have

$$\left| q_n - \frac{2q_1 - q_2}{[1 + q_1 - q_2 + 2(q_1 - q_2)^2]^n} \right| \leq n\Delta_2^H(1 - q_1)^2$$

- $\Delta_2^H = 3.3 + \frac{9}{n} + \left[ 15.51n(1 - q_1) + \frac{561}{n} \right] (1 - q_1)$ .

## THEOREM [AMĂRIOAREI, 2012]

For  $x$  such that  $\mathbb{P}(Z_1 > x) = 1 - q_1 < 0.1$  and  $n > 3$  we have

$$\left| q_n - \frac{2q_1 - q_2}{[1 + q_1 - q_2 + 2(q_1 - q_2)^2]^n} \right| \leq n\Delta_2(1 - q_1)^2$$

- $\Delta_2 = F(q_1, n) = 1 + \frac{3}{n} + \left[ K(1 - q_1) + \frac{\Gamma(1 - q_1)}{n} \right] (1 - q_1)$ .

- Increased range of applicability
- Sharper error bounds

# THE MAIN RESULT

## THEOREM [HAIMAN, 1999]

For  $x$  such that  $\mathbb{P}(Z_1 > x) = 1 - q_1 < 0.025$  and  $n > 3$  we have

$$\left| q_n - \frac{2q_1 - q_2}{[1 + q_1 - q_2 + 2(q_1 - q_2)^2]^n} \right| \leq n\Delta_2^H(1 - q_1)^2$$

- $\Delta_2^H = 3.3 + \frac{9}{n} + \left[ 15.51n(1 - q_1) + \frac{561}{n} \right] (1 - q_1)$ .

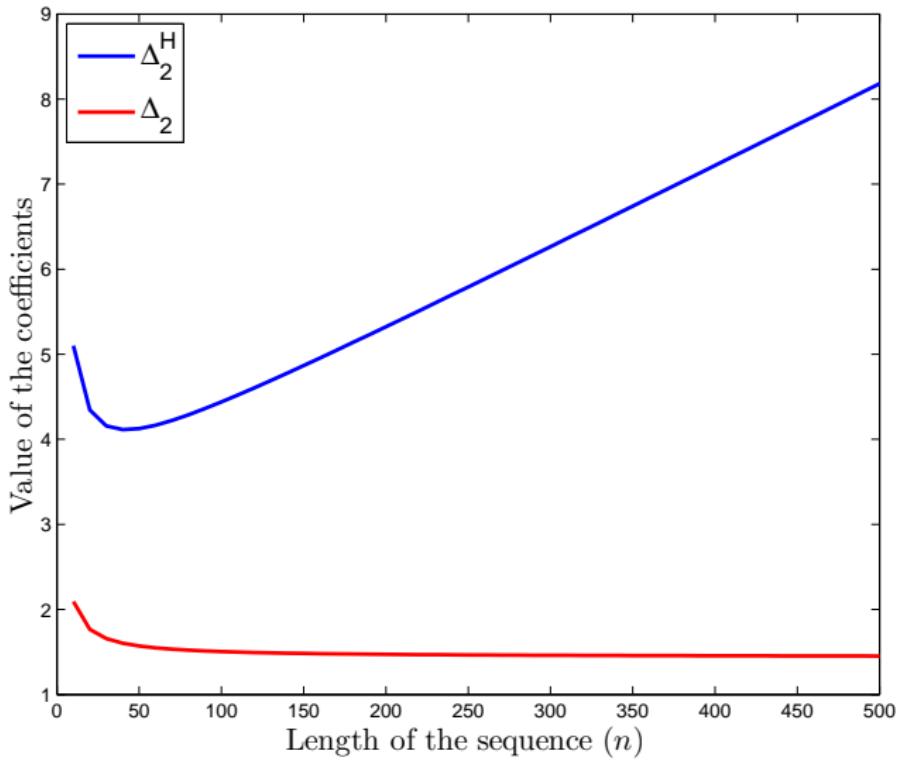
## THEOREM [AMĂRIOAREI, 2012]

For  $x$  such that  $\mathbb{P}(Z_1 > x) = 1 - q_1 < 0.1$  and  $n > 3$  we have

$$\left| q_n - \frac{2q_1 - q_2}{[1 + q_1 - q_2 + 2(q_1 - q_2)^2]^n} \right| \leq n\Delta_2(1 - q_1)^2$$

- $\Delta_2 = F(q_1, n) = 1 + \frac{3}{n} + \left[ K(1 - q_1) + \frac{\Gamma(1 - q_1)}{n} \right] (1 - q_1)$ .

- Increased range of applicability
- Sharper error bounds

DIFFERENCE BETWEEN THE RESULTS:  $1 - q_1 = 0.025$ 

# OUTLINE

## 1 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (I.I.D. MODEL)

- Framework
- Extremes of 1-dependent stationary sequences
- **Scan statistics and 1-dependent sequences**
- Simulation methods and computational aspects
- Numerical examples

## 2 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (BLOCK-FACTOR MODEL)

- Model and discussion
- Applications

## 3 CONCLUSIONS AND PERSPECTIVES

- Conclusions
- Perspectives

## 4 REFERENCES

# Scan statistics and 1-dependent sequences

# THE KEY IDEA

## MAIN OBSERVATION

The scan statistic r.v. can be viewed as a maximum of a sequence of 1-dependent stationary r.v..

- The idea:
  - one dimensional scan statistic: [Haiman, 2000], [Haiman, 2007]
  - two dimensional scan statistic: [Haiman and Preda, 2002], [Haiman and Preda, 2006]
  - three dimensional scan statistic: [Amărioarei and Preda, 2013a]

# $S_m(\mathbf{T})$ VIEWED AS MAXIMUM OF 1-DEPENDENT R.V.'S

Let  $L_j = \frac{T_j}{m_j - 1}$ ,  $j \in \{1, 2, \dots, d\}$ , be positive integers

- Define for each  $k_1 \in \{1, 2, \dots, L_1 - 1\}$  the random variables

$$Z_{k_1} = \max_{\substack{(k_1-1)(m_1-1)+1 \leq i_1 \leq k_1(m_1-1) \\ 1 \leq j \leq (L_j-1)(m_j-1) \\ j \in \{2, \dots, d\}}} Y_{i_1, i_2, \dots, i_d}$$

- $(Z_j)_j$  is 1-dependent and stationary

- Observe

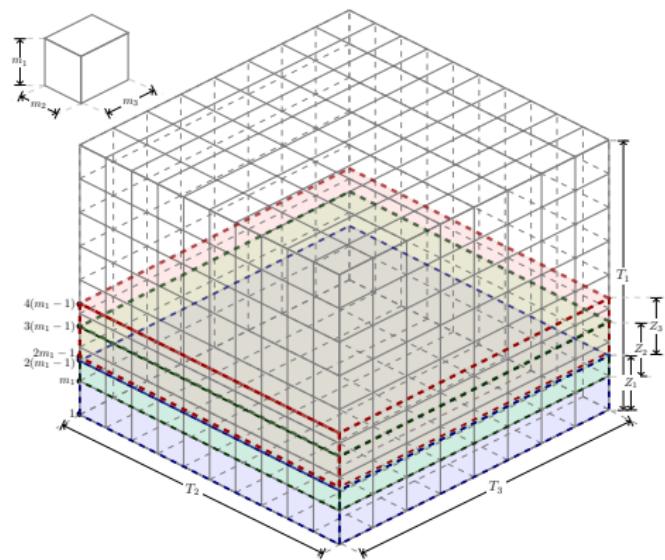
$$S_m(\mathbf{T}) = \max_{1 \leq k_1 \leq L_1 - 1} Z_{k_1}$$

EXAMPLE ( $d = 1$ )

$$\underbrace{X_1, X_2, \dots, X_{m_1-1}}_{Z_1}, \overbrace{\underbrace{X_{m_1}, \dots, X_{2(m_1-1)}}^{Z_2}, \underbrace{X_{2m_1-1}, \dots, X_{3(m_1-1)}}_{Z_3}, \dots, X_{4(m_1-1)}}$$

# $S_m(\mathbf{T})$ VIEWED AS MAXIMUM OF 1-DEPENDENT R.V.'S

EXAMPLE ( $d = 3$ )



# APPROXIMATION PROCESS

Define for  $t_1 \in \{2, 3\}$ ,

$$Q_{t_1} = Q_{t_1}(\tau) = \mathbb{P} \left( \bigcap_{k_1=1}^{t_1-1} \{Z_{k_1} \leq \tau\} \right) = \mathbb{P} \left( \max_{\substack{1 \leq i_1 \leq (t_1-1)(m_1-1) \\ 1 \leq i_j \leq (L_j-1)(m_j-1) \\ j \in \{2, \dots, d\}}} Y_{i_1, i_2, \dots, i_d} \leq \tau \right)$$

If  $1 - Q_2 \leq 0.1$  then

$$\left| Q_m(T) - \frac{2Q_2 - Q_3}{[1+Q_2 - Q_3 + 2(Q_2 - Q_3)^2]^{L_1-1}} \right| \leq (L_1 - 1)F(Q_2, L_1 - 1)(1 - Q_2)^2$$

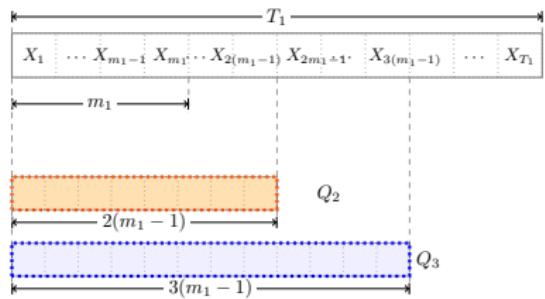
EXAMPLE ( $d = 1$ )

- The approximation

$$\mathbb{P}(S_{m_1}(T_1) \leq \tau) \approx \frac{2Q_2 - Q_3}{[1+Q_2 - Q_3 + 2(Q_2 - Q_3)^2]^{L_1-1}}$$

- Approximation error, about

$$(L_1 - 1)F(Q_2, L_1 - 1)(1 - Q_2)^2$$



# APPROXIMATION PROCESS

The approximation of  $S_m(\mathbf{T})$  is an iterative process. The  $s$  step,  $1 \leq s \leq d$ , becomes:

- Let

$$Q_{t_1, t_2, \dots, t_s} = Q_{t_1, t_2, \dots, t_s}(\tau) = \mathbb{P} \left( \begin{array}{l} \max_{\substack{1 \leq i_l \leq (t_l-1)(m_l-1) \\ l \in \{1, \dots, s\}}} Y_{i_1, i_2, \dots, i_d} \leq \tau \\ \max_{\substack{1 \leq i_j \leq (L_j-1)(m_j-1) \\ j \in \{s+1, \dots, d\}}} Y_{i_1, i_2, \dots, i_d} \leq \tau \end{array} \right)$$

- Define for  $t_l \in \{2, 3\}$ ,  $l \in \{1, \dots, s-1\}$  and  $k_s \in \{1, 2, \dots, L_s - 1\}$

$$Z_{k_s}^{(t_1, t_2, \dots, t_{s-1})} = \max_{\substack{1 \leq i_l \leq (t_l-1)(m_l-1) \\ l \in \{1, 2, \dots, s-1\}}} Y_{i_1, i_2, \dots, i_d} \quad (k_s-1)(m_s-1)+1 \leq i_s \leq k_s(m_s-1) \\ \max_{\substack{1 \leq i_j \leq (L_j-1)(m_j-1) \\ j \in \{s+1, \dots, d\}}} Y_{i_1, i_2, \dots, i_d}$$

- $\{Z_1^{(t_1, t_2, \dots, t_{s-1})}, \dots, Z_{L_s-1}^{(t_1, t_2, \dots, t_{s-1})}\}$  forms a 1-dependent stationary sequence
- If we take  $H(x, y, m) = \frac{2x-y}{[1+x-y+2(x-y)^2]^{m-1}}$ , then we have the approximation

$$\left| Q_{t_1, \dots, t_{s-1}} - H(Q_{t_1, \dots, t_{s-1}, 2}, Q_{t_1, \dots, t_{s-1}, 3}, L_s) \right| \leq (L_s - 1) F(Q_{t_1, \dots, t_{s-1}, 2}, L_s - 1) (1 - Q_{t_1, \dots, t_{s-1}, 2})^2$$

