

The Discrete New XLindley Distribution And The Associated Autoregressive Process

R. Maya^a, P. Jodrá^b, S. Aswathy^c and M. R. Irshad^c¹

^a Department of Statistics, University College, Thiruvananthapuram 695034, Kerala, India.

^b Departamento de Métodos Estadísticos, EINA, Universidad de Zaragoza, 50018 Zaragoza, Spain.;
Email: pjodra@unizar.es

^c Department of Statistics, Cochin University of Science and Technology, Cochin 682 022, Kerala, India.; Email: aswathysree22@cusat.ac.in and irshadmr@cusat.ac.in

Abstract

The continuous new XLindley distribution was introduced by Nawel et al. (2023) as a special case of the polynomial exponential distribution proposed by Beghriche et al. (2022). The current paper introduces the one-parameter discrete analogue distribution of the new XLindley model and studies its main statistical properties. In particular, closed-form expressions are provided for the moment generating function, mean, variance, quantile function, hazard rate function and mean residual life. Moreover, the new distribution has discrete increasing failure rate and both overdispersed and underdispersed count data can be handled. The estimation of the unknown parameter can be performed by the maximum likelihood method and a Monte Carlo simulation study reveals that this method provides satisfactory estimates. Additionally, a first-order integer-valued autoregressive process is constructed from the discrete distribution and, via a simulation study, the conditional maximum likelihood method is recommended for estimation purposes. In order to assess the usefulness in practical applications, the proposed distribution and the associated first-order autoregressive process are compared to other competing distributions and processes, using to this end several real data sets. In the context of statistical quality control, finally a cumulative sum control chart is developed for monitoring the process mean. To illustrate its usefulness, both simulation and real data analysis are performed.

Keywords: New XLindley distribution, discrete distribution, Lambert W function, IFR property, INAR(1) process, data analysis, CUSUM chart.

1. Introduction

Over the last years, there has been a considerable development of new discrete probability distributions which are used to model count data in fields such as biology, epidemiology, finance and social sciences, for example, modelling the number of species in an ecological community, the number of customer purchases, the number of accidents in a region, the number of votes received by the candidates in an election, among other situations. Generally, count data represent the number of random events that have occurred in a specific time interval and frequently may be interpreted as outcomes of an underlying count process in continuous time. Moreover, in order to simplify their

¹Corresponding author. Email: irshadmr@cusat.ac.in

analysis and interpretation, it is interesting to point out that in many applications the discretization of continuous data is considered more suitable than examining the data on a real scale. Such data are usually non-negative integer values, including zero counts, without an inherent upper bound and, accordingly, these data are modelled by discrete probability distributions.

The classical discrete models such as the geometric, negative binomial and Poisson distributions have been traditionally used to model count data but it is well-known that these distributions have some limitations. For example, the Poisson distribution models count data showing equidispersion but count data frequently also exhibit overdispersion or underdispersion, which has led to the development of more flexible models during the last decades. In this regard, different methods of generating discrete probability distributions as analogues of continuous probability distributions have been introduced in the statistical literature, such as the infinite series discretization method (Good (1953), Kulasekera and Tonkyn (1992), Sato et al. (1999)), the survival discretization approach (Nakagawa and Osaki (1975)), the reversed hazard function discretization method (Ghosh et al. (2013)) and the compound two-phase method (Chakraborty (2015)), among others. The reader is referred to Chakraborty (2015) for a survey on discretization methods of continuous distributions, where their main differences and implications can be found.

One of the most popular discretization techniques and possibly the easiest method of construction is the survival discretization approach. More precisely, assuming that X is a continuous random variable with survival function $S(x) = P(X > x)$, the probability mass function (pmf) of its analogue discrete random variable X with support the set of non-negative integer numbers is given by

$$P(X = x) = S(x) - S(x + 1), \quad x = 0, 1, 2, \dots$$

This procedure provides a discrete distribution that retains the functional shape of the survival function corresponding to the continuous distribution. Some examples of distributions obtained by this approach are the discrete Weibull (Nakagawa and Osaki (1975)), discrete Pareto (Krishna and Pundir (2009)), discrete generalized Pareto (Haj Ahmad and Almetwally (2022)), discrete inverse Weibull (Jazi et al. (2010)), discrete Lindley (Gómez-Déniz and Calderín-Ojeda (2011)) (see also (Al-Babtain et al. (2021))), discrete Burr-Hatke (El-Morshedy et al. (2020)), discrete pseudo Lindley (Irshad et al. (2021)), discrete Bilal (Altun et al. (2022)), discrete half-logistic (Teamah (2024)), discrete Teissier (Irshad et al. (2023)), discrete Marshall–Olkin length biased exponential (Aljohani et al. (2023)), discrete Marshall–Olkin inverted Topp–Leone (Almetwally (2022)) and discrete alpha power inverse Weibull (Alotaibi et al. (2023)) among others.

Furthermore, count data are also used to model time series when the frequency of events is observed over consecutive time intervals. In this context, the well-known integer-valued first-order autoregressive process INAR(1) is very popular for modelling count time series in many applied contexts, for example, analysing the number of insurance claims or the number of daily reported cases of an infectious disease. The INAR(1) process is suitable for counting processes in which an element of the process at time t can be either the survival of an element of the process at time $t - 1$ or the outcome of an innovation process. Pioneer works dealing with an INAR(1) process with Poisson innovations were given by McKenzie (1985) and Al-Osh and Alzaid (1987). Subsequently, in order to model both overdispersed and underdispersed count data, several INAR(1) models with different innovation processes have been introduced in the statistical literature, such as

the INAR(1)G with geometric innovations ([Aghababaei Jazi et al. \(2022\)](#)), INAR(1)PQX with Poisson quasi-xgamma innovations ([Altun et al. \(2021\)](#)), INAR(1)BL with Bell innovations ([Huang and Zhu \(2021\)](#)) and INAR(1)DB with Bilal innovations ([Altun et al. \(2022\)](#)), among others.

The utilization of count-type data becomes prevalent across various sectors, offering a viable method for evaluating performance when faced with challenges in accurately measuring certain quality characteristics on a numerical scale. For example, in practical scenarios, consider situations where statistical quality-control applications rely on data representation involving the count of defects in manufactured items or the frequency of service errors within designated time frames. In this scenario, control charts are used in quality control to detect deviations or variations in a process that could affect the quality of the output and, accordingly, effective monitoring is crucial for promptly detecting changes in the mean. In statistical process control, the cumulative sum (CUSUM) chart is a tool widely employed to identify abnormal changes in the mean within a manufacturing process. Traditional control charts often assume the independence of observations, an assumption that may not always hold in practical applications. Hence, the use of control charts becomes necessary in models for integer-valued time series data. In this regard, see the CUSUM charts for monitoring the mean of an INAR(1) process with Poisson marginals ([Weiss and Testik \(2009\)](#)), of an INAR(1) process with geometrically inflated Poisson innovations ([Li et al. \(2022\)](#)), of an INAR(1) process with Katz family innovations ([Kim and Lee \(2017\)](#)) and also for monitoring correlated Poisson counts with an excessive number of zeros ([Rakitzis et al. \(2017\)](#)).