# ILLUSTRATION FOR $d = 2$

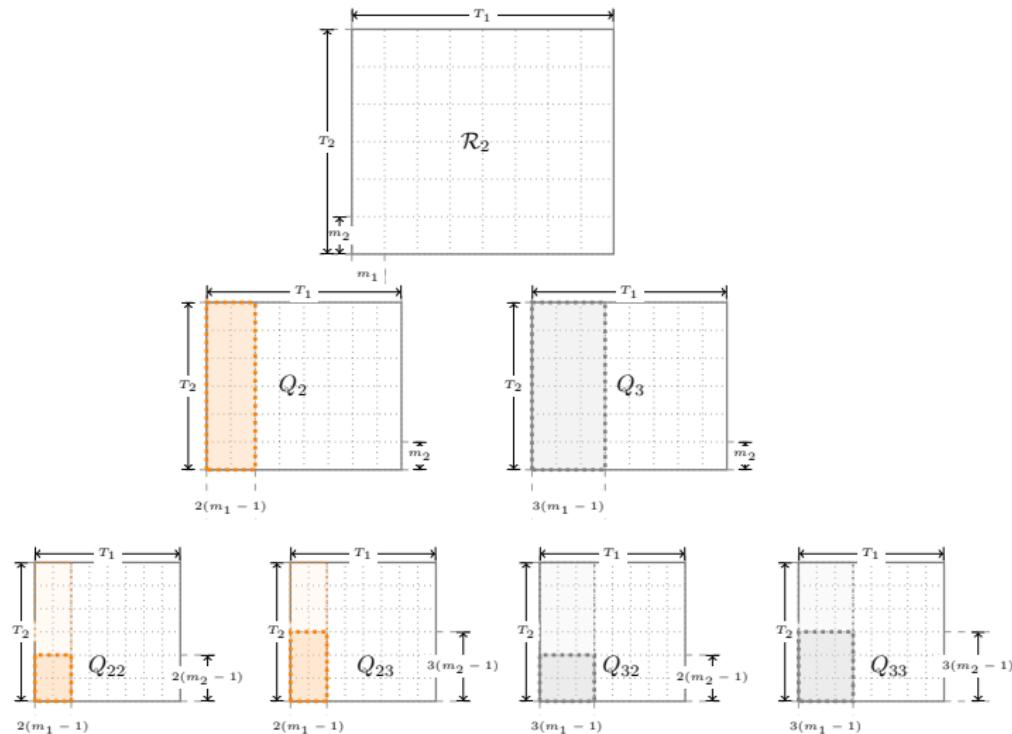
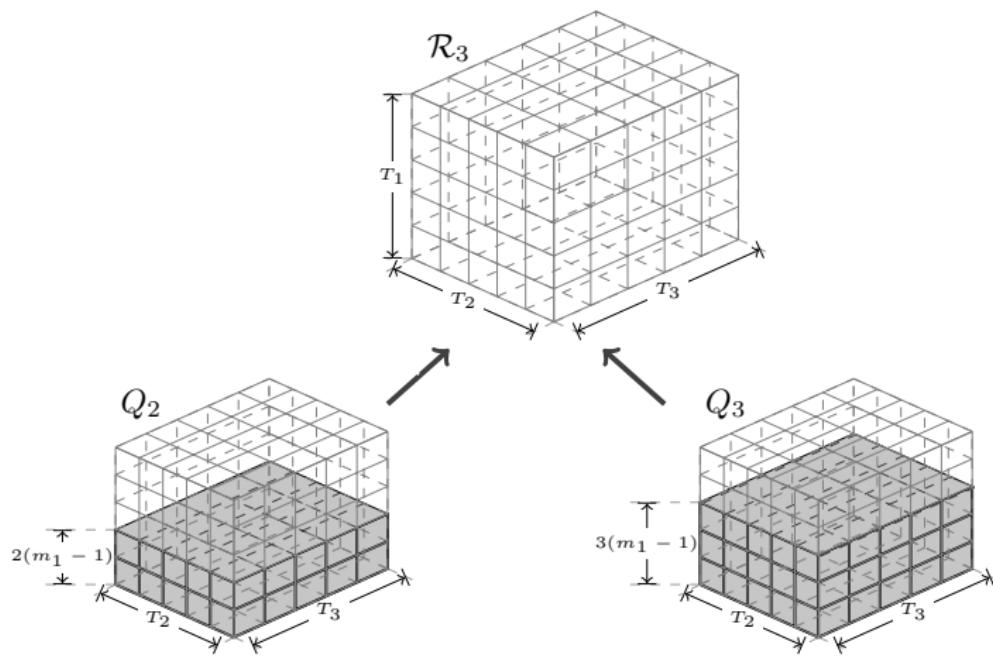
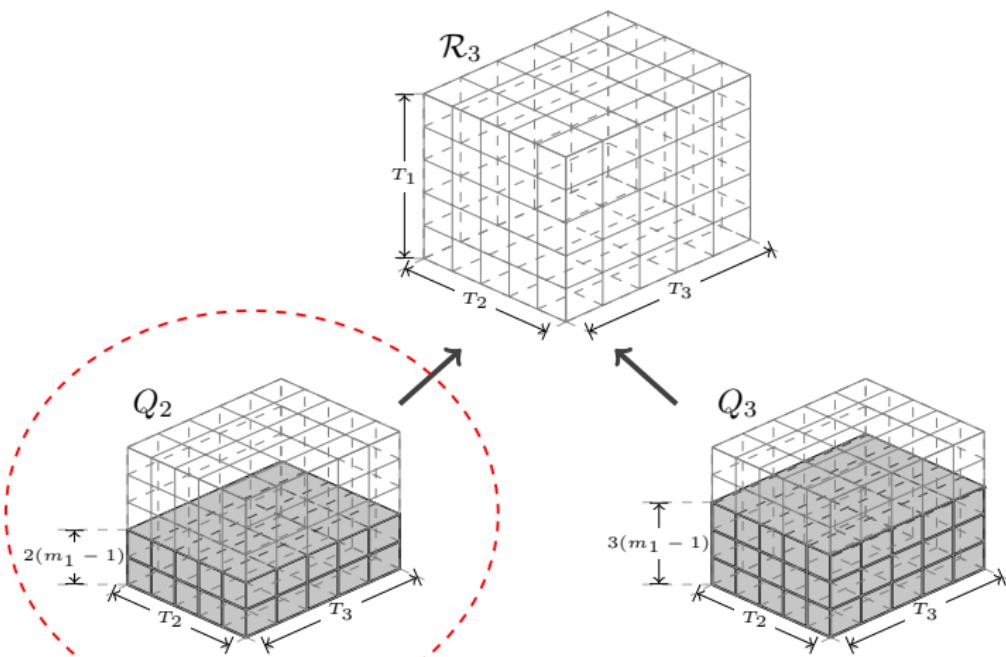
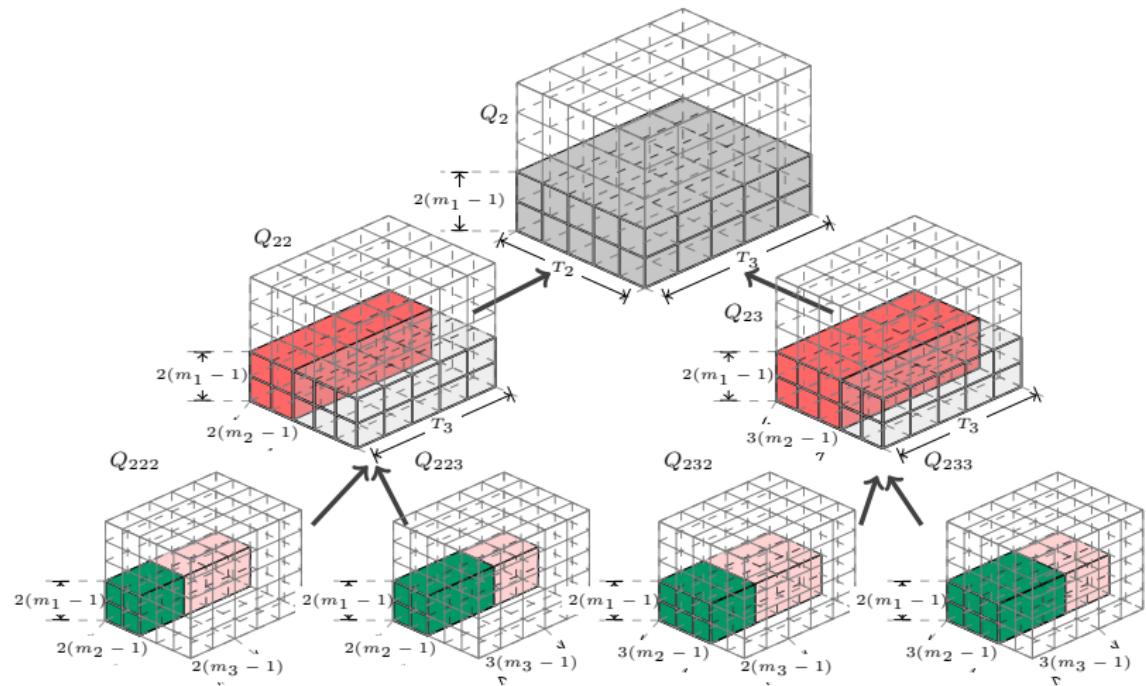


ILLUSTRATION FOR  $d = 3$ 

# ILLUSTRATION FOR $d = 3$



# ILLUSTRATION FOR $d = 3$



# ERROR BOUNDS

Let  $\gamma_{t_1, \dots, t_d} = Q_{t_1, \dots, t_d}$ , with  $t_j \in \{2, 3\}$ ,  $j \in \{1, \dots, d\}$ , and define

$$\gamma_{t_1, \dots, t_{s-1}} = H(\gamma_{t_1, \dots, t_{s-1}, 2}, \gamma_{t_1, \dots, t_{s-1}, 3}, L_s), \quad 2 \leq s \leq d$$

Denote with  $\hat{Q}_{t_1, \dots, t_d}$  the estimated value of  $Q_{t_1, \dots, t_d}$  and define

$$\hat{Q}_{t_1, \dots, t_{s-1}} = H(\hat{Q}_{t_1, \dots, t_{s-1}, 2}, \hat{Q}_{t_1, \dots, t_{s-1}, 3}, L_s), \quad 2 \leq s \leq d$$

## OBJECTIVE

$$Q_m(T) \approx H(\hat{Q}_2, \hat{Q}_3, L_1)$$

We observe that

$$|Q_m(T) - H(\hat{Q}_2, \hat{Q}_3, L_1)| \leq |Q_m(T) - H(\gamma_2, \gamma_3, L_1)| + |H(\gamma_2, \gamma_3, L_1) - H(\hat{Q}_2, \hat{Q}_3, L_1)|$$

The quantities  $\hat{Q}_{t_1, \dots, t_d}$  will be estimated by Monte Carlo simulations.

[Error bounds](#)

# ERROR BOUNDS

Let  $\gamma_{t_1, \dots, t_d} = Q_{t_1, \dots, t_d}$ , with  $t_j \in \{2, 3\}$ ,  $j \in \{1, \dots, d\}$ , and define

$$\gamma_{t_1, \dots, t_{s-1}} = H(\gamma_{t_1, \dots, t_{s-1}, 2}, \gamma_{t_1, \dots, t_{s-1}, 3}, L_s), \quad 2 \leq s \leq d$$

Denote with  $\hat{Q}_{t_1, \dots, t_d}$  the estimated value of  $Q_{t_1, \dots, t_d}$  and define

$$\hat{Q}_{t_1, \dots, t_{s-1}} = H(\hat{Q}_{t_1, \dots, t_{s-1}, 2}, \hat{Q}_{t_1, \dots, t_{s-1}, 3}, L_s), \quad 2 \leq s \leq d$$

## OBJECTIVE

$$Q_m(T) \approx H(\hat{Q}_2, \hat{Q}_3, L_1)$$

We observe that

$$|Q_m(T) - H(\hat{Q}_2, \hat{Q}_3, L_1)| \leq \underbrace{|Q_m(T) - H(\gamma_2, \gamma_3, L_1)|}_{E_{app}(d)} + \underbrace{|H(\gamma_2, \gamma_3, L_1) - H(\hat{Q}_2, \hat{Q}_3, L_1)|}_{E_{sf}(d)}$$

The quantities  $\hat{Q}_{t_1, \dots, t_d}$  will be estimated by Monte Carlo simulations.

► Error bounds

# ERROR BOUNDS

Let  $\gamma_{t_1, \dots, t_d} = Q_{t_1, \dots, t_d}$ , with  $t_j \in \{2, 3\}$ ,  $j \in \{1, \dots, d\}$ , and define

$$\gamma_{t_1, \dots, t_{s-1}} = H(\gamma_{t_1, \dots, t_{s-1}, 2}, \gamma_{t_1, \dots, t_{s-1}, 3}, L_s), \quad 2 \leq s \leq d$$

Denote with  $\hat{Q}_{t_1, \dots, t_d}$  the estimated value of  $Q_{t_1, \dots, t_d}$  and define

$$\hat{Q}_{t_1, \dots, t_{s-1}} = H(\hat{Q}_{t_1, \dots, t_{s-1}, 2}, \hat{Q}_{t_1, \dots, t_{s-1}, 3}, L_s), \quad 2 \leq s \leq d$$

## OBJECTIVE

$$Q_m(T) \approx H(\hat{Q}_2, \hat{Q}_3, L_1)$$

We observe that

$$\left| Q_m(T) - H(\hat{Q}_2, \hat{Q}_3, L_1) \right| \leq \underbrace{|Q_m(T) - H(\gamma_2, \gamma_3, L_1)|}_{E_{app}(d) \leq E_{sapp}(d)} + \underbrace{|H(\gamma_2, \gamma_3, L_1) - H(\hat{Q}_2, \hat{Q}_3, L_1)|}_{E_{sf}(d)}$$

The quantities  $\hat{Q}_{t_1, \dots, t_d}$  will be estimated by Monte Carlo simulations.

► Error bounds

# OUTLINE

## 1 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (I.I.D. MODEL)

- Framework
- Extremes of 1-dependent stationary sequences
- Scan statistics and 1-dependent sequences
- Simulation methods and computational aspects**
- Numerical examples

## 2 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (BLOCK-FACTOR MODEL)

- Model and discussion
- Applications

## 3 CONCLUSIONS AND PERSPECTIVES

- Conclusions
- Perspectives

## 4 REFERENCES

# Simulation methods and computational aspects

# NAIVE HIT-OR-MISS MC

## OBJECTIVE

Find an estimate for  $\mathbb{P}_{H_0}(S_m(\mathbf{T}) \geq \tau)$ .

### Algorithm 1 Classical Monte Carlo algorithm for scan statistics

**Begin**

    Repeat for each  $k$  from 1 to  $ITER$  (iterations number)

1: Generate  $\mathbf{X}^{(k)} = \left\{ X_{s_1, s_2, \dots, s_d}^{(k)}, 1 \leq s_j \leq T_j, 1 \leq j \leq d \right\}$  under  $H_0$

2: Compute the  $d$ -dimensional scan statistics  $S_m^{(k)}(\mathbf{T})$  over  $\mathbf{X}^{(k)}$

End Repeat

Return

$$\widehat{p}_{MC} = \frac{1}{ITER} \sum_{i=1}^{ITER} \mathbf{1}_{\{S_m^{(i)}(\mathbf{T}) \geq \tau\}}, \quad \widehat{s.e.MC} = \sqrt{\frac{\widehat{p}_{MC}(1 - \widehat{p}_{MC})}{ITER}}$$

the unbiased direct Monte Carlo estimate and its consistent standard error estimate.

**End**

- computationally intensive since just a fraction of the generated observations will cause a rejection
- needs a large number of replications in order to reduce the standard error estimate to an acceptable level (especially for  $d \geq 2$ )

# IMPORTANCE SAMPLING FOR SCAN STATISTICS

## IDEA BEHIND IMPORTANCE SAMPLING

Find a good change of measure that leads to an efficient sampling process.

The method was previously used for solving the problem of:

- union count: [Frigessi and Vercellis, 1984], [Fishman, 1996]
- exceeding probabilities: [Naiman and Wynn, 1997]
- scan statistics: [Naiman and Priebe, 2001], [Malley et al., 2002]

We are interested in evaluating the probability

$$\mathbb{P}_{H_0}(S_{\mathbf{m}}(\mathbf{T}) \geq \tau) = \mathbb{P}\left(\bigcup_{i_1=1}^{T_1-m_1+1} \dots \bigcup_{i_d=1}^{T_d-m_d+1} E_{i_1, \dots, i_d}\right) = \int G(\mathbf{x})f(\mathbf{x}) d\mathbf{x}$$

where  $E_{i_1, \dots, i_d} = \{Y_{i_1, \dots, i_d} \geq \tau\}$ ,  $G(\mathbf{x}) = \mathbf{1}_E(\mathbf{x})$ ,  $E = \bigcup_{i_1=1}^{T_1-m_1+1} \dots \bigcup_{i_d=1}^{T_d-m_d+1} E_{i_1, \dots, i_d}$  and  $f$  is the joint density of  $Y_{i_1, \dots, i_d}$  under  $H_0$ .

# IMPORTANCE SAMPLING FOR SCAN STATISTICS

We introduce the change of measure

$$g(\mathbf{x}) = \sum_{j_1=1}^{\tau_1-m_1+1} \cdots \sum_{j_d=1}^{\tau_d-m_d+1} \left\{ \frac{\mathbb{P}(E_{j_1, \dots, j_d})}{B(d)} \right\} \left\{ \frac{1_{E_{j_1, \dots, j_d}} f(\mathbf{x})}{\mathbb{P}(E_{j_1, \dots, j_d})} \right\}$$

and we observe that  $\mathbb{P}_{H_0}(S_m(\mathbf{T}) \geq \tau) = B(d)\rho(d)$

- the Bonferroni upper bound  $B(d)$

$$B(d) = \sum_{i_1=1}^{\tau_1-m_1+1} \cdots \sum_{i_d=1}^{\tau_d-m_d+1} \mathbb{P}(E_{i_1, \dots, i_d})$$

- the correction factor  $\rho(d)$  between 0 and 1

$$\rho(d) = \sum_{i_1=1}^{\tau_1-m_1+1} \cdots \sum_{i_d=1}^{\tau_d-m_d+1} p_{i_1, \dots, i_d} \int \frac{1}{C(Y)} d\mathbb{P}_{H_0}(\cdot | E_{i_1, \dots, i_d})$$

where

$$p_{i_1, \dots, i_d} = \frac{1}{(\tau_1-m_1+1) \cdots (\tau_d-m_d+1)}, \quad C(Y) = \sum_{i_1=1}^{\tau_1-m_1+1} \cdots \sum_{i_d=1}^{\tau_d-m_d+1} 1_{E_{i_1, \dots, i_d}}$$

# IMPORTANCE SAMPLING FOR SCAN STATISTICS

## Algorithm 2 Importance Sampling Algorithm for Scan Statistics

**Begin**

Repeat for each  $k$  from 1 to  $ITER$  (iterations number)

- 1: Generate uniformly the  $d$ -tuple  $(i_1^{(k)}, \dots, i_d^{(k)})$  from the set  $\{1, \dots, T_1 - m_1 + 1\} \times \dots \times \{1, \dots, T_d - m_d + 1\}$ .
- 2: Given the  $d$ -tuple  $(i_1^{(k)}, \dots, i_d^{(k)})$ , generate a sample of the random field  $\tilde{\mathbf{X}}^{(k)} = \{\tilde{X}_{s_1, s_2, \dots, s_d}^{(k)}\}$ , with  $s_j \in \{1, \dots, T_j\}$  and  $j \in \{1, \dots, d\}$ , from the conditional distribution of  $\mathbf{X}$  given  $\left\{ Y_{i_1^{(k)}, \dots, i_d^{(k)}} \geq \tau \right\}$ .
- 3: Take  $c_k = C(\tilde{\mathbf{X}}^{(k)})$  the number of all  $d$ -tuple  $(i_1, \dots, i_d)$  for which  $\tilde{Y}_{i_1, \dots, i_d} \geq \tau$  and put  $\hat{\rho}_k(d) = \frac{1}{c_k}$ .

End Repeat

Return

$$\hat{\rho}(d) = \frac{1}{ITER} \sum_{k=1}^{ITER} \hat{\rho}_k(d), \quad Var[\hat{\rho}(d)] \approx \frac{1}{ITER-1} \sum_{k=1}^{ITER} \left( \hat{\rho}_k(d) - \frac{1}{ITER} \sum_{k=1}^{ITER} \hat{\rho}_k(d) \right)^2$$

**End**

# IMPLEMENTATION PROBLEMS

Algorithm 2 presents two main difficulties:

- A) being able to sample from the conditional distribution of  $\mathbf{X}$  given  $\left\{ Y_{i_1^{(k)}, \dots, i_d^{(k)}} \geq \tau \right\}$  in **Step 2**
- B) the number of locality statistics that exceed the predetermined threshold is supposed to be found in a *reasonable* time

Partial solutions were found for:

- A) binomial, Poisson and Gaussian model
- B) cumulative counts or *fast spatial scan* techniques (see [Neil, 2006], [Neil, 2012])

▶ Scan 1d for normal data

# OUTLINE

## 1 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (I.I.D. MODEL)

- Framework
- Extremes of 1-dependent stationary sequences
- Scan statistics and 1-dependent sequences
- Simulation methods and computational aspects
- Numerical examples

## 2 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (BLOCK-FACTOR MODEL)

- Model and discussion
- Applications

## 3 CONCLUSIONS AND PERSPECTIVES

- Conclusions
- Perspectives

## 4 REFERENCES

# Numerical examples

# EXAMPLES FOR $d = 1, 2, 3$ WHEN $X_{i_1, \dots, i_d} \sim \mathcal{B}(n, p)$

TABLE 1 :  $n = 1, p = 0.005, m_1 = 10, T_1 = 1000, It_{App} = 10^4$

$\tau$	Exact	Glaz et al. Product-type	Our Approximation	Approximation Error	Lower Bound	Upper Bound
1	0.810209	0.810216	0.810404	0.001111	0.809903	0.810439
2	0.995764	0.995764	0.995764	$3 \times 10^{-7}$	0.995764	0.995764
3	0.999950	0.999950	0.999950	$4 \times 10^{-11}$	0.999950	0.999950

TABLE 2 :  $n = 5, p = 0.002, m_1 = 5, m_2 = 10, T_1 = 50, T_2 = 80, It_{App} = 10^4$

$\tau$	$\hat{\mathbb{P}}(S \leq \tau)$	Glaz et al. Product-type	Our Approximation	Total Error	Lower Bound	Upper Bound
4	0.894654	0.873256	0.893724	0.037136	0.803422	0.944318
5	0.988003	0.986249	0.988144	0.002125	0.981418	0.993451
6	0.998963	0.998847	0.998963	0.000152	0.998543	0.999401
7	0.999926	0.999919	0.999925	$9 \times 10^{-6}$	0.999903	0.999955
8	0.999995	0.999995	0.999995	$5 \times 10^{-7}$	0.999994	0.999997

TABLE 3 :  $n = 1, p = 0.0001, m_1 = m_2 = m_3 = 5, T_1 = T_2 = T_3 = 60, It_{App} = 10^5$

$\tau$	$\hat{\mathbb{P}}(S \leq \tau)$	Glaz et al. Product-type	Our Approximation	Approximation Error	Simulation Error	Total Error
2	0.993294	0.993241	0.993192	0.000010	0.001367	0.001377
3	0.999963	0.999964	0.999963	0.000000	0.000005	0.000005
4	0.999999	0.999999	0.999999	0.000000	$2 \times 10^{-9}$	$2 \times 10^{-9}$

# EXAMPLES FOR $d = 1, 2, 3$ WHEN $X_{i_1, \dots, i_d} \sim \mathcal{B}(n, p)$

TABLE 1 :  $n = 1, p = 0.005, m_1 = 10, T_1 = 1000, It_{App} = 10^4$

$\tau$	Exact	Glaz et al. Product-type	Our Approximation	Approximation Error	Lower Bound	Upper Bound
1	0.810209	0.810216	0.810404	0.001111	0.809903	0.810439
2	0.995764	0.995764	0.995764	$3 \times 10^{-7}$	0.995764	0.995764
3	0.999950	0.999950	0.999950	$4 \times 10^{-11}$	0.999950	0.999950

TABLE 2 :  $n = 5, p = 0.002, m_1 = 5, m_2 = 10, T_1 = 50, T_2 = 80, It_{App} = 10^4$

$\tau$	$\hat{\mathbb{P}}(S \leq \tau)$	Glaz et al. Product-type	Our Approximation	Total Error	Lower Bound	Upper Bound
4	0.894654	0.873256	0.893724	0.037136	0.803422	0.944318
5	0.988003	0.986249	0.988144	0.002125	0.981418	0.993451
6	0.998963	0.998847	0.998963	0.000152	0.998543	0.999401
7	0.999926	0.999919	0.999925	$9 \times 10^{-6}$	0.999903	0.999955
8	0.999995	0.999995	0.999995	$5 \times 10^{-7}$	0.999994	0.999997

TABLE 3 :  $n = 1, p = 0.0001, m_1 = m_2 = m_3 = 5, T_1 = T_2 = T_3 = 60, It_{App} = 10^5$

$\tau$	$\hat{\mathbb{P}}(S \leq \tau)$	Glaz et al. Product-type	Our Approximation	Approximation Error	Simulation Error	Total Error
2	0.993294	0.993241	0.993192	0.000010	0.001367	0.001377
3	0.999963	0.999964	0.999963	0.000000	0.000005	0.000005
4	0.999999	0.999999	0.999999	0.000000	$2 \times 10^{-9}$	$2 \times 10^{-9}$

# EXAMPLES FOR $d = 1, 2, 3$ WHEN $X_{i_1, \dots, i_d} \sim \mathcal{B}(n, p)$

TABLE 1 :  $n = 1, p = 0.005, m_1 = 10, T_1 = 1000, It_{App} = 10^4$

$\tau$	Exact	Glaz et al. Product-type	Our Approximation	Approximation Error	Lower Bound	Upper Bound
1	0.810209	0.810216	0.810404	0.001111	0.809903	0.810439
2	0.995764	0.995764	0.995764	$3 \times 10^{-7}$	0.995764	0.995764
3	0.999950	0.999950	0.999950	$4 \times 10^{-11}$	0.999950	0.999950

TABLE 2 :  $n = 5, p = 0.002, m_1 = 5, m_2 = 10, T_1 = 50, T_2 = 80, It_{App} = 10^4$

$\tau$	$\hat{\mathbb{P}}(S \leq \tau)$	Glaz et al. Product-type	Our Approximation	Total Error	Lower Bound	Upper Bound
4	0.894654	0.873256	0.893724	0.037136	0.803422	0.944318
5	0.988003	0.986249	0.988144	0.002125	0.981418	0.993451
6	0.998963	0.998847	0.998963	0.000152	0.998543	0.999401
7	0.999926	0.999919	0.999925	$9 \times 10^{-6}$	0.999903	0.999955
8	0.999995	0.999995	0.999995	$5 \times 10^{-7}$	0.999994	0.999997

TABLE 3 :  $n = 1, p = 0.0001, m_1 = m_2 = m_3 = 5, T_1 = T_2 = T_3 = 60, It_{App} = 10^5$

$\tau$	$\hat{\mathbb{P}}(S \leq \tau)$	Glaz et al. Product-type	Our Approximation	Approximation Error	Simulation Error	Total Error
2	0.993294	0.993241	0.993192	0.000010	0.001367	0.001377
3	0.999963	0.999964	0.999963	0.000000	0.000005	0.000005
4	0.999999	0.999999	0.999999	0.000000	$2 \times 10^{-9}$	$2 \times 10^{-9}$

# EXAMPLES FOR $d = 1, 2, 3$ WHEN $X_{i_1, \dots, i_d} \sim \mathcal{B}(n, p)$

TABLE 1 :  $n = 1, p = 0.005, m_1 = 10, T_1 = 1000, It_{App} = 10^4$

$\tau$	Exact	Glaz et al. Product-type	Our Approximation	Approximation Error	Lower Bound	Upper Bound
1	0.810209	0.810216	0.810404	0.001111	0.809903	0.810439
2	0.995764	0.995764	0.995764	$3 \times 10^{-7}$	0.995764	0.995764
3	0.999950	0.999950	0.999950	$4 \times 10^{-11}$	0.999950	0.999950