Recently, a novel family of distributions has been introduced in [Beghriche et al. \(2022\)](#), the so-called one-parameter polynomial exponential distribution. It is noteworthy that both the XLindley distribution [Chouia \(2021\)](#) and the new XLindley distribution ([Nawel et al. \(2023\)](#)) belong to this emerging family. In addition, some variations of the XLindley distribution have also been proposed, such as the inverse XLindley distribution ([Beghriche et al. \(2023\)](#)) and the discrete XLindley distribution ([Eldeeb et al. \(2023\)](#)). With respect to the new XLindley distribution, it constitutes a novel one-parameter distribution that combines the Lindley and exponential distributions with potential applications across various fields (see [Nawel et al. \(2023\)](#)) and, as far as we know, no additional versions of this distribution have been introduced in the statistical literature. In the current paper, a discrete variant of the new XLindley distribution is proposed. It is interesting to note that its pmf involves only the exponential function which implies that its main statistical properties can be expressed in closed form, by contrast to the aforementioned discrete XLindley distribution whose pmf depends of the logarithm function and this fact adds complexity to the model. Among others, a rationale for choosing the continuous new XLindley distribution for discretization lies in its unique combination of features from both the Lindley and exponential distributions. In particular, this distribution shows a heightened risk rate along with a decreasing average residual life function. As a novel one-parameter distribution, it offers potential benefits across various fields including biology, engineering, astronomy, actuarial science, and medicine.

The aim of this paper is to study the new one-parameter discrete distribution obtained from the continuous new XLindley (NXL) distribution as well as the associated INAR(1) process. Moreover, a CUSUM control chart is also developed for monitoring the mean of the proposed process. To be more precise, the probability density function (pdf) and the cumulative distribution function (cdf) of the NXL distribution with parameter $\theta > 0$ are

given, respectively, by

$$f(x) = \frac{\theta}{2} (1 + \theta x) e^{-\theta x}, \quad x > 0$$

and

$$F(x) = 1 - \left(1 + \frac{1}{2} \theta x\right) e^{-\theta x}, \quad x > 0. \quad (1.1)$$

In particular, the discrete analogue of the above NXL distribution will be obtained by applying the survival discretization method and will be called the discrete new XLindley (DNXL) distribution.

Before going further, it is important to note that the DNXL distribution should not be confused with the discrete XLindley distribution recently proposed by [Eldeeb et al. \(2023\)](#). It must also be remarked that the DNXL distribution can be characterized as a particular case of the discrete Pseudo Lindley (DPsL) distribution introduced by [Irshad et al. \(2021\)](#). The main motivations to study this special case are the following.

- The DNXL distribution maintains similar statistical properties using only one parameter and, therefore, constitutes a parsimonious alternative to the DPsL model.
- The DNXL distribution has the capability to handle both underdispersed and overdispersed count data with a single parameter and also has an increasing hazard rate function.
- The estimation of parameters of the DNXL distribution and its associated INAR(1) process is simpler. For example, the conditional least squares method does not provide acceptable estimates for an INAR(1) process with DPsL innovations (see [Irshad et al. \(2021\)](#)) and, by contrast, as will be seen, it is the recommended estimated method for an INAR(1) process with DNXL innovations.
- The simplicity of having just one parameter makes the development of the CUSUM chart for detecting increasing shifts in the process mean levels more straightforward.

The remainder of this paper is organized as follows. Section 2 introduces the DNXL distribution and provides some statistical properties, emphasizing the statistical measures that can be expressed analytically. The parameter estimation is carried out in Section 3 together with a Monte Carlo simulation study to assess the performance of different estimation methods. Section 4 introduces an INAR(1) process with DNXL innovations and a simulation study is conducted to choose the most suitable method to estimate its parameters. The usefulness of the DNXL distribution and the associated INAR(1) process in practical applications are illustrated in Section 5, where different discrete distributions and INAR(1) processes are fitted to several real data sets. Section 6 studies a CUSUM control chart to detect mean increases and its effectiveness is evaluated via numerical simulation and real data analysis. Finally, Section 7 summarizes the main conclusions together with future research.

2. The discrete new XLindley distribution

In this section, the discrete analogue of the NXL distribution is studied together with its main properties. It should be noted that notations and some results are similar to those in [Irshad et al. \(2021\)](#).

2.1. Definition

The DNXL distribution is obtained by applying the survival discretization method to the NXL distribution defined by Eq. (1.1). As result, the pmf of a random variable X having DNXL distribution with parameter $\theta > 0$ is given by

$$p(x; \theta) = \frac{1}{2} e^{-x\theta} \left[(2 + \theta x)(1 - e^{-\theta}) - \theta e^{-\theta} \right], \quad x = 0, 1, 2, \dots \quad (2.1)$$

Consequently, the cdf and survival function of X are given, respectively, by

$$F(x; \theta) = 1 - \frac{1}{2} e^{-(x+1)\theta} [2 + (x+1)\theta], \quad x = 0, 1, 2, \dots \quad (2.2)$$

and

$$S(x; \theta) = \frac{1}{2} e^{-(x+1)\theta} [2 + (x+1)\theta], \quad x = 0, 1, 2, \dots$$

Hereafter, a random variable with pmf given by Eq. (2.1) will be referred to as a random variable having DNXL distribution with parameter $\theta > 0$. Figure 1 represents pmf (2.1) for different values of θ . The figure suggests that the DNXL distribution is unimodal and exhibits a right-skewed pattern.

2.2. Mode

The unimodality of a random variable having DNXL distribution is proved in the next results.

Proposition 2.1

The pmf of a DNXL distribution is log-concave for any $\theta > 0$.

Proof. The pmf of a non-negative discrete random variable is said to be log-concave if satisfies the inequality $p(x; \theta)^2 \geq p(x-1; \theta)p(x+1; \theta)$ for $x = 0, 1, 2, \dots$. Taking into account Eq. (2.1), it can be checked the following

$$p(x; \theta)^2 - p(x-1; \theta)p(x+1; \theta) = \frac{1}{4} \theta^2 e^{-2\theta x} \left[1 + e^{-\theta}(e^{-\theta} - 2) \right] \geq 0,$$

which holds for any $\theta > 0$ and implies the stated result. \square

Corollary 2.1

The DNXL distribution is unimodal.

Proof. It is well-known that a log-concave pmf is strongly unimodal (cf. Keilson and Gerber (1971)). Therefore, by virtue of Proposition 2.1 the result follows by applying Theorem 3 in Keilson and Gerber (1971). \square

Corollary 2.2

Let X be a random variable having DNXL distribution with parameter $\theta > 0$. The mode of X is $x = 0$.

Proof. The DNXL distribution is unimodal by virtue of Corollary 2.1 and additionally it can be proved that $p(x+1; \theta) - p(x; \theta) < 0$ for $x = 0, 1, \dots$. To see this, from Eq. (2.1) note that the above difference can be written as follows

$$p(x+1; \theta) - p(x; \theta) = -\frac{e^\theta - 1}{2e^{2\theta}e^{\theta x}} \left[(2 + \theta x)e^\theta - (2 + x)\theta - 2 \right]$$