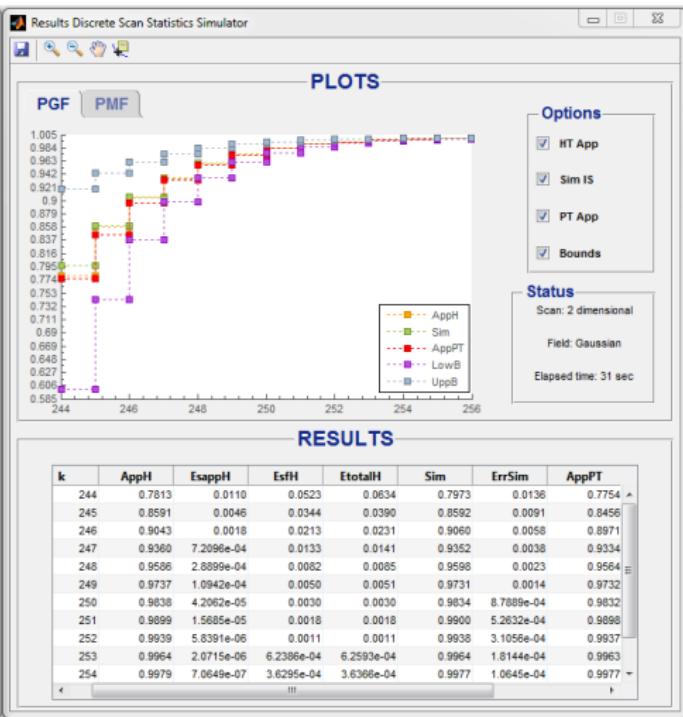
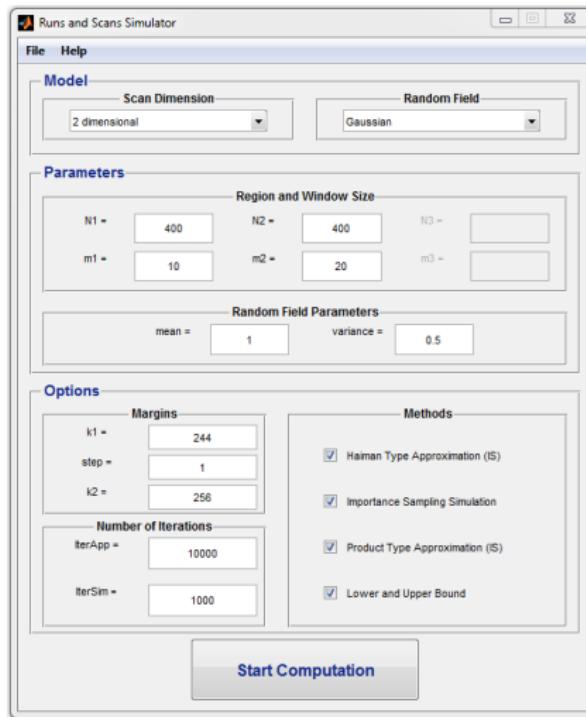
TABLE 2 :  $n = 5, p = 0.002, m_1 = 5, m_2 = 10, T_1 = 50, T_2 = 80, It_{App} = 10^4$

$\tau$	$\hat{\mathbb{P}}(S \leq \tau)$	Glaz et al. Product-type	Our Approximation	Total Error	Lower Bound	Upper Bound
4	0.894654	0.873256	0.893724	0.037136	0.803422	0.944318
5	0.988003	0.986249	0.988144	0.002125	0.981418	0.993451
6	0.998963	0.998847	0.998963	0.000152	0.998543	0.999401
7	0.999926	0.999919	0.999925	$9 \times 10^{-6}$	0.999903	0.999955
8	0.999995	0.999995	0.999995	$5 \times 10^{-7}$	0.999994	0.999997

TABLE 3 :  $n = 1, p = 0.0001, m_1 = m_2 = m_3 = 5, T_1 = T_2 = T_3 = 60, It_{App} = 10^5$

$\tau$	$\hat{\mathbb{P}}(S \leq \tau)$	Glaz et al. Product-type	Our Approximation	Approximation Error	Simulation Error	Total Error
2	0.993294	0.993241	0.993192	0.000010	0.001367	0.001377
3	0.999963	0.999964	0.999963	0.000000	0.000005	0.000005
4	0.999999	0.999999	0.999999	0.000000	$2 \times 10^{-9}$	$2 \times 10^{-9}$

# MATLAB GUI APPLICATION



# OUTLINE

## 1 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (I.I.D. MODEL)

- Framework
- Extremes of 1-dependent stationary sequences
- Scan statistics and 1-dependent sequences
- Simulation methods and computational aspects
- Numerical examples

## 2 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (BLOCK-FACTOR MODEL)

- Model and discussion
- Applications

## 3 CONCLUSIONS AND PERSPECTIVES

- Conclusions
- Perspectives

## 4 REFERENCES

# Model and discussion

# DEFINITION OF A BLOCK-FACTOR

## $k$ BLOCK-FACTOR

The sequence  $(Z_n)_{n \geq 1}$  of random variables with state space  $S_W$  is said to be  $k$  block-factor of the sequence  $(Y_n)_{n \geq 1}$  with state space  $S_Y$  if there is a measurable function  $f : S_Y^k \rightarrow S_W$  such that

$$Z_n = f(Y_n, Y_{n+1}, \dots, Y_{n+k-1}), \forall n \geq 1.$$

## EXAMPLE (2 BLOCK-FACTORS)

- $Z_n = Y_n + Y_{n+1}$ ,  $n \geq 1$  for  $f(x, y) = x + y$
- $Z_n = Y_n Y_{n+1}$ ,  $n \geq 1$  for  $f(x, y) = xy$

## OBSERVATION

If a sequence  $(Z_n)_{n \geq 1}$  of random variables is a  $k$  block-factor, then the sequence is  $(k - 1)$ -dependent.

# INTRODUCING THE MODEL

For each  $1 \leq j \leq d$ ,  $d \geq 1$ , let  $\tilde{T}_j$ ,  $x_1^{(j)}$ ,  $x_2^{(j)}$ ,  $c_j = x_1^{(j)} + x_2^{(j)} + 1$ ,  $T_j = \tilde{T}_j - c_j + 1$  and  $2 \leq m_j \leq T_j$ ,  $1 \leq j \leq d$  be nonnegative integers.

- The rectangular region,  $\tilde{\mathcal{R}}_d = [0, \tilde{T}_1] \times [0, \tilde{T}_2] \times \cdots \times [0, \tilde{T}_d]$
- $\tilde{X}_{s_1, s_2, \dots, s_d}$ ,  $1 \leq s_j \leq \tilde{T}_j$ ,  $j \in \{1, 2, \dots, d\}$  be i.i.d. r.v.'s

To each  $d$ -tuple  $(s_1, \dots, s_d)$ , with  $s_j \in \left\{x_1^{(j)} + 1, \dots, \tilde{T}_j - x_2^{(j)}\right\}$ ,  $j \in \{1, \dots, d\}$ , associate a  $d$ -way tensor  $\mathfrak{X}_{s_1, \dots, s_d} \in \mathbb{R}^{c_1 \times \cdots \times c_d}$

$$\mathfrak{X}_{s_1, \dots, s_d}(j_1, \dots, j_d) = \tilde{X}_{s_1 - x_1^{(1)} - 1 + j_1, \dots, s_d - x_1^{(d)} - 1 + j_d}$$

where  $(j_1, \dots, j_d) \in \{1, \dots, c_1\} \times \cdots \times \{1, \dots, c_d\}$ .

Let  $\Pi : \mathbb{R}^{c_1 \times \cdots \times c_d} \rightarrow \mathbb{R}$  be a measurable real valued function and define, for all  $1 \leq s_j \leq T_j$ ,  $1 \leq j \leq d$ , the *block-factor type* model

$$X_{s_1, \dots, s_d} = \Pi \left( \mathfrak{X}_{s_1 + x_1^{(1)}, \dots, s_d + x_1^{(d)}} \right)$$

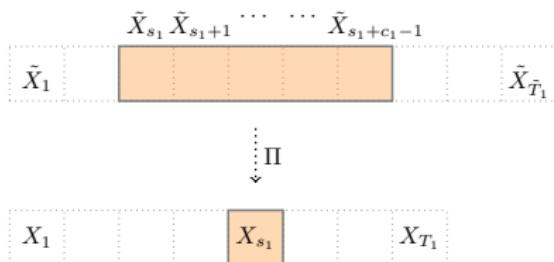
- for  $d = 2$ : [Amărioarei and Preda, 2013b] and [Amărioarei and Preda, 2014]

# EXAMPLES FOR ONE AND TWO DIMENSIONS

EXAMPLE ( $d = 1$ )

$$\mathbf{x}_{\mathbf{s}_1} = \left[ \tilde{X}_{s_1 - x_1^{(1)}}, \dots, \tilde{X}_{s_1 + x_2^{(1)}} \right]$$

$$\mathbf{x}_{\mathbf{s}_1} = \Pi \left( \mathbf{x}_{s_1 + x_1^{(1)}} \right) = \Pi \left( \tilde{X}_{s_1}, \dots, \tilde{X}_{s_1 + c_1 - 1} \right)$$

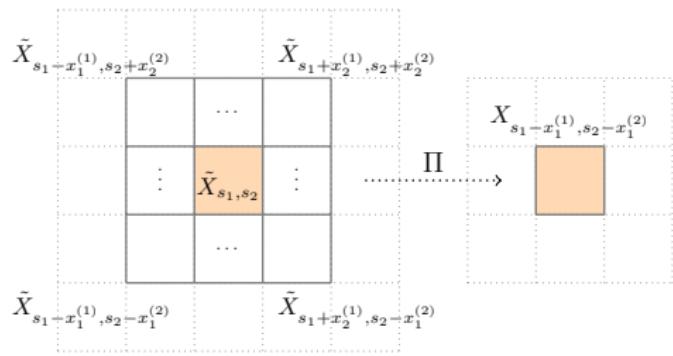


EXAMPLE ( $d = 2$ )

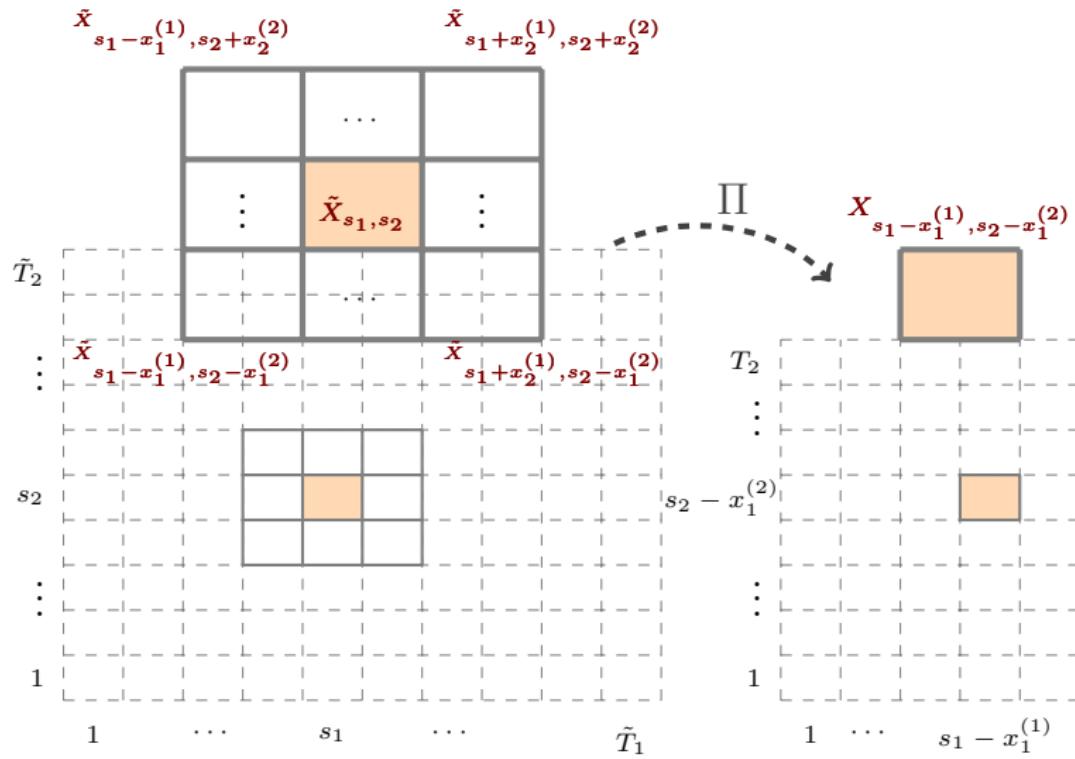
$$\mathbf{x}_{\mathbf{s}_1, \mathbf{s}_2} =$$

$$\begin{pmatrix} \tilde{X}_{s_1 - x_1^{(1)}, s_2 - x_1^{(2)}} & \cdots & \tilde{X}_{s_1 + x_2^{(1)}, s_2 - x_1^{(2)}} \\ \vdots & \ddots & \vdots \\ \tilde{X}_{s_1 - x_1^{(1)}, s_2 + x_2^{(2)}} & \cdots & \tilde{X}_{s_1 + x_2^{(1)}, s_2 + x_2^{(2)}} \end{pmatrix}$$

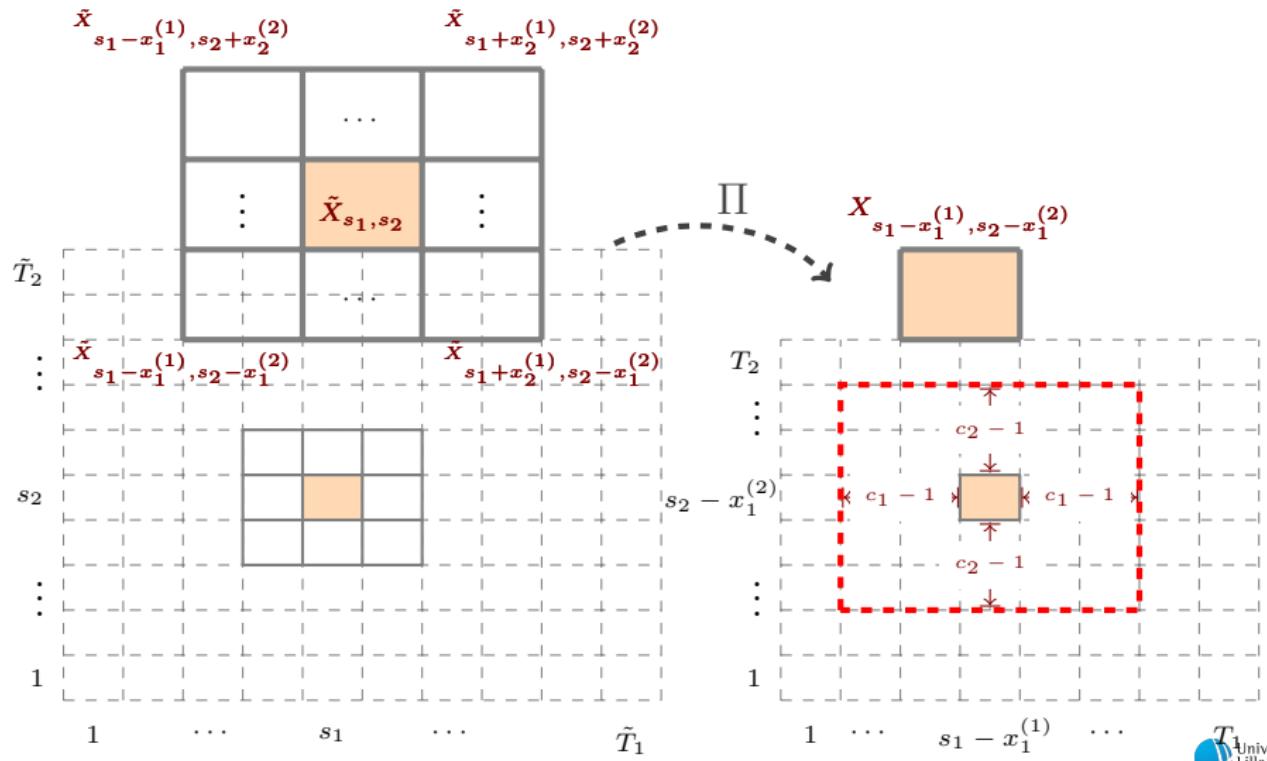
$$\mathbf{x}_{\mathbf{s}_1, \mathbf{s}_2} = \Pi \left( \mathbf{x}_{s_1 + x_1^{(1)}, s_2 + x_1^{(2)}} \right)$$