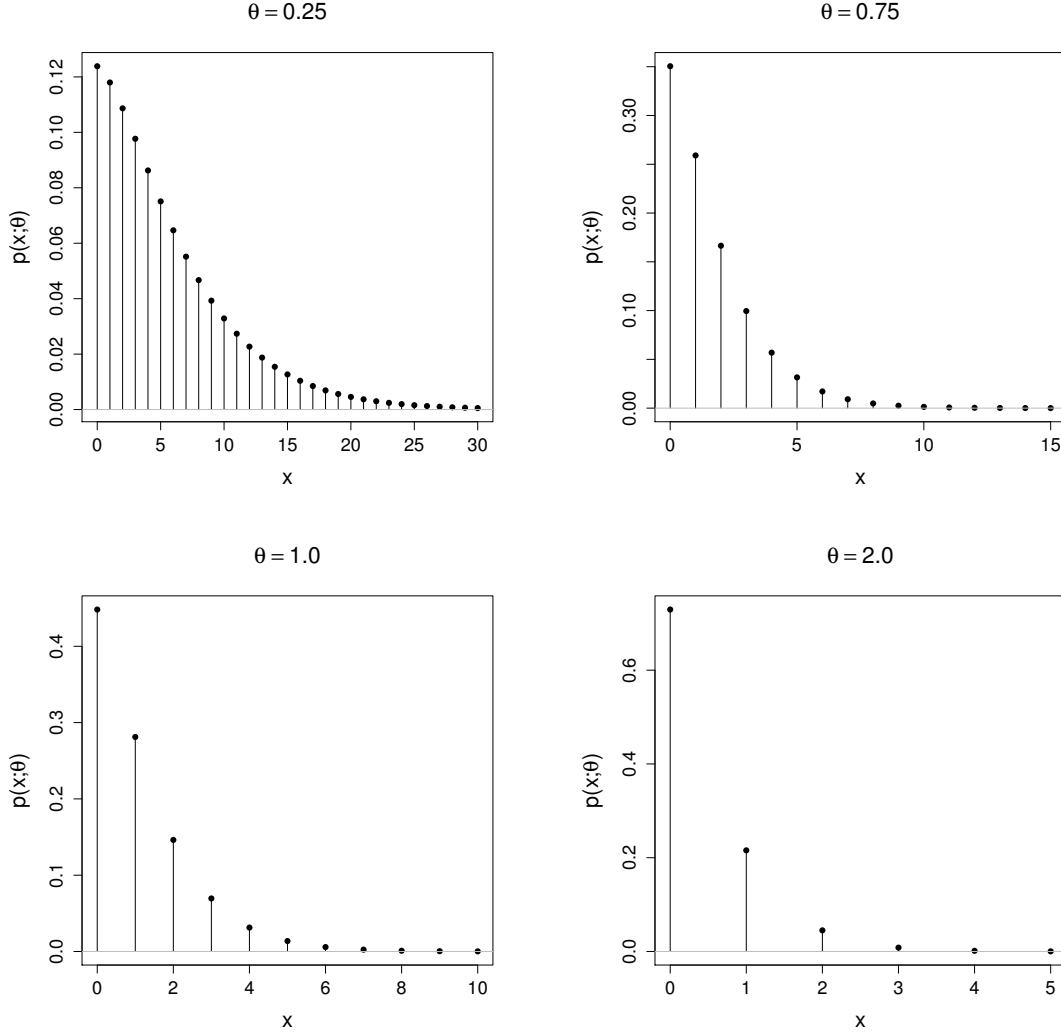


Figure 1: Pmf of the DNXL distribution for different values of θ .

and it is easy to check the inequality $(2 + \theta x)e^\theta - (2 + x)\theta - 2 > 0$ for any $\theta > 0$ and $x = 0, 1, \dots$, which implies the result. \square

2.3. Moment generating function, moments, skewness and kurtosis coefficients

The moment generating function of a random variable having DNXL distribution is given in closed form in the next result, which is obtained from routine calculations.

Proposition 2.2

Let X be a random variable having DNXL distribution with parameter $\theta > 0$. The moment generating function of X is given by

$$M(t) = \frac{2e^t + e^\theta[(\theta - 2)e^t + 2e^\theta - \theta - 2]}{2(e^\theta - e^t)^2}, \quad -\infty < t < \infty.$$

From Proposition 2.2, the non-central moments of the DNXL distribution can be calculated taking into account the well-known formula $E[X^k] = M^{(k)}(0)$ for $k = 1, 2, \dots$, where $M^{(k)}(0)$ denotes the k -th derivative of $M(t) = E[e^{tX}]$ evaluated at $t = 0$. As a consequence, elementary calculations provide the following result.

Corollary 2.3

Let X be a random variable having DNXL distribution with parameter $\theta > 0$. The mean and variance of X are the following:

$$(i) \ E(X) = \frac{(\theta + 2)e^\theta - 2}{2(e^\theta - 1)^2},$$

$$(ii) \ Var(X) = \frac{e^\theta \left[(2(\theta + 2)e^\theta - \theta^2 - 8)e^\theta - 2(\theta - 2) \right]}{4(e^\theta - 1)^4}.$$

The dispersion index (DI) measures the variability or spread of data, providing insights into the distribution patterns. From Corollary 2.3, the DI of a random variable X having DNXL distribution with parameter $\theta > 0$ is the following

$$DI(X) = \frac{Var(X)}{E(X)} = \frac{e^\theta \left[(2(\theta + 2)e^\theta - \theta^2 - 8)e^\theta - 2(\theta - 2) \right]}{2[(\theta + 2)e^\theta - 2](e^\theta - 1)^2}. \quad (2.3)$$

Explicit expressions for the skewness coefficient defined by $\gamma_1 = E[(X - \mu)^3]/\sigma^3$ and the kurtosis coefficient defined by $\gamma_2 = E[(X - \mu)^4]/\sigma^4 - 3$, where $\mu = E(X)$ and $\sigma = \sqrt{Var(X)}$, can also be obtained for the DNXL distribution. However, their expressions are too involved and are omitted here for the sake of space. Table 1 displays the numerical values of the mean, variance, DI, skewness and kurtosis coefficients for different values of θ . From these results, it can be seen that the mean and variance decrease quickly as θ increases and also that the skewness and kurtosis coefficients increase quickly as θ increases. Moreover, clearly the DNXL distribution is positively skewed and can handle both overdispersed and underdispersed count data.

Table 1: Mean, variance, DI, skewness and kurtosis for different values of θ .

| θ | $E(X)$ | $Var(X)$ | $DI(X)$ | γ_1 | γ_2 |
|----------|---------|----------|---------|------------|------------|
| 0.25 | 5.51042 | 27.95809 | 5.07367 | 1.62518 | 3.80987 |
| 0.5 | 2.52098 | 6.95740 | 2.75986 | 1.64165 | 3.85312 |
| 1 | 1.04231 | 1.70491 | 1.63569 | 1.71371 | 4.04748 |
| 1.5 | 0.56450 | 0.72938 | 1.29209 | 1.85244 | 4.45297 |
| 2 | 0.33753 | 0.38592 | 1.14337 | 2.07751 | 5.20811 |
| 3 | 0.13510 | 0.13967 | 1.03383 | 2.86291 | 8.76364 |
| 5 | 0.02385 | 0.02384 | 0.99942 | 6.46945 | 41.82915 |
| 7 | 0.00411 | 0.00410 | 0.99913 | 15.57796 | 242.25703 |

2.4. Quantile function

The quantile function (qf) of an integer-valued random variable Y can be defined as follows

$$Q(u) = \inf\{y \in \mathbb{Z} : F(y) \geq u\}, \quad 0 < u < 1, \quad (2.4)$$

where F denotes the cdf of Y and \mathbb{Z} the set of integer numbers. This definition implies that the u -th quantile of Y is a unique integer value. There exists an alternative definition if the u -th quantile is defined as any y that satisfies the inequalities $F(y_-) \leq u \leq F(y)$,

where $F(y_-)$ stands for the left limit of F at y , which implies that the u -th quantile is an interval on the real line. From now on, definition (2.4) will be used in order to overcome drawbacks arising from the non-uniqueness.

A useful property of the DNXL distribution is that the quantile function can be expressed in closed form in terms of the Lambert W function, which is a multivalued complex function defined as the solution of the following equation

$$W(z)e^{W(z)} = z, \quad z \in \mathbb{C}. \quad (2.5)$$

The Lambert W function has only two real branches for real numbers $x \geq -1/e$. The principal branch is denoted by $W_0(x)$ and takes values in the interval $[-1, \infty)$ for $x \geq -1/e$. The negative branch is denoted by $W_{-1}(x)$ and takes values in the interval $(-\infty, -1]$ for $x \in [-1/e, 0)$. In the proof of the next result the following elementary properties concerning W_{-1} are used: $W_{-1}(-1/e) = -1$, $W_{-1}(x)$ is decreasing as x increases and $W_{-1}(x) \rightarrow -\infty$ as $x \rightarrow 0$. The reader is referred to Corless et al. (1996) for a comprehensive survey on the Lambert W function and to Jodrá (2010) for highlighting the importance of W_{-1} in relation with the Lindley distribution.

Proposition 2.3

Let X be a random variable having DNXL distribution with parameter $\theta > 0$. The qf of X is given by

$$Q(u; \theta) = \left\lceil -1 - \frac{2}{\theta} - \frac{1}{\theta} W_{-1} \left(2(u-1)e^{-2} \right) \right\rceil, \quad 0 < u < 1, \quad (2.6)$$

where $\lceil \cdot \rceil$ denotes the ceiling function.

Proof. For any $\theta > 0$ and $u \in (0, 1)$, the equation $F(x; \theta) = u$ can be written as follows

$$-e^{-(2+(x+1)\theta)}(2 + \theta(x+1)) = 2(u-1)e^{-2}, \quad x = 0, 1, \dots$$

By virtue of Eq. (2.5), the solution of the above equation is

$$W \left(2(u-1)e^{-2} \right) = -(2 + \theta(x+1)), \quad x = 0, 1, \dots \quad (2.7)$$

Additionally, it is easy to check the inequalities $-1/e < 2(u-1)e^{-2} < 0$ for any $u \in (0, 1)$ and $-(2 + \theta(x+1)) < -2$ for any $\theta > 0$ and $x = 0, 1, \dots$. As a consequence, taking into account the aforementioned properties of W_{-1} , the real branch of W in Eq. (2.7) is the negative branch. Finally, $Q(u; \theta)$ is obtained in closed form by solving x in Eq. (2.7) and taking into account definition (2.4). \square

From a computational point of view, pseudo-random data from the DNXL distribution can be easily computer-generated from formula (2.6) by applying the inverse transform method, that is, $Q(u; \theta)$ is evaluated in a pseudo-random value u generated from the standard uniform distribution and the result is a pseudo-random value from the DNXL distribution. In this respect, it is worth to mention that the real branches of the Lambert W function are available in usual computer algebra systems and programming languages; for example, in the R programming language the package `gsl` implements the functions `lambert_W0` for W_0 and `lambert_Wm1` for W_{-1} . Therefore, random samples from the DNXL distribution can be computer-generated in a straightforward manner.

In addition, the qfs of the extreme order statistics of the DNXL distribution can be written in closed form by virtue of Proposition 2.3. Let X_1, \dots, X_n be n independent and identically distributed (iid) random variables having DNXL distribution with parameter θ and denote the minimum and maximum order statistics by $X_{(1)} = \min\{X_1, \dots, X_n\}$ and $X_{(n)} = \max\{X_1, \dots, X_n\}$, respectively. The cdfs of $X_{(1)}$ and $X_{(n)}$ are the following, $F_{X_{(1)}}(x; \theta) = 1 - [1 - F(x; \theta)]^n$ and $F_{X_{(n)}}(x; \theta) = [F(x; \theta)]^n$, where $F(x; \theta)$ is given by Eq. (2.2). From these expressions, the qfs of $X_{(1)}$ and $X_{(n)}$ can also be written in closed form, in particular, $Q_{X_{(1)}}(u; \theta) = Q(1 - (1 - u)^{1/n}; \theta)$ and $Q_{X_{(n)}}(u; \theta) = Q(u^{1/n}; \theta)$, $0 < u < 1$, where Q is given by Eq. (2.6).

2.5. Reliability properties

Next, some reliability measures of the DNXL distribution are provided in closed form, specifically the failure rate function and the mean residual life. Furthermore, it is also seen that the DNXL distribution has increasing failure rate (IFR).

The failure rate function (hrf) is defined as the conditional probability of failure of a device at age x , given that it has not failed by time $x - 1$. This is an important function because uniquely characterizes a probability distribution. Let X be a random variable having DNXL distribution with parameter $\theta > 0$. The hrf of X is the following

$$\begin{aligned} h(x; \theta) &= P(X = x | X \geq x) = \frac{p(x; \theta)}{S(x - 1; \theta)} \\ &= 1 - \frac{[2 + (1 + x)\theta]e^{-\theta}}{2 + \theta x}, \quad x = 0, 1, 2, \dots \end{aligned}$$

Figure 2 represents the hrf of the DNXL distribution for some values of θ and suggests that the hrf is an increasing function.

The next result establishes that the DNXL distribution has discrete IFR. Some relationships among IFR and other discrete ageing classes such as NBU (new better than used) can be found in Lai and Xie (2006).

Proposition 2.4

The DNXL distribution has discrete IFR.

Proof. The result is a direct consequence of Proposition 2.1 since a log-concave distribution has discrete IFR (see Lemma 5.8 in Barlow and Proschan (1975)). \square

The residual life random variable at age x is defined as the remaining lifetime beyond that age and is denoted by $(X - x | X \geq x)$. The mean residual life at x of a random variable X having DNXL distribution with parameter $\theta > 0$ is given below

$$\begin{aligned} \mu(x; \theta) &= E(X - x | X \geq x) = \frac{1}{S(x - 1; \theta)} \sum_{j=x}^{\infty} S(j; \theta) \\ &= \frac{(\theta x + \theta + 2)e^{\theta} - \theta x - 2}{(\theta x + 2)(e^{\theta} - 1)^2}, \quad x = 0, 1, 2, \dots \end{aligned}$$

In addition, as a direct consequence of Proposition 2.4, the following result holds.

Corollary 2.4

The mean residual life of the DNXL distribution is a decreasing function.