# DEPENDENCY STRUCTURE ( $d = 2$ )



# DEPENDENCY STRUCTURE ( $d = 2$ )



# APPROXIMATION: IDEA

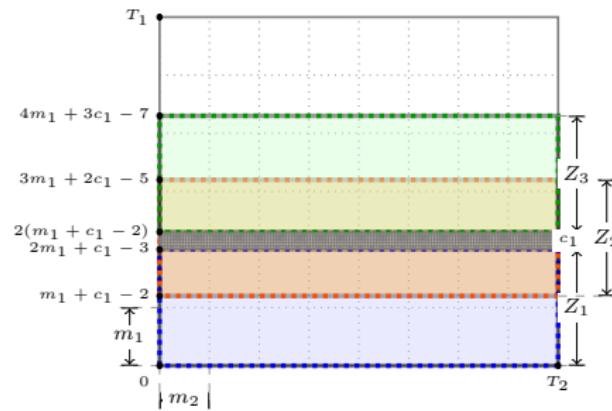
Let  $L_j = \frac{\tilde{T}_j}{m_j + c_j - 2}$ ,  $j \in \{1, 2, \dots, d\}$ , be positive integers

- Define for each  $k_1 \in \{1, 2, \dots, L_1 - 1\}$  the random variables

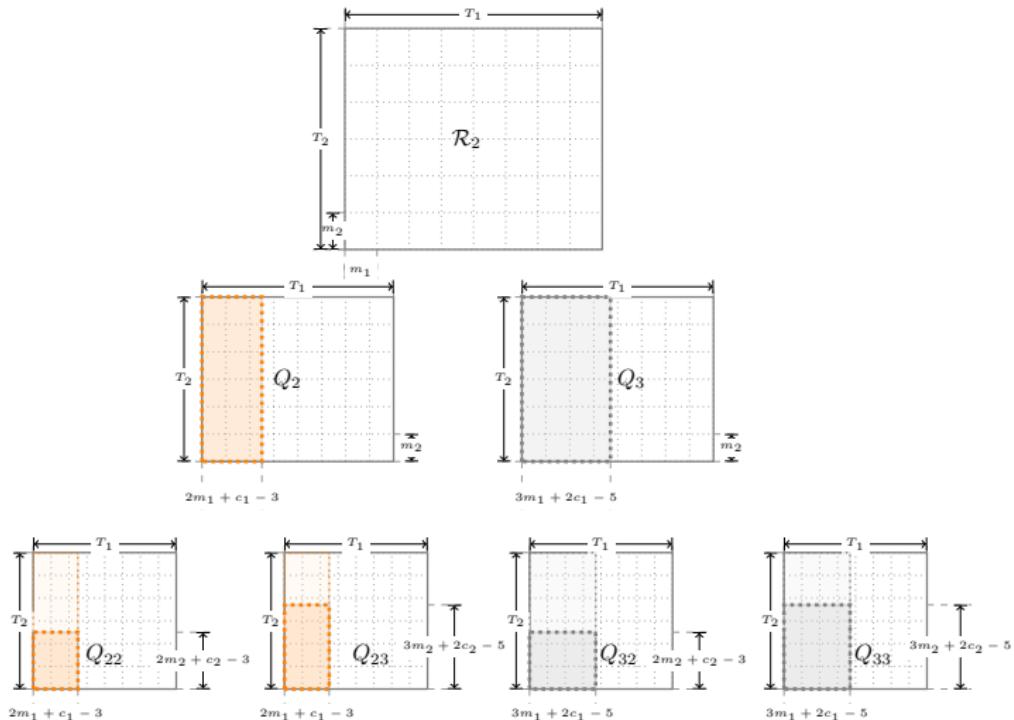
$$Z_{k_1} = \max_{\substack{(k_1-1)(m_1+c_1-2)+1 \leq i_1 \leq k_1(m_1+c_1-2) \\ 1 \leq i_j \leq (L_j-1)(m_j+c_j-2) \\ j \in \{2, \dots, d\}}} Y_{i_1, i_2, \dots, i_d}$$

- $(Z_j)_j$  is 1-dependent, stationary and  $S_m(\mathbf{T}) = \max_{1 \leq k_1 \leq L_1 - 1} Z_{k_1}$

EXAMPLE (1-DEPENDENCE OF  $(Z_j)_j$  FOR  $d = 2$ )



# APPROXIMATION PROCESS ( $d = 2$ )



► Error bounds

# OUTLINE

## 1 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (I.I.D. MODEL)

- Framework
- Extremes of 1-dependent stationary sequences
- Scan statistics and 1-dependent sequences
- Simulation methods and computational aspects
- Numerical examples

## 2 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (BLOCK-FACTOR MODEL)

- Model and discussion
- Applications

## 3 CONCLUSIONS AND PERSPECTIVES

- Conclusions
- Perspectives

## 4 REFERENCES

# Applications

# LONGEST INCREASING RUN

Let  $(\tilde{X}_n)_{n \geq 1}$  be a sequence of i.i.d. r.v.'s with the common distribution  $G$ .

## INCREASING RUN

A subsequence  $(\tilde{X}_k, \dots, \tilde{X}_{k+l-1})$  forms an *increasing run* of length  $l \geq 1$ , starting at position  $k \geq 1$ , if

$$\tilde{X}_{k-1} > \tilde{X}_k < \tilde{X}_{k+1} < \cdots < \tilde{X}_{k+l-1} > \tilde{X}_{k+l}$$

## NOTATIONS

- $M_{\tilde{T}_1}$  = the length of the longest increasing run among the first  $\tilde{T}_1$  r.v.'s
- $L_{\tilde{T}_1}$  = the length of the longest run of ones among the first  $\tilde{T}_1$  r.v.'s

The asymptotic distribution was studied

- $G$  continuous distribution: [Pittel, 1981], [Révész, 1983], [Grill, 1987], [Novak, 1992], etc.
- $G$  discrete distribution: [Csaki and Foldes, 1996], [Grabner et al., 2003], [Eryilmaz, 2006], etc.

# LONGEST INCREASING RUN

## SCAN STATISTICS APPROACH

Let  $d = 1$ ,  $c_1 = 2$ ,  $T_1 = \tilde{T}_1 - 1$  and define  $\Pi : \mathbb{R}^2 \rightarrow \mathbb{R}$  by

$$\Pi(x, y) = \begin{cases} 1, & \text{if } x < y \\ 0, & \text{otherwise} \end{cases}$$

- the block-factor model becomes:  $X_{s_1} = \mathbf{1}_{\tilde{X}_{s_1} < \tilde{X}_{s_1+1}}$

EXAMPLE ( $\tilde{X}_{s_1} \sim \mathcal{U}(0, 1)$ ,  $\tilde{T}_1 = 10$ )

$\tilde{X}_{s_1} : 0.79 \quad 0.31 \quad 0.52 \quad 0.16 \quad 0.60 \quad 0.26 \quad 0.65 \quad 0.68 \quad 0.74 \quad 0.45$

$X_{s_1} :$

We have

$$\mathbb{P}\left(M_{\tilde{T}_1} \leq m_1\right) = \mathbb{P}\left(L_{T_1} < m_1\right) = \mathbb{P}(S_{m_1}(T_1) < m_1), \text{ for } m_1 \geq 1$$

# LONGEST INCREASING RUN

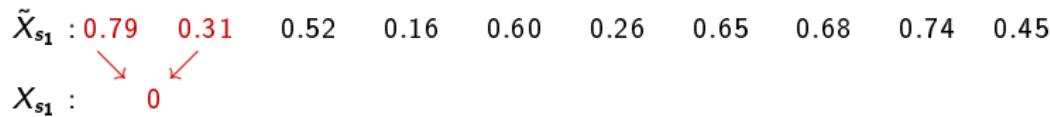
## SCAN STATISTICS APPROACH

Let  $d = 1$ ,  $c_1 = 2$ ,  $T_1 = \tilde{T}_1 - 1$  and define  $\Pi : \mathbb{R}^2 \rightarrow \mathbb{R}$  by

$$\Pi(x, y) = \begin{cases} 1, & \text{if } x < y \\ 0, & \text{otherwise} \end{cases}$$

- the block-factor model becomes:  $X_{s_1} = \mathbf{1}_{\tilde{X}_{s_1} < \tilde{X}_{s_1+1}}$

EXAMPLE ( $\tilde{X}_{s_1} \sim \mathcal{U}(0, 1)$ ,  $\tilde{T}_1 = 10$ )



We have

$$\mathbb{P}\left(M_{\tilde{T}_1} \leq m_1\right) = \mathbb{P}\left(L_{T_1} < m_1\right) = \mathbb{P}(S_{m_1}(T_1) < m_1), \text{ for } m_1 \geq 1$$

# LONGEST INCREASING RUN

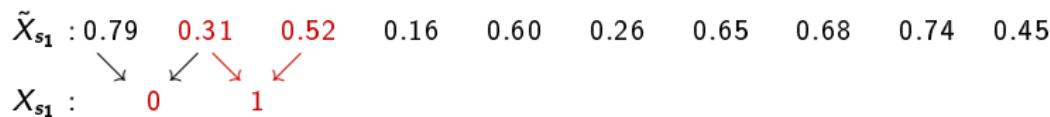
## SCAN STATISTICS APPROACH

Let  $d = 1$ ,  $c_1 = 2$ ,  $T_1 = \tilde{T}_1 - 1$  and define  $\Pi : \mathbb{R}^2 \rightarrow \mathbb{R}$  by

$$\Pi(x, y) = \begin{cases} 1, & \text{if } x < y \\ 0, & \text{otherwise} \end{cases}$$

- the block-factor model becomes:  $X_{s_1} = \mathbf{1}_{\tilde{X}_{s_1} < \tilde{X}_{s_1+1}}$

EXAMPLE ( $\tilde{X}_{s_1} \sim \mathcal{U}(0, 1)$ ,  $\tilde{T}_1 = 10$ )



We have

$$\mathbb{P}\left(M_{\tilde{T}_1} \leq m_1\right) = \mathbb{P}\left(L_{T_1} < m_1\right) = \mathbb{P}(S_{m_1}(T_1) < m_1), \text{ for } m_1 \geq 1$$

# LONGEST INCREASING RUN

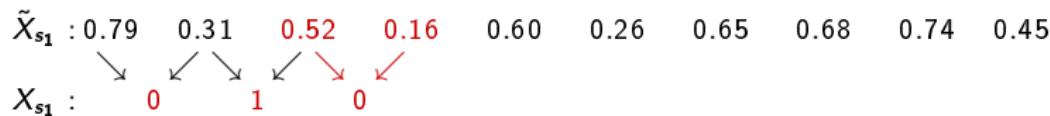
## SCAN STATISTICS APPROACH

Let  $d = 1$ ,  $c_1 = 2$ ,  $T_1 = \tilde{T}_1 - 1$  and define  $\Pi : \mathbb{R}^2 \rightarrow \mathbb{R}$  by

$$\Pi(x, y) = \begin{cases} 1, & \text{if } x < y \\ 0, & \text{otherwise} \end{cases}$$

- the block-factor model becomes:  $X_{s_1} = \mathbf{1}_{\tilde{X}_{s_1} < \tilde{X}_{s_1+1}}$

EXAMPLE ( $\tilde{X}_{s_1} \sim \mathcal{U}(0, 1)$ ,  $\tilde{T}_1 = 10$ )



We have

$$\mathbb{P}\left(M_{\tilde{T}_1} \leq m_1\right) = \mathbb{P}\left(L_{T_1} < m_1\right) = \mathbb{P}(S_{m_1}(T_1) < m_1), \text{ for } m_1 \geq 1$$

# LONGEST INCREASING RUN

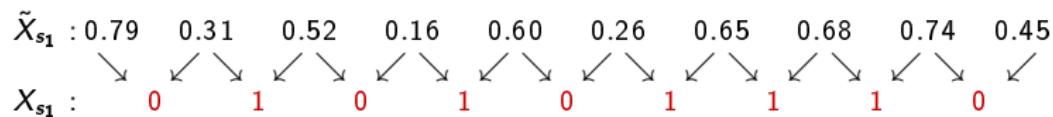
## SCAN STATISTICS APPROACH

Let  $d = 1$ ,  $c_1 = 2$ ,  $T_1 = \tilde{T}_1 - 1$  and define  $\Pi : \mathbb{R}^2 \rightarrow \mathbb{R}$  by

$$\Pi(x, y) = \begin{cases} 1, & \text{if } x < y \\ 0, & \text{otherwise} \end{cases}$$

- the block-factor model becomes:  $X_{s_1} = \mathbf{1}_{\tilde{X}_{s_1} < \tilde{X}_{s_1+1}}$