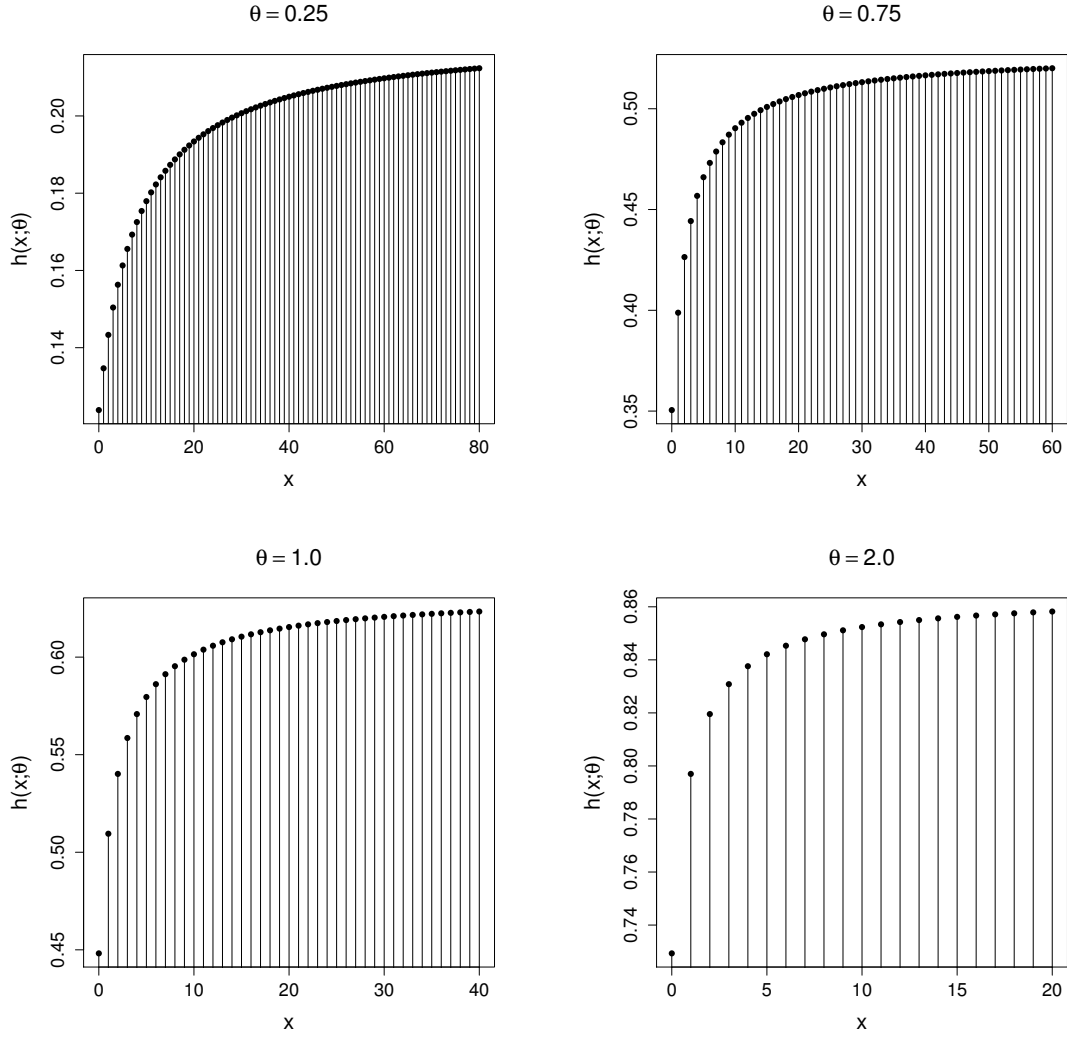


Figure 2: Hrf of the DNXL distribution for different values of θ .

3. Parameter estimation of the DNXL distribution

This section describes the estimation of θ applying the maximum likelihood (ML) method and the moments method (MM) together with their performance assessed via a Monte Carlo simulation study. The least squares method was also applied but this procedure provided estimates far apart from the true value of θ and, for this reason, its description together with the corresponding simulation study are omitted for the sake of space.

3.1. Maximum likelihood method

Let X_1, X_2, \dots, X_n be a random sample of size n from the DNXL distribution with unknown parameter θ and denote by x_1, x_2, \dots, x_n the observed values. From the likelihood function $L(\theta) = \prod_{i=1}^n p(x_i; \theta)$, the log-likelihood function is given by

$$\log L(\theta) = -n \log 2 - \theta \sum_{i=1}^n x_i + \sum_{i=1}^n \log \left[(2 + \theta x_i)(1 - e^{-\theta}) - \theta e^{-\theta} \right]. \quad (3.1)$$

The ML estimate of θ is obtained by maximizing $\log L(\theta)$ with respect to θ . To this end, the first derivative of $\log L(\theta)$ with respect to θ is the following

$$\frac{\partial}{\partial \theta} \log L(\theta) = -n\bar{x} + \sum_{i=1}^n \frac{[1 + \theta + (\theta - 1)x_i]e^{-\theta} + x_i}{(2 + \theta x_i)(1 - e^{-\theta}) - \theta e^{-\theta}}, \quad (3.2)$$

where $\bar{x} = (1/n) \sum_{i=1}^n x_i$. The ML estimate $\hat{\theta}$ of θ is the solution of Eq. (3.2) set equal to zero if $(\partial^2/\partial\theta^2) \log L(\theta)|_{\theta=\hat{\theta}} < 0$. Moreover, the Hessian matrix $H(\theta) = (\partial^2/\partial\theta^2) \log L(\theta)$ can be used to determine the uniqueness of the ML estimate (see Mäkeläinen et al. (1981) for full details).

As can be seen, Eq. (3.2) does not have an explicit solution and is too complicated to be solved by numerical methods. Hence, in order to calculate a ML estimate of θ for an observed sample it is necessary to solve the associated optimization problem, that is, $\max \log L(\theta)$ subject to the constrain $\theta > 0$. More computational details are given in Subsection 3.3.

3.2. Method of moments

From Corollary 2.3, the MM estimate of θ is obtained by solving the following equation

$$\frac{\theta e^\theta + 2e^\theta - 2}{2(e^\theta - 1)^2} = \bar{x}, \quad (3.3)$$

where \bar{x} denotes the sample mean of the observed values. From the above equation, it is clear that the MM estimator of θ cannot be given in closed form and for an observed sample the MM estimate of θ must be calculated numerically by solving Eq. (3.3). Some computational details are given in Subsection 3.3.

3.3. Computational considerations and simulation results

A Monte Carlo simulation study was carried out to evaluate the performance of the ML and MM methods. Let N be the number of random samples generated. For each simulated random sample j , let $\hat{\theta}_j$ be the resulting estimate of θ . In order to compare both estimation methods, the following quantities were computed: mean value $\bar{\theta} = \frac{1}{N} \sum_{j=1}^N \hat{\theta}_j$, Bias = $\bar{\theta} - \theta$, mean relative estimate MRE = $\frac{1}{\theta N} \sum_{j=1}^N |\hat{\theta}_j - \theta|$ and mean-square error MSE = $\frac{1}{N} \sum_{j=1}^N (\hat{\theta}_j - \theta)^2$. In particular, the simulation results reported in this subsection were obtained by generating $N = 10^4$ random samples of different sample sizes n for several values of θ . Pseudo-random data from the DNXL distribution were computer-generated by means of Eq. (2.6) and more specifically using the function `lambert_Wm1` available in the package `gs1` in the R programming language.

With respect to the ML method, the optimization problem was solved by the Brent algorithm for one-dimensional problems, which is available in the function `constrOptim` in the R programming language. The algorithm requires the objective function given by Eq. (3.1), the gradient function given by Eq. (3.2) and a feasible initial point in the parametric space of θ . In this regard, via simulation studies was checked that this point has no influence in the resulting estimates. Table 2 summarizes some simulation results and it can be seen that Bias, MRE and MSE decrease as n increases, as was expected. As a result, the performance of ML estimate proves to be consistently reliable. From the numerical results, it can be concluded that the ML method provides acceptable estimates of θ . In this regard, note that larger sample sizes are needed to obtain accurate estimates

as the value of θ increases, which is an expected behaviour since the number of zeros in the sample increases quickly as θ increases. Additionally, asymptotic confidence intervals for θ can be calculated based on the asymptotic normal approximation for the ML estimate $\hat{\theta}$; with this aim, it is useful to note that the asymptotic variance of $\hat{\theta}$ cannot be expressed in closed form and it should be evaluated numerically.

Table 2: Simulation results for the ML method for different values of θ and n .

| | $\theta = 0.25$ | | | | $\theta = 0.50$ | | | | $\theta = 0.75$ | | | |
|------------|-----------------|---------|---------|---------|-----------------|----------|---------|---------|-----------------|----------|---------|---------|
| | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE |
| $n = 25$ | 0.25876 | 0.00876 | 0.14666 | 0.00232 | 0.51789 | 0.01789 | 0.14607 | 0.00914 | 0.77786 | 0.02786 | 0.14768 | 0.02116 |
| $n = 50$ | 0.25432 | 0.00432 | 0.10063 | 0.00104 | 0.50932 | 0.00932 | 0.10237 | 0.00429 | 0.763435 | 0.01343 | 0.10182 | 0.00956 |
| $n = 75$ | 0.25240 | 0.00240 | 0.08233 | 0.00068 | 0.50550 | 0.00550 | 0.08352 | 0.00280 | 0.75786 | 0.00786 | 0.08268 | 0.00625 |
| $n = 100$ | 0.25195 | 0.00195 | 0.07030 | 0.00049 | 0.50422 | 0.00422 | 0.07104 | 0.00201 | 0.75662 | 0.00662 | 0.07083 | 0.00457 |
| $n = 200$ | 0.25085 | 0.00085 | 0.04956 | 0.00024 | 0.50182 | 0.00182 | 0.05039 | 0.00101 | 0.75376 | 0.00376 | 0.05064 | 0.00231 |
| $n = 500$ | 0.25065 | 0.00065 | 0.03109 | 0.00009 | 0.50068 | 0.00068 | 0.03162 | 0.00039 | 0.75097 | 0.00097 | 0.03124 | 0.00086 |
| $n = 1000$ | 0.25020 | 0.00020 | 0.02236 | 0.00004 | 0.50064 | 0.00064 | 0.02228 | 0.00019 | 0.75096 | 0.00096 | 0.02270 | 0.00045 |
| | $\theta = 1.0$ | | | | $\theta = 1.5$ | | | | $\theta = 2.0$ | | | |
| | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE |
| $n = 25$ | 1.03657 | 0.03657 | 0.15141 | 0.04016 | 1.56320 | 0.06320 | 0.15421 | 0.09528 | 2.09837 | 0.09837 | 0.16388 | 0.19367 |
| $n = 50$ | 1.01719 | 0.01719 | 0.10337 | 0.01771 | 1.52934 | 0.02934 | 0.10574 | 0.04170 | 2.04679 | 0.04679 | 0.11025 | 0.08230 |
| $n = 75$ | 1.01105 | 0.01105 | 0.08402 | 0.01141 | 1.52063 | 0.02063 | 0.08599 | 0.02720 | 2.02610 | 0.02610 | 0.08838 | 0.05096 |
| $n = 100$ | 1.00844 | 0.00844 | 0.07286 | 0.00849 | 1.51411 | 0.01411 | 0.07469 | 0.02031 | 2.02066 | 0.02066 | 0.07702 | 0.03840 |
| $n = 200$ | 1.00395 | 0.00395 | 0.05051 | 0.00407 | 1.50582 | 0.00582 | 0.05133 | 0.0094 | 2.01012 | 0.01012 | 0.05333 | 0.01806 |
| $n = 500$ | 1.00205 | 0.00205 | 0.03203 | 0.00162 | 1.50340 | 0.00340 | 0.03266 | 0.00377 | 2.00319 | 0.00319 | 0.03337 | 0.00704 |
| $n = 1000$ | 1.00085 | 0.00085 | 0.02239 | 0.00079 | 1.50208 | 0.00208 | 0.02293 | 0.00185 | 2.00257 | 0.00257 | 0.02422 | 0.00370 |
| | $\theta = 3.0$ | | | | $\theta = 4.0$ | | | | $\theta = 5.0$ | | | |
| | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE |
| $n = 25$ | 3.14018 | 0.14018 | 0.17453 | 0.44195 | 3.83689 | -0.16311 | 0.13325 | 0.35366 | 4.16743 | -0.83257 | 0.16651 | 0.86226 |
| $n = 50$ | 3.09655 | 0.09655 | 0.12782 | 0.26431 | 4.11237 | 0.11237 | 0.13273 | 0.43995 | 4.72353 | -0.27647 | 0.09801 | 0.37353 |
| $n = 75$ | 3.06057 | 0.06057 | 0.09984 | 0.15574 | 4.1315 | 0.13154 | 0.12058 | 0.39790 | 4.94081 | -0.05919 | 0.10663 | 0.38252 |
| $n = 100$ | 3.05246 | 0.05246 | 0.08701 | 0.11448 | 4.12453 | 0.12453 | 0.10883 | 0.33287 | 5.04764 | 0.04764 | 0.10614 | 0.40108 |
| $n = 200$ | 3.02643 | 0.02643 | 0.06077 | 0.05429 | 4.05342 | 0.05342 | 0.07252 | 0.14130 | 5.12067 | 0.12067 | 0.09188 | 0.36434 |
| $n = 500$ | 3.00936 | 0.00936 | 0.03795 | 0.02080 | 4.02300 | 0.02300 | 0.04427 | 0.05062 | 5.05303 | 0.05303 | 0.05575 | 0.13245 |
| $n = 1000$ | 3.00361 | 0.00361 | 0.02643 | 0.00995 | 4.01026 | 0.01026 | 0.03150 | 0.02523 | 5.02055 | 0.02055 | 0.03841 | 0.06032 |

With respect to the MM method, Eq. (3.3) can be solved numerically by applying the function `uniroot` available in the package `rootSolve` in the R programming language. It should be noted that the MM method does not provide a valid estimate of θ if all values of the sample are equal to zero, because in that case the solution of Eq. (3.3) is $\theta = 0$. Table 3 reports some simulation results and it must be remarked that the random samples with all values equal to zero were disregarded in order to obtain this table. In that case, it can be seen that the MM method also produces satisfactory estimates of θ .

Remark 3.1

A comparison of the MSE values in Tables 2 and 3 suggests that the ML and MM methods provide very similar results. However, an MM estimate of θ cannot be obtained if all values in the sample are equal to zero and the ML method should be used in that case.

Table 3: Simulation results for the MM method for different values of θ and n .