EXAMPLE ( $\tilde{X}_{s_1} \sim \mathcal{U}(0, 1)$ ,  $\tilde{T}_1 = 10$ )



We have

$$\mathbb{P}\left(M_{\tilde{T}_1} \leq m_1\right) = \mathbb{P}\left(L_{T_1} < m_1\right) = \mathbb{P}(S_{m_1}(T_1) < m_1), \text{ for } m_1 \geq 1$$

# LONGEST INCREASING RUN

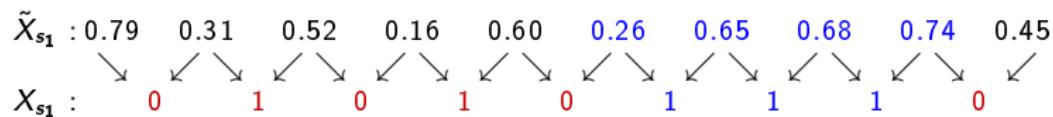
## SCAN STATISTICS APPROACH

Let  $d = 1$ ,  $c_1 = 2$ ,  $T_1 = \tilde{T}_1 - 1$  and define  $\Pi : \mathbb{R}^2 \rightarrow \mathbb{R}$  by

$$\Pi(x, y) = \begin{cases} 1, & \text{if } x < y \\ 0, & \text{otherwise} \end{cases}$$

- the block-factor model becomes:  $X_{s_1} = \mathbf{1}_{\tilde{X}_{s_1} < \tilde{X}_{s_1+1}}$

EXAMPLE ( $\tilde{X}_{s_1} \sim \mathcal{U}(0, 1)$ ,  $\tilde{T}_1 = 10$ )



We have

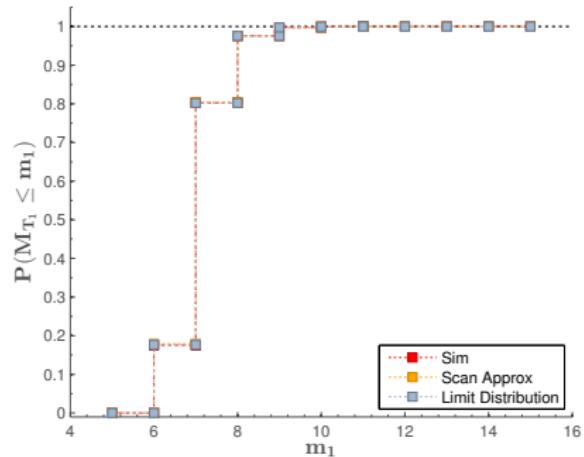
$$\mathbb{P}\left(M_{\tilde{T}_1} \leq m_1\right) = \mathbb{P}\left(L_{T_1} < m_1\right) = \mathbb{P}(S_{m_1}(T_1) < m_1), \text{ for } m_1 \geq 1$$

# LONGEST INCREASING RUN: NUMERICAL RESULTS

For  $\tilde{X}_{s_1} \sim \mathcal{U}([0, 1])$ , [Novak, 1992] showed that

$$\max_{1 \leq m_1 \leq T_1} \left| \mathbb{P}(L_{T_1} < m_1) - e^{-T_1 \frac{m_1+1}{(m_1+2)!}} \right| = \mathcal{O}\left(\frac{\ln T_1}{T_1}\right)$$

$m_1$	Sim	AppH	$E_{total}(1)$	LimApp
5	0.00000700	0.00000733	0.14860299	0.00000676
6	0.17567262	0.17937645	0.01089628	0.17620431
7	0.80257424	0.80362353	0.00110990	0.80215088
8	0.97548510	0.97566460	0.00011579	0.97550345
9	0.99749821	0.99751049	0.00001114	0.99749792
10	0.99977074	0.99977183	0.00000098	0.99977038
11	0.99998075	0.99998083	0.00000008	0.99998073
12	0.99999851	0.99999851	0.00000001	0.99999851
13	0.99999989	0.99999989	0.00000000	0.99999989
14	0.99999999	0.99999999	0.00000000	0.99999999
15	1.00000000	1.00000000	0.00000000	1.00000000

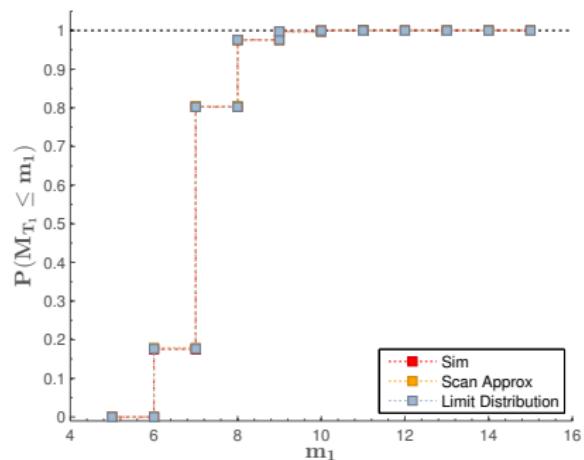


# LONGEST INCREASING RUN: NUMERICAL RESULTS

For  $\tilde{X}_{s_1} \sim \mathcal{U}([0, 1])$ , [Novak, 1992] showed that

$$\max_{1 \leq m_1 \leq T_1} \left| \mathbb{P}(L_{T_1} < m_1) - e^{-T_1 \frac{m_1+1}{(m_1+2)!}} \right| = \mathcal{O}\left(\frac{\ln T_1}{T_1}\right)$$

$m_1$	Sim	AppH	$E_{total}(1)$	LimApp
5	0.00000700	0.00000733	0.14860299	0.00000676
6	0.17567262	0.17937645	0.01089628	0.17620431
7	0.80257424	0.80362353	0.00110990	0.80215088
8	0.97548510	0.97566460	0.00011579	0.97550345
9	0.99749821	0.99751049	0.00001114	0.99749792
10	0.99977074	0.99977183	0.00000098	0.99977038
11	0.99998075	0.99998083	0.00000008	0.99998073
12	0.99999851	0.99999851	0.00000001	0.99999851
13	0.99999989	0.99999989	0.00000000	0.99999989
14	0.99999999	0.99999999	0.00000000	0.99999999
15	1.00000000	1.00000000	0.00000000	1.00000000



# MOVING AVERAGE OF ORDER $q$

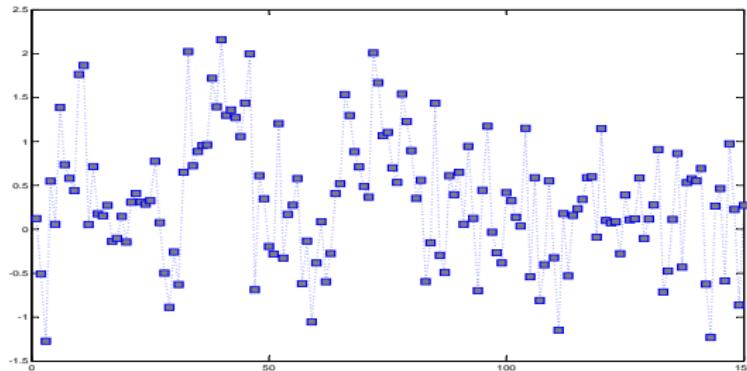
Let  $(\tilde{X}_n)_{n \geq 1}$  be a sequence of i.i.d.  $\mathcal{N}(0, \sigma^2)$  r.v.'s.

## MA( $q$ )

The sequence  $(X_n)_{n \geq 1}$  is said to be an *moving average of order  $q$*  (MA( $q$ )) if

$$X_{s_1} = a_1 \tilde{X}_{s_1} + a_2 \tilde{X}_{s_1+1} + \cdots + a_{q+1} \tilde{X}_{s_1+q}, \quad s_1 \geq 1,$$

and  $(a_1, \dots, a_{q+1}) \in \mathbb{R}^{q+1}$  not all zero.



# MOVING AVERAGE OF ORDER $q$

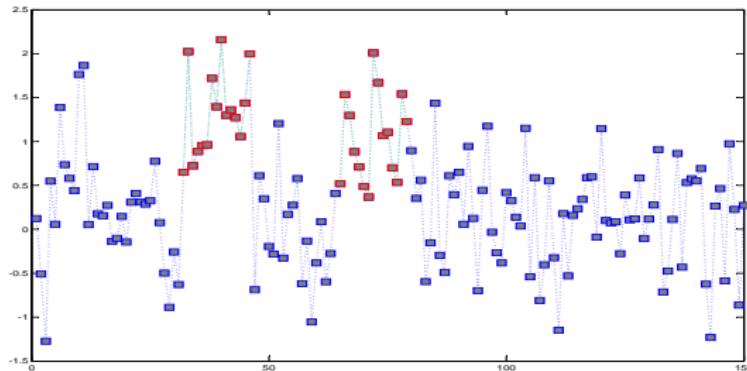
Let  $(\tilde{X}_n)_{n \geq 1}$  be a sequence of i.i.d.  $\mathcal{N}(0, \sigma^2)$  r.v.'s.

## MA( $q$ )

The sequence  $(X_n)_{n \geq 1}$  is said to be an *moving average of order  $q$*  (MA( $q$ )) if

$$X_{s_1} = a_1 \tilde{X}_{s_1} + a_2 \tilde{X}_{s_1+1} + \cdots + a_{q+1} \tilde{X}_{s_1+q}, \quad s_1 \geq 1,$$

and  $(a_1, \dots, a_{q+1}) \in \mathbb{R}^{q+1}$  not all zero.



# MOVING AVERAGE OF ORDER $q$

## SCAN STATISTICS APPROACH

Let  $d = 1$ ,  $x_1^{(1)} = 0$ ,  $x_2^{(1)} = q$  thus  $c_1 = q + 1$ ,  $T_1 = \tilde{T}_1 - q$  and take for  $s_1 \in \{1, \dots, T_1\}$ , the 1-way tensor  $\mathcal{X}_{s_1}$

$$\mathcal{X}_{s_1} = (\tilde{X}_{s_1}, \tilde{X}_{s_1+1}, \dots, \tilde{X}_{s_1+q})$$

and define the block-factor  $\Pi : \mathbb{R}^{q+1} \rightarrow \mathbb{R}$

$$\Pi(x_1, \dots, x_{q+1}) = a_1 x_1 + a_2 x_2 + \dots + a_{q+1} x_{q+1}.$$

## EXAMPLE ( $MA(2)$ )

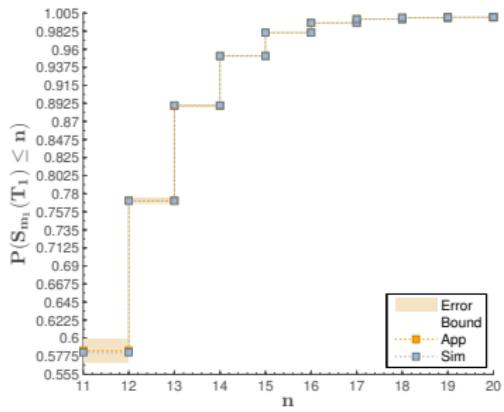
Let  $T_1 = 1000$ ,  $m_1 = 20$ ,  $\tilde{X}_{s_1} \sim \mathcal{N}(0, 1)$  and consider the  $MA(2)$

$$X_{s_1} = 0.3\tilde{X}_{s_1} + 0.1\tilde{X}_{s_1+1} + 0.5\tilde{X}_{s_1+2}$$

- Product-type approximation for  $MA(2)$ : [Wang and Glaz, 2013] and [Wang, 2013]

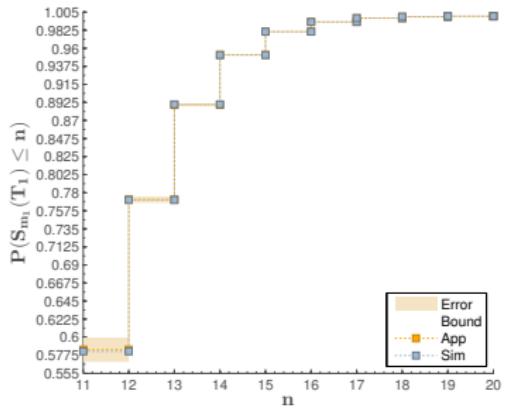
# MOVING AVERAGE OF ORDER $q$ : NUMERICAL RESULTS

$\tau$	Sim	AppPT	AppH	$E_{app}(1)$	$E_{sf}(1)$	$E_{total}(1)$
11	0.582252	0.589479	0.584355	0.011503	0.003653	0.015156
12	0.770971	0.773700	0.771446	0.002319	0.001691	0.004010
13	0.889986	0.890009	0.889431	0.000434	0.000733	0.001167
14	0.951529	0.954536	0.951723	0.000073	0.000297	0.000370
15	0.980653	0.982433	0.980675	0.000011	0.000113	0.000124
16	0.992827	0.993690	0.992791	0.000001	0.000040	0.000042
17	0.997486	0.995471	0.997499	0.000000	0.000013	0.000014
18	0.999186	0.999411	0.999188	0.000000	0.000004	0.000004
19	0.999754	0.999717	0.999754	0.000000	0.000001	0.000001
20	0.999930	1	0.999930	0.000000	0.000000	0.000000



# MOVING AVERAGE OF ORDER $q$ : NUMERICAL RESULTS

$\tau$	Sim	AppPT	AppH	$E_{app}(1)$	$E_{sf}(1)$	$E_{total}(1)$
11	0.582252	0.589479	0.584355	0.011503	0.003653	0.015156
12	0.770971	0.773700	0.771446	0.002319	0.001691	0.004010
13	0.889986	0.890009	0.889431	0.000434	0.000733	0.001167
14	0.951529	0.954536	0.951723	0.000073	0.000297	0.000370
15	0.980653	0.982433	0.980675	0.000011	0.000113	0.000124
16	0.992827	0.993690	0.992791	0.000001	0.000040	0.000042
17	0.997486	0.995471	0.997499	0.000000	0.000013	0.000014
18	0.999186	0.999411	0.999188	0.000000	0.000004	0.000004
19	0.999754	0.999717	0.999754	0.000000	0.000001	0.000001
20	0.999930	1	0.999930	0.000000	0.000000	0.000000



# OUTLINE

## 1 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (I.I.D. MODEL)

- Framework
- Extremes of 1-dependent stationary sequences
- Scan statistics and 1-dependent sequences
- Simulation methods and computational aspects
- Numerical examples

## 2 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (BLOCK-FACTOR MODEL)

- Model and discussion
- Applications

## 3 CONCLUSIONS AND PERSPECTIVES

- Conclusions
- Perspectives

## 4 REFERENCES

# CONCLUSIONS

In this talk:

- improved a result concerning extremes of 1-dependent sequences
- introduced the multidimensional discrete scan statistics
- introduced a new model of dependence based on block-factor constructions
- presented a unified method for estimating the distribution of the multidimensional discrete scan statistics both for the i.i.d model and the block-factor model
- illustrated an importance sampling algorithm that increases the efficiency of the proposed approximation

# OUTLINE

## 1 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (I.I.D. MODEL)

- Framework
- Extremes of 1-dependent stationary sequences
- Scan statistics and 1-dependent sequences
- Simulation methods and computational aspects
- Numerical examples

## 2 MULTIDIMENSIONAL DISCRETE SCAN STATISTICS (BLOCK-FACTOR MODEL)

- Model and discussion
- Applications

## 3 CONCLUSIONS AND PERSPECTIVES

- Conclusions
- Perspectives

## 4 REFERENCES

# FUTURE WORK

Extend the results to

- multidimensional continuous scan statistics
- multidimensional conditional scan statistics

Investigate

- other dependent models
- the influence of the shape of the scanning window
- power of scan statistic based tests under different models
- scan statistics on graphs



Amărioarei, A. (2012).