| | $\theta = 0.25$ | | | | $\theta = 0.50$ | | | | $\theta = 0.75$ | | | |
|------------|-----------------|---------|---------|---------|-----------------|----------|---------|---------|-----------------|----------|---------|---------|
| | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE |
| $n = 25$ | 0.25865 | 0.00865 | 0.14698 | 0.00232 | 0.51663 | 0.01663 | 0.14790 | 0.0093 | 0.77523 | 0.02523 | 0.14814 | 0.02128 |
| $n = 50$ | 0.25388 | 0.00388 | 0.10212 | 0.00106 | 0.50794 | 0.00794 | 0.10264 | 0.00432 | 0.76133 | 0.01133 | 0.10223 | 0.00969 |
| $n = 75$ | 0.25274 | 0.00274 | 0.08343 | 0.00070 | 0.50634 | 0.00634 | 0.08278 | 0.00281 | 0.75888 | 0.00888 | 0.08263 | 0.00619 |
| $n = 100$ | 0.25193 | 0.00193 | 0.07079 | 0.00050 | 0.50444 | 0.00444 | 0.07158 | 0.00204 | 0.75541 | 0.00541 | 0.07291 | 0.00479 |
| $n = 200$ | 0.25112 | 0.00112 | 0.05017 | 0.00024 | 0.50178 | 0.00178 | 0.05020 | 0.00099 | 0.75341 | 0.00341 | 0.05065 | 0.00231 |
| $n = 500$ | 0.25039 | 0.00039 | 0.03186 | 0.00009 | 0.50104 | 0.00104 | 0.03185 | 0.00040 | 0.75099 | 0.00099 | 0.03139 | 0.00086 |
| $n = 1000$ | 0.250118 | 0.00011 | 0.02221 | 0.00004 | 0.50038 | 0.00038 | 0.02221 | 0.00019 | 0.75043 | 0.00043 | 0.02271 | 0.00045 |
| | $\theta = 1.0$ | | | | $\theta = 1.5$ | | | | $\theta = 2.0$ | | | |
| | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE |
| $n = 25$ | 1.03668 | 0.03668 | 0.15101 | 0.03956 | 1.55653 | 0.05653 | 0.15652 | 0.09855 | 2.09269 | 0.09269 | 0.16381 | 0.19514 |
| $n = 50$ | 1.01685 | 0.01685 | 0.10324 | 0.01764 | 1.52613 | 0.02613 | 0.10478 | 0.04083 | 2.04342 | 0.04342 | 0.11174 | 0.08546 |
| $n = 75$ | 1.01054 | 0.01054 | 0.08354 | 0.01139 | 1.52014 | 0.02014 | 0.08627 | 0.02744 | 2.02627 | 0.02627 | 0.08935 | 0.05186 |
| $n = 100$ | 1.00831 | 0.00831 | 0.07303 | 0.00851 | 1.51400 | 0.01400 | 0.07427 | 0.02022 | 2.02036 | 0.02036 | 0.07659 | 0.03812 |
| $n = 200$ | 1.00404 | 0.00404 | 0.05092 | 0.00409 | 1.50746 | 0.00746 | 0.05179 | 0.00969 | 2.00797 | 0.00797 | 0.05385 | 0.01838 |
| $n = 500$ | 1.00192 | 0.00192 | 0.03233 | 0.00164 | 1.50124 | 0.00124 | 0.03251 | 0.00377 | 2.00568 | 0.00568 | 0.03397 | 0.00727 |
| $n = 1000$ | 1.00103 | 0.00103 | 0.02270 | 0.00081 | 1.50081 | 0.00081 | 0.02327 | 0.00191 | 2.00174 | 0.00174 | 0.02377 | 0.00357 |
| | $\theta = 3.0$ | | | | $\theta = 4.0$ | | | | $\theta = 5.0$ | | | |
| | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE | $\bar{\theta}$ | Bias | MRE | MSE |
| $n = 25$ | 3.12759 | 0.12759 | 0.17228 | 0.42835 | 3.82251 | -0.17749 | 0.13480 | 0.36542 | 4.15704 | -0.84296 | 0.16859 | 0.88463 |
| $n = 50$ | 3.10283 | 0.10283 | 0.12924 | 0.27039 | 4.11189 | 0.11189 | 0.13264 | 0.43903 | 4.73517 | -0.26483 | 0.09611 | 0.35913 |
| $n = 75$ | 3.06588 | 0.06588 | 0.10287 | 0.16536 | 4.13399 | 0.13399 | 0.12029 | 0.39783 | 4.95443 | -0.04557 | 0.10555 | 0.37636 |
| $n = 100$ | 3.04803 | 0.04803 | 0.08714 | 0.11521 | 4.11810 | 0.11810 | 0.10638 | 0.32102 | 5.05335 | 0.05335 | 0.10783 | 0.41198 |
| $n = 200$ | 3.02497 | 0.02497 | 0.06036 | 0.05453 | 4.05369 | 0.05369 | 0.07206 | 0.14167 | 5.12436 | 0.12436 | 0.09139 | 0.36719 |
| $n = 500$ | 3.01092 | 0.01092 | 0.03712 | 0.01985 | 4.02682 | 0.02682 | 0.04515 | 0.05313 | 5.05261 | 0.05261 | 0.05585 | 0.13415 |
| $n = 1000$ | 3.00282 | 0.00282 | 0.02647 | 0.00994 | 4.00933 | 0.00933 | 0.03144 | 0.02515 | 5.02843 | 0.02843 | 0.03850 | 0.06041 |

4. The first-order integer-valued autoregressive process with DNXL innovations

This section introduces a new INAR(1) process based on the binomial thinning operator together with innovations having a DNXL distribution, which is useful to handle both overdispersed and underdispersed count data in time series. Some statistical properties and the estimation of parameters of this process are also described. It should be noted that notations and some results are similar to those in [Irshad et al. \(2021\)](#).

4.1. Stochastic representation and properties

Before going further, let X be a non-negative integer-valued random variable. Recall that the binomial thinning operator ‘ \circ ’ is defined as a sum of X iid Bernoulli random variables W_j with probability success $\alpha \in (0, 1)$. Based on this operator, an INAR(1) process is defined by the following recursive equation

$$X_t = \alpha \circ X_{t-1} + \epsilon_t, \quad t \in \mathbb{Z}^+,$$

where $\alpha \circ X_{t-1} = \sum_{j=1}^{X_{t-1}} W_j$, $\{\epsilon_t\}_{t \in \mathbb{Z}^+}$ denotes a sequence of iid non-negative random variables with mean $E(\epsilon_t) = \mu_\epsilon$ and finite variance $\text{Var}(\epsilon_t) = \sigma_\epsilon^2$ and, additionally, ϵ_t is independent of the random variables W_j . The sequence $\{\epsilon_t\}_{t \in \mathbb{Z}^+}$ is the so-called innovation

process which can follow different probability distributions. With the preceding notations, the one step transition probability matrix of an INAR(1) process is the following

$$P(X_t = k | X_{t-1} = l) = \sum_{i=1}^{\min\{k,l\}} \binom{l}{i} \alpha^i (1-\alpha)^{l-i} P(\epsilon_t = k-i), \quad k, l \geq 0.$$

The new INAR(1) process is based on the binomial thinning operator and assumes that the innovation process $\{\epsilon_t\}_{t \in \mathbb{Z}^+}$ has a DNXL distribution with parameter $\theta > 0$. Accordingly, the one step transition probability matrix is given by

$$P(X_t = k | X_{t-1} = l) = \sum_{i=1}^{\min\{k,l\}} \binom{l}{i} \alpha^i (1-\alpha)^{l-i} \frac{1}{2} e^{-(k-i)\theta} \left[(\theta(k-i) + 2)(1 - e^{-\theta}) - \theta e^{-\theta} \right], \quad (4.1)$$

where $k, l \geq 0$. This process will be denoted by INAR(1)DNXL. Also the stationary marginal density of $\{X_t\}$ is given by

$$\begin{aligned} P(X_t = k) &= \sum_{l=0}^{\infty} P(X_t = k | X_{t-1} = l) P(X_{t-1} = l) \\ &= \sum_{l=0}^{\infty} \sum_{i=0}^{\min(k,l)} \binom{l}{i} \alpha^i (1-\alpha)^{l-i} \frac{1}{2} e^{-(k-i)\theta} \left[(\theta(k-i) + 2)(1 - e^{-\theta}) - \theta e^{-\theta} \right]. \end{aligned} \quad (4.2)$$

Taking into account Corollary 2.3, the next result provides the mean and variance of the INAR(1)DNXL process together with the conditional expectation and variance as well as the covariance and the correlation coefficient. The proof is omitted since it follows similar lines as Theorem 1 in Qi et al. (2019) (see also Weiß for further details).

Proposition 4.1

The INAR(1)DNXL process $\{X_t\}_{t \in \mathbb{Z}^+}$ has the following properties:

- (i) $E(X_t) = \frac{\mu_\epsilon}{1-\alpha} = \frac{(\theta+2)e^\theta - 2}{(1-\alpha)(e^\theta - 1)^2}$.
- (ii) $Var(X_t) = \frac{\sigma_\epsilon^2 + \alpha\mu_\epsilon}{1-\alpha^2} = \frac{1}{4(e^\theta - 1)^4(\alpha^2 - 1)} \left[4\alpha + 2((1-\alpha)\theta - 2(3\alpha + 1))e^\theta + (\theta(4\alpha + \theta) + 4(3\alpha + 2))e^{2\theta} - (2(\alpha + 1)\theta + 4(\alpha + 1))e^{3\theta} \right]$.
- (iii) $E(X_t | X_{t-1}) = \alpha X_{t-1} + \mu_\epsilon = \alpha X_{t-1} + \frac{(\theta+2)e^\theta - 2}{2(e^\theta - 1)^2}$.
- (iv) $Var(X_t | X_{t-1}) = \alpha(1-\alpha)X_{t-1} + \sigma_\epsilon^2 = \alpha(1-\alpha)X_{t-1} + \frac{\left[(2(\theta+2)e^\theta - \theta^2 - 8)e^\theta - 2(\theta-2) \right] e^\theta}{4(e^\theta - 1)^4}$.
- (v) $DI(X_t) = \frac{DI(\epsilon) + \alpha}{1+\alpha}$ and $DI(\epsilon)$ is given by Eq. (2.3).