Approximation for the distribution of extremes of one dependent stationary sequences of random variables.

*arXiv:1211.5456v1, submitted.*



Amărioarei, A. and Preda, C. (2013a).

Approximation for the distribution of three-dimensional discrete scan statistic.

*Methodol Comput Appl Probab.*



Amărioarei, A. and Preda, C. (2013b).

Approximations for two-dimensional discrete scan statistics in some dependent models.

In *Proceedings, 15th Applied Stochastic Models and Data Analysis (ASMDA2013)*.



Amărioarei, A. and Preda, C. (2014).

Approximations for two-dimensional discrete scan statistics in some block-factor type dependent models.

*Journal of Statistical Planning and Inference*, 151-152:107–120.



Boutsikas, M. V. and Koutras, M. V. (2000).

Reliability approximation for Markov chain imbeddable systems.

*Methodol. Comput. Appl. Probab.*, 2:393–411.

-  Boutsikas, M. V. and Koutras, M. V. (2003).  
Bounds for the distribution of two-dimensional binary scan statistics.  
*Probab. Eng. Inform. Sci.*, 17:509–525.
-  Chen, J. and Glaz, J. (1996).  
Two-dimensional discrete scan statistics.  
*Statist. Probab. Lett.*, 31:59–68.
-  Chen, J. and Glaz, J. (1997).  
Approximations and inequalities for the distribution of a scan statistic for 0-1 Bernoulli trials.  
*Advances in the Theory and Practice of Statistics*, 1:285–298.
-  Csaki, E. and Foldes, A. (1996).  
On the length of theh longest monnotone block.  
*Studio Scientiarum Mathematicarum Hungarica*, 31:35–46.
-  Devroye, L. (1986).  
*Non uniform random variate generation*.  
Springer-Verlag, New York.
-  Eryilmaz, S. (2006).

A note on runs of geometrically distributed random variables.

*Discrete Mathematics*, 306:1765–1770.



Fishman, G. (1996).

*Monte Carlo: Concepts, Algorithms and Applications.*

Springer Series in Operations Research. Springer-Verlag, New York.



Frigessi, A. and Vercellis, C. (1984).

An analysis of Monte Carlo algorithms for counting problems.

*Department of Mathematics, University of Milan.*



Fu, J. (2001).

Distribution of the scan statistic for a sequence of bivariate trials.

*J. Appl. Probab.*, 38:908–916.



Gao, T., Ebneshahrashoob, M., and Wu, M. (2005).

An efficient algorithm for exact distribution of discrete scan statistics.

*Methodol. Comput. Appl. Probab.*, 7:1423–1436.



Genz, A. and Bretz, F. (2009).

*Computation of Multivariate Normal and T Probabilities.*

Springer-Verlag, New York.



Glaz, J. (1990).

A comparison of product-type and Bonferroni-type inequalities in presence of dependence.

In *Symposium on Dependence in Probability and Statistics.*, volume 16 of *IMS Lecture Notes-Monograph Series*, pages 223–235. IMS Lecture Notes.



Glaz, J. and Naus, J. (1991).

Tight bounds and approximations for scan statistic probabilities for discrete data.  
*Annals of Applied Probability*, 1:306–318.



Glaz, J., Naus, J., and Wallenstein, S. (2001).

*Scan statistics.*

Springer Series in Statistics. Springer-Verlag, New York.



Grabner, P., Knopfmacher, A., and Prodinger, H. (2003).

Combinatorics of geometrically distributed random variables: run statistics.  
*Theoret. Comput. Sci.*, 297:261–270.



Grill, K. (1987).

Erdos-Révész type bounds for the length of the longest run from a stationary mixing sequence.

*Probab. Theory Relat. Fields*, 75:169–179.



Guerriero, M., Glaz, J., and Sen, R. (2010).

Approximations for a three dimensional scan statistic.

*Methodol. Comput. Appl. Probab.*, 12:731–747.



Haiman, G. (1999).

First passage time for some stationary processes.

*Stochastic Process. Appl.*, 80:231–248.



Haiman, G. (2000).

Estimating the distributions of scan statistics with high precision.

*Extremes*, 3:349–361.



Haiman, G. (2007).

Estimating the distribution of one-dimensional discrete scan statistics viewed as extremes of 1-dependent stationary sequences.

*J. Statist. Plann. Inference*, 137:821–828.



Haiman, G. and Preda, C. (2002).

A new method for estimating the distribution of scan statistics for a two-dimensional Poisson process.

*Methodol. Comput. Appl. Probab.*, 4:393–407.



Haiman, G. and Preda, C. (2006).

Estimation for the distribution of two-dimensional discrete scan statistics.

*Methodol. Comput. Appl. Probab.*, 8:373–381.



Malley, J., Naiman, D. Q., and Bailey-Wilson, J. (2002).

A compressive method for genome scans.

*Human Heredity*, 54:174–185.



Naiman, D. Q. and Priebe, C. E. (2001).

Computing scan statistic  $p$  values using importance sampling, with applications to genetics and medical image analysis.

*J. Comput. Graph. Statist.*, 10:296–328.



Naiman, D. Q. and Wynn, P. (1997).

Abstract tubes, improved inclusion exclusion identities and inequalities and importance sampling.

*The Annals of Statistics*, 25:1954–1983.



Naus, J. (1974).

Probabilities for a generalized birthday problem.

*Journal of American Statistical Association*, 69:810–815.



Naus, J. (1982).

Approximations for distributions of scan statistics.

*Journal of American Statistical Association*, 77:177–183.



Neil, D. (2006).

*Detection of spatial and spatio-temporal clusters.*

PhD thesis, School of Computer Science, Carnegie Mellon University.



Neil, D. (2012).

Fast subset scan for spatial pattern detection.

*Journal of the Royal Statistical Society*, 74(2):337–360.



Novak, S. (1992).

Longest runs in a sequence of  $m$ -dependent random variables.

*Probab. Theory Relat. Fields*, 91:269–281.



Pittel, B. (1981).

Limiting behavior of a process of runs.

*Ann. Probab.*, 9:119–129.



Révész, P. (1983).

Three problems on the length of increasing runs.

*Stochastic Process. Appl.*, 5:169–179.



Shi, J., Siegmund, D., and Yakir, B. (2007).

Importance sampling for estimating  $p$  values in linkage analysis.

*Journal of American Statistical Association*, 102:929–937.



Wang, X. (2013).

Scan statistics for normal data.

PhD thesis, University of Connecticut.



Wang, X. and Glaz, J. (2013).

A variable window scan statistic for  $MA(1)$  process.

In *Proceedings, 15th Applied Stochastic Models and Data Analysis (ASMDA 2013)*, pages 905–912.

thank you!

# PRODUCT-TYPE APPROXIMATIONS

- One dimensional scan statistics

$$\mathbb{P}(S_{m_1}(T_1) \leq \tau) \approx Q(2m_1) \left[ \frac{Q(3m_1)}{Q(2m_1)} \right]^{\frac{T_1}{m_1} - 2},$$

- Two dimensional scan statistics

$$\mathbb{P}(S_{m_1, m_2}(T_1, T_2) \leq \tau) \approx \frac{Q(m_1+1, m_2+1)^{(T_1-m_1)(T_2-m_2)}}{Q(m_1+1, m_2)^{(T_1-m_1)(T_2-m_2-1)}} \times \frac{Q(m_1, 2m_2-1)^{(T_1-m_1-1)(T_2-2m_2)}}{Q(m_1, 2m_2)^{(T_1-m_1-1)(T_2-2m_2+1)}}$$

- Three dimensional scan statistics

$$\begin{aligned} \mathbb{P}(S_{m_1, m_2, m_3}(T_1, T_2, T_3) \leq \tau) &\approx \\ \frac{Q(m_1+1, m_2+1, m_3+1)^{(T_1-m_1)(T_2-m_2)(T_3-m_3)} Q(m_1+1, m_2, m_3)^{(T_1-m_1)(T_2-m_2-1)(T_3-m_3-1)}}{Q(m_1, m_2, m_3)^{(T_1-m_1-1)(T_2-m_2-1)(T_3-m_3-1)} Q(m_1+1, m_2+1, m_3)^{(T_1-m_1)(T_2-m_2)(T_3-m_3-1)}} \times \\ \frac{Q(m_1, m_2+1, m_3)^{(T_1-m_1-1-1)(T_2-m_2)(T_3-m_3-1)} Q(m_1, m_2, m_3+1)^{(T_1-m_1-1)(T_2-m_2-1)(T_3-m_3-1)}}{Q(m_1+1, m_2, m_3+1)^{(T_1-m_1-1)(T_2-m_2-1)(T_3-m_3-3)} Q(m_1, m_2+1, m_3+1)^{(T_1-m_1-1)(T_2-m_2)(T_3-m_3-3)}} \end{aligned}$$

# SELECTED VALUES FOR $K(\cdot)$ AND $\Gamma(\cdot)$

TABLE 4 : Selected values for  $K(\cdot)$  and  $\Gamma(\cdot)$

$1 - q_1$	$K(1 - q_1)$	$\Gamma(1 - q_1)$
0.1	38.63	480.69
0.05	21.28	180.53
0.025	17.56	145.20
0.01	15.92	131.43

◀ Return

# SELECTED VALUES FOR $K(\cdot)$ AND $\Gamma(\cdot)$

TABLE 4 : Selected values for  $K(\cdot)$  and  $\Gamma(\cdot)$

$1 - q_1$	$K(1 - q_1)$	$\Gamma(1 - q_1)$
0.1	38.63	480.69
0.05	21.28	180.53
<b>0.025</b>	<b>17.56</b>	<b>145.20</b>
0.01	15.92	131.43

◀ Return

# ERROR BOUNDS: APPROXIMATION ERROR

## APPROXIMATION ERROR

$$E_{app}(d) = \sum_{s=1}^d (L_1 - 1) \cdots (L_s - 1) \sum_{t_1, \dots, t_{s-1} \in \{2, 3\}} F_{t_1, \dots, t_{s-1}} \left( 1 - \gamma_{t_1, \dots, t_{s-1}, 2} + B_{t_1, \dots, t_{s-1}, 2} \right)^2,$$

where for  $2 \leq s \leq d$

$$F_{t_1, \dots, t_{s-1}} = F(Q_{t_1, \dots, t_{s-1}, 2}, L_s - 1), \quad F = F(Q_2, L_1 - 1),$$

$$B_{t_1, \dots, t_{s-1}} = (L_s - 1) \left[ F_{t_1, \dots, t_{s-1}} \left( 1 - \gamma_{t_1, \dots, t_{s-1}, 2} + B_{t_1, \dots, t_{s-1}, 2} \right)^2 + \sum_{t_s \in \{2, 3\}} B_{t_1, \dots, t_s} \right],$$

$$B_{t_1, \dots, t_{d-1}} = (L_d - 1) F_{t_1, \dots, t_{d-1}} \left( 1 - \gamma_{t_1, \dots, t_{d-1}, 2} + B_{t_1, \dots, t_{d-1}, 2} \right)^2, \quad B_{t_1, \dots, t_d} = 0,$$

and for  $s = 1$ :

$$\sum_{t_1, t_0 \in \{2, 3\}} x = x, \quad F_{t_1, t_0} = F, \quad \gamma_{t_1, t_0, 2} = \gamma_2 \text{ and } B_{t_1, t_0, 2} = B_2.$$

[Return](#)

# ERROR BOUNDS: SIMULATION ERRORS

## SIMULATION ERRORS

$$E_{sf}(d) = (L_1 - 1) \dots (L_d - 1) \sum_{t_1, \dots, t_d \in \{2, 3\}} \beta_{t_1, \dots, t_d}$$

$$\begin{aligned} E_{sapp}(d) = & \sum_{s=1}^d (L_1 - 1) \dots (L_s - 1) \sum_{t_1, \dots, t_{s-1} \in \{2, 3\}} F_{t_1, \dots, t_{s-1}} \left( 1 - \hat{Q}_{t_1, \dots, t_{s-1}, 2} \right. \\ & \left. + A_{t_1, \dots, t_{s-1}, 2} + C_{t_1, \dots, t_{s-1}, 2} \right)^2 \end{aligned}$$

where for  $2 \leq s \leq d$

$$A_{t_1, \dots, t_{s-1}} = (L_s - 1) \dots (L_d - 1) \sum_{t_s, \dots, t_d \in \{2, 3\}} \beta_{t_1, \dots, t_d}, \quad A_{t_1, \dots, t_d} = \beta_{t_1, \dots, t_d}$$

$$\begin{aligned} C_{t_1, \dots, t_{s-1}} = & (L_s - 1) \left[ F_{t_1, \dots, t_{s-1}} \left( 1 - \hat{Q}_{t_1, \dots, t_{s-1}, 2} + A_{t_1, \dots, t_{s-1}, 2} + C_{t_1, \dots, t_{s-1}, 2} \right)^2 \right. \\ & \left. + \sum_{t_s \in \{2, 3\}} C_{t_1, \dots, t_s} \right] \end{aligned}$$

[◀ Return](#)

# DISCRETE SCAN STATISTICS FOR NORMAL DATA

Consider  $d = 1$  and let  $2 \leq m_1 \leq T_1$ ,  $m_1$  and  $T_1$  be positive integers

- $X_{s_1} \sim \mathcal{N}(\mu, \sigma^2)$  are i.i.d.,  $1 \leq s_1 \leq T_1$

The variables  $Y_{i_1} = \sum_{s_1=i_1}^{i_1+m_1-1} X_{s_1}$  follow a multivariate normal distribution with mean  $\bar{\mu} = m_1\mu$  and covariance matrix  $\Sigma = (\Sigma_{i_1,j_1})$

$$\Sigma_{i_1,j_1} = \text{Cov}[Y_{i_1}, Y_{j_1}] = \begin{cases} (m_1 - |i_1 - j_1|) \sigma^2 & , |i_1 - j_1| < m_1 \\ 0 & , \text{otherwise.} \end{cases}$$

[◀ Return](#)

## STEP 2 IN ALGORITHM 2

**Step 2** requires to sample:

- $Y_{i_1^{(k)}}$  from the tail distribution  $\mathbb{P}\left(Y_{i_1^{(k)}} \geq \tau\right)$  ([Devroye, 1986])
- for the other indices, from the conditional distribution given  $\left\{Y_{i_1^{(k)}} \geq \tau\right\}$

For  $\mathbf{W}_1 = \left(Y_1, \dots, Y_{i_1^{(k)}-1}\right)$  and  $\mathbf{W}_2 = \left(Y_{i_1^{(k)}+1}, \dots, Y_{T_1-m_1+1}\right)$

$$\overline{\mathbf{W}}_1 = \mathbf{W}_1 | (Y_{i_1^{(k)}} = t) \sim \mathcal{N}(\mu_{w_1|t}, \Sigma_{w_1|t}) \text{ and } \overline{\mathbf{W}}_2 = \mathbf{W}_2 | (Y_{i_1^{(k)}} = t) \sim \mathcal{N}(\mu_{w_2|t}, \Sigma_{w_2|t})$$

where for  $i \in \{1, 2\}$ ,

$$\mu_{w_i|t} = \mathbb{E}[\mathbf{W}_i] + \frac{1}{Var[Y_{i_1^{(k)}}]} Cov[\mathbf{W}_i, Y_{i_1^{(k)}}](t - \mathbb{E}[Y_{i_1^{(k)}}]),$$

$$\Sigma_{w_i|t} = Cov(\mathbf{W}_i) - \frac{1}{Var[Y_{i_1^{(k)}}]} Cov[\mathbf{W}_i, Y_{i_1^{(k)}}] Cov^T[\mathbf{W}_i, Y_{i_1^{(k)}}].$$

[◀ Return](#)

# CUMULATIVE COUNTS METHOD

## IDEA

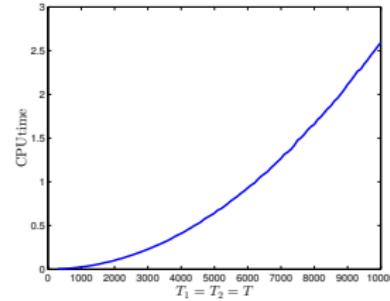
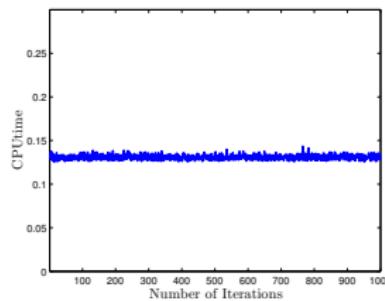
Precompute a matrix of cumulative counts  $M$  using dynamic programming and express the variables of interest as differences.

- efficiently searches for the locality statistics over  $\mathcal{R}_d$  in constant time

EXAMPLE ( $d = 2$ ,  $T_1 = T_2 = T$ ,  $m_1 = m_2 = m$ )

The matrix  $M$  has the entries  $M(i, j) = \sum_{k=1}^i \sum_{l=1}^j X_{k,l}$ , so the locality statistic is

$$Y_{i_1, i_2} = M(i_1 + m - 1, i_2 + m - 1) - M(i_1 + m - 1, i_2 - 1) - M(i_1 - 1, i_2 + m - 1) + M(i_1 - 1, i_2 - 1)$$



## ALTERNATIVE APPROACHES

Several other methods were proposed:

- I) [Genz and Bretz, 2009] developed a quasi Monte Carlo algorithm for numerically approximate the distribution of a multivariate normal, the algorithm was implemented in R and Matlab ([Wang and Glaz, 2013], [Wang, 2013])
- II) [Shi et al., 2007] introduced another IS algorithm (Algo 3)
  - idea: imbed the probability measure under  $H_0$  into an exponential family

► Details Algo 3

To measure the efficiency of the methods we evaluate the *relative efficiency* introduced by [Malley et al., 2002]

$$\text{Rel Eff} = \frac{\sigma_{\text{method 1}}^2 \times \text{CPU Time}_{\text{method 1}}}{\sigma_{\text{method 2}}^2 \times \text{CPU Time}_{\text{method 2}}}$$

# IS ALGORITHM [SHI ET AL., 2007]

## Algorithm 3 Second Importance Sampling Algorithm for Scan Statistics

Take  $d\mathbb{P}_{\xi, \mathbf{r}_1} = \frac{e^{\xi Y_{\mathbf{r}_1}}}{\mathbb{E}_{H_0}[e^{\xi Y_{\mathbf{r}_1}}]} d\mathbb{P}_{H_0}$  and compute

$$\xi \approx \frac{\tau}{m_1 \sigma^2} - \frac{\mu}{\sigma^2}, \quad \mathbb{E}_{\xi, \mathbf{r}_1} [Y_{\mathbf{i}_1}] = \xi \text{Cov}_{H_0} [Y_{\mathbf{i}_1}, Y_{\mathbf{r}_1}] + m_1 \mu, \quad \text{Cov}_{\xi, \mathbf{r}_1} [Y_{\mathbf{i}_1}, Y_{\mathbf{j}_1}] = \text{Cov}_{H_0} [Y_{\mathbf{i}_1}, Y_{\mathbf{j}_1}]$$

Repeat for each  $k$  from 1 to  $ITER$  (iterations number)

- 1: Generate uniformly  $i_1^{(k)}$  from the set  $\{1, \dots, T_1 - m_1 + 1\}$ .
- 2: Given  $i_1^{(k)}$ , generate the Gaussian process  $Y_{\mathbf{i}_1}$  according to the new measure  $d\mathbb{P}_{\xi, i_1^{(k)}}$ .
- 3: Compute  $\hat{\rho}_k(1)$  based on

$$\hat{\rho}_k(1) = \sum_{j_1=1}^{T_1-m_1+1} e^{\xi Y_{j_1} - m_1 (\mu \xi + \frac{\sigma^2 \xi^2}{2})} \mathbf{1}_{\{S_{m_1}(T_1) \geq \tau\}}$$

End Repeat  
Return

$$\bar{\rho}(1) = \frac{1}{ITER} \sum_{k=1}^{ITER} \hat{\rho}_k(1), \quad \text{Var} [\hat{\rho}(1)] \approx \frac{1}{ITER-1} \sum_{k=1}^{ITER} \left( \hat{\rho}_k(1) - \frac{1}{ITER} \sum_{k=1}^{ITER} \hat{\rho}_k(1) \right)^2$$

[◀ Return](#)

# NUMERICAL RESULTS

All the results are compared with respect to Algo 2 for  $ITER = 10000$

TABLE 5 : Algorithm [Genz and Bretz, 2009], IS (Algo 2) and the relative efficiency (Rel Eff)

$T_1$	$m_1$	$\tau$	Genz	Err Genz	IS Algo 2	Err Algo 2	Rel Eff
200	15	12	0.932483	0.000732	0.933215	0.000743	7
500	25	18	0.976117	0.000460	0.975797	0.000425	518
750	30	24	0.998454	0.000125	0.998493	0.000024	688
800	40	30	0.999752	0.000029	0.999742	0.000004	617

TABLE 6 : Naive Monte Carlo (MC), IS (Algo 2) and the relative efficiency (Rel Eff)

$T_1$	$m_1$	$\tau$	MC	Err MC	IS Algo 2	Err Algo 2	Rel Eff
200	15	12	0.932624	0.000694	0.933215	0.000743	15
500	25	18	0.975880	0.000425	0.975797	0.000425	33
750	30	24	0.998515	0.000061	0.998493	0.000024	101
800	40	30	0.999741	0.000009	0.999742	0.000004	602

# NUMERICAL RESULTS

TABLE 7 : IS algorithms (Algo 2 and Algo 2) and the relative efficiency (Rel Eff)

$T_1$	$m_1$	$\tau$	IS Algo 2	Err Algo 2	IS Algo 2	Err Algo 2	Rel Eff
200	15	12	0.932744	0.000839	0.933215	0.000743	3
500	25	18	0.976105	0.000448	0.975797	0.000425	3.5
750	30	24	0.998508	0.000032	0.998493	0.000024	3.5
800	40	30	0.999740	0.000006	0.999742	0.000004	3.6

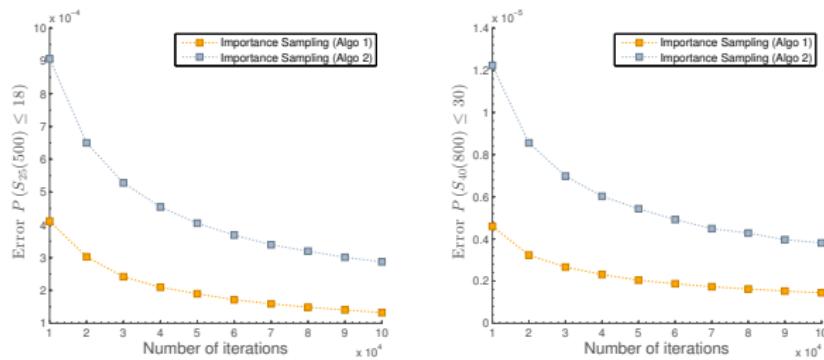


FIGURE 1 : The evolution of simulation error in IS Algorithm 2 and IS Algorithm 2

# ERROR BOUNDS: APPROXIMATION ERROR

## APPROXIMATION ERROR

$$E_{app}(d) = \sum_{s=1}^d (L_1 - 1) \cdots (L_s - 1) \sum_{t_1, \dots, t_{s-1} \in \{2, 3\}} F_{t_1, \dots, t_{s-1}} \left( 1 - \gamma_{t_1, \dots, t_{s-1}, 2} + B_{t_1, \dots, t_{s-1}, 2} \right)^2,$$

where for  $2 \leq s \leq d$

$$F_{t_1, \dots, t_{s-1}} = F(Q_{t_1, \dots, t_{s-1}, 2}, L_s - 1), \quad F = F(Q_2, L_1 - 1),$$

$$B_{t_1, \dots, t_{s-1}} = (L_s - 1) \left[ F_{t_1, \dots, t_{s-1}} \left( 1 - \gamma_{t_1, \dots, t_{s-1}, 2} + B_{t_1, \dots, t_{s-1}, 2} \right)^2 + \sum_{t_s \in \{2, 3\}} B_{t_1, \dots, t_s} \right],$$

$$B_{t_1, \dots, t_{d-1}} = (L_d - 1) F_{t_1, \dots, t_{d-1}} \left( 1 - \gamma_{t_1, \dots, t_{d-1}, 2} + B_{t_1, \dots, t_{d-1}, 2} \right)^2, \quad B_{t_1, \dots, t_d} = 0,$$

and for  $s = 1$ :

$$\sum_{t_1, t_0 \in \{2, 3\}} x = x, \quad F_{t_1, t_0} = F, \quad \gamma_{t_1, t_0, 2} = \gamma_2 \text{ and } B_{t_1, t_0, 2} = B_2.$$

[Return](#)

# ERROR BOUNDS: SIMULATION ERRORS

## SIMULATION ERRORS

$$E_{sf}(d) = (L_1 - 1) \dots (L_d - 1) \sum_{t_1, \dots, t_d \in \{2, 3\}} \beta_{t_1, \dots, t_d}$$

$$E_{sapp}(d) = \sum_{s=1}^d (L_1 - 1) \dots (L_s - 1) \sum_{t_1, \dots, t_{s-1} \in \{2, 3\}} F_{t_1, \dots, t_{s-1}} \left( 1 - \hat{Q}_{t_1, \dots, t_{s-1}, 2} \right. \\ \left. + A_{t_1, \dots, t_{s-1}, 2} + C_{t_1, \dots, t_{s-1}, 2} \right)^2$$

where for  $2 \leq s \leq d$

$$A_{t_1, \dots, t_{s-1}} = (L_s - 1) \dots (L_d - 1) \sum_{t_s, \dots, t_d \in \{2, 3\}} \beta_{t_1, \dots, t_d}, \quad A_{t_1, \dots, t_d} = \beta_{t_1, \dots, t_d}$$

$$C_{t_1, \dots, t_{s-1}} = (L_s - 1) \left[ F_{t_1, \dots, t_{s-1}} \left( 1 - \hat{Q}_{t_1, \dots, t_{s-1}, 2} + A_{t_1, \dots, t_{s-1}, 2} + C_{t_1, \dots, t_{s-1}, 2} \right)^2 \right. \\ \left. + \sum_{t_s \in \{2, 3\}} C_{t_1, \dots, t_s} \right]$$

[Return](#)