(vi) $Cov(X_t, X_{t+k}) = \alpha^k Var(X_t)$ for $k \geq 1$, where $Var(X_t)$ is given in (ii).

(vii) $\rho_k = Corr(X_t, X_{t+k}) = \alpha^k$ for $k \geq 1$.

4.2. Estimation of parameters

In order to estimate the unknown parameters of an INAR(1) process, the conditional maximum likelihood (CML) method (see [Sprott \(1983\)](#) for further details), the conditional least squares (CLS) method (see [Klimko and Nelson \(1978\)](#)) and the Yule–Walker (YW) method (see [Al-Osh and Alzaid \(1987\)](#)) are the most frequently utilized in practice. Next, these methods are described for an INAR(1)DNXL process.

4.2.1. Conditional maximum likelihood method

Let X_1, X_2, \dots, X_T be a random sample from an INAR(1)DNXL process and denote by x_1, x_2, \dots, x_T the observed values. The log-likelihood function of the INAR(1)DNXL process is the following

$$\begin{aligned} \log L(\alpha, \theta) &= \sum_{t=2}^T \log(P(X_t = k | X_{t-1} = l)) \\ &= \sum_{t=2}^T \log \left(\sum_{i=1}^{\min\{X_t, X_{t-1}\}} \binom{X_{t-1}}{i} \alpha^i (1 - \alpha)^{X_{t-1}-i} \frac{1}{2} e^{-(X_t-i)\theta} \right. \\ &\quad \left. \times [(\theta(X_t - i) + 2)(1 - e^{-\theta}) - \theta e^{-\theta}] \right). \end{aligned} \quad (4.3)$$

The partial derivatives of $\log L(\alpha, \theta)$ with respect to α and θ are given by

$$\begin{aligned} \frac{\partial}{\partial \alpha} \log L(\alpha, \theta) &= \sum_{t=2}^T \frac{1}{A(\alpha, \theta)} \sum_{i=1}^{\min\{X_t, X_{t-1}\}} \binom{X_{t-1}}{i} \frac{1}{2} e^{-(X_t-i)\theta} \\ &\quad [(\theta(X_t - i) + 2)(1 - e^{-\theta}) - \theta e^{-\theta}] (1 - \alpha)^{X_{t-1}-i} \alpha^i \left(\frac{i}{\alpha} - \frac{X_{t-1} - i}{1 - \alpha} \right), \end{aligned} \quad (4.4)$$

$$\begin{aligned} \frac{\partial}{\partial \theta} \log L(\alpha, \theta) &= \sum_{t=2}^T \frac{1}{A(\alpha, \theta)} \sum_{i=1}^{\min\{X_t, X_{t-1}\}} \binom{X_{t-1}}{i} \alpha^i (1 - \alpha)^{X_{t-1}-i} \frac{1}{2} e^{-(X_t-i)\theta} \\ &\quad [-(X_t - i)(\theta i - \theta X_t - 1) + e^{-\theta}(X_t - i + 1)(\theta X_t - \theta i + \theta + 1)], \end{aligned} \quad (4.5)$$

where

$$A(\alpha, \theta) = \sum_{i=1}^{\min\{X_t, X_{t-1}\}} \binom{X_{t-1}}{i} \alpha^i (1 - \alpha)^{X_{t-1}-i} \frac{1}{2} e^{-(X_t-i)\theta} [(\theta(X_t - i) + 2)(1 - e^{-\theta}) - \theta e^{-\theta}].$$

The CML estimators of α and θ , say $\hat{\alpha}$ and $\hat{\theta}$, are obtained by maximizing Eq. (4.3). However, the above system of partial derivatives of $\log L(\alpha, \theta)$ with respect to each parameter set equal to zero does not have an explicit solution and then, for an observed sample, the CML estimates must be obtained numerically from the associated optimization problem, that is, $\max \log L(\alpha, \theta)$ subject to the constraints $0 < \alpha < 1$ and $\theta > 0$. This problem can be solved by the BFGS algorithm available in the function `optim` in the R programming language. In order to apply this algorithm, the objective function is given by Eq. (4.3) and the gradient function by Eqs. (4.4) and (4.5).

4.2.2. Conditional least squares method

Taking into account the notations in the preceding subsection, the CLS estimators of parameters α and θ are obtained by minimizing the following function

$$S(\alpha, \theta) = \sum_{t=2}^T (X_t - E(X_t|X_{t-1}))^2, \quad (4.6)$$

where $E(X_t|X_{t-1})$ is given in Proposition 4.1(iii). The derivatives of Eq. (4.6) with respect to α and θ set equal to zero, written in terms of the observed values, are the following:

$$\begin{aligned} \frac{\partial}{\partial \alpha} S(\alpha, \theta) &= \sum_{t=2}^T x_t x_{t-1} - \alpha \sum_{t=2}^T x_{t-1}^2 - \frac{(\theta + 2)e^\theta - 2}{2(e^\theta - 1)^2} \sum_{t=2}^T x_{t-1} = 0, \\ \frac{\partial}{\partial \theta} S(\alpha, \theta) &= \sum_{t=2}^T x_t - \alpha \sum_{t=2}^T x_{t-1} - (T-1) \frac{(\theta + 2)e^\theta - 2}{2(e^\theta - 1)^2} = 0. \end{aligned} \quad (4.7)$$

The CLS estimates of α and θ correspond to the solution of system (4.7), which cannot be expressed in closed form. Accordingly, the CLS estimates are obtained numerically from the associated optimization problem, that is, $\min S(\alpha, \theta)$ subject to the constraints $0 < \alpha < 1$ and $\theta > 0$. This problem can be solved by the aforementioned BFGS algorithm using the function `optim`. In order to apply this algorithm, the objective function is given by Eq. (4.6) and the gradient function by Eq. (4.7).

4.2.3. Yule-Walker method

With the notations in the preceding subsections, taking into account that the autocorrelation function (ACF) of an INAR(1) process at lag k is $\rho_k = \alpha^k$, the YW estimate of α , say $\hat{\alpha}_{YW}$, is the following

$$\hat{\alpha}_{YW} = \frac{\sum_{t=2}^T (x_t - \bar{x})(x_{t-1} - \bar{x})}{\sum_{t=1}^T (x_t - \bar{x})^2},$$

where $\bar{x} = (1/T) \sum_{t=1}^T x_t$. Then, by virtue of Proposition 4.1(i), the YW estimate of θ , say $\hat{\theta}_{YW}$, is obtained by solving the following equation

$$\frac{(\theta + 2)e^\theta - 2}{2(e^\theta - 1)^2(1 - \hat{\alpha}_{YW})} = \bar{x}. \quad (4.8)$$

An explicit solution of $\hat{\theta}_{YW}$ cannot be obtained from Eq. (4.8) and this equation must be solved numerically. To this end, the aforementioned function `uniroot` can be used.

4.3. Simulation results

A Monte Carlo simulation study was carried out to evaluate the performance of the CML, CLS and YW estimates of an INAR(1)DNXL process. With this aim, $N = 1000$ random samples of this process were generated for different sample sizes n and several values of α and θ . Table 4 summarizes some simulation results providing the mean estimates (Estimates), Bias and MSE for each estimation method. It can be seen that the Bias and MSE of the CML estimates tend to zero more quickly than those of the CLS and YW methods, both for small and large sample sizes. Therefore, the CML estimates

perform better than the YW and CLS estimates and, consequently, the CML method is recommended for the estimation of parameters of an INAR(1)DNXL process.

Table 4: Simulation results for the INAR(1)DNXL process.

| Parameters | n | CML | | | CLS | | | YW | | | |
|------------|------|-----------|--------|--------|-----------|--------|--------|-----------|--------|--------|--------|
| | | Estimates | Bias | MSE | Estimates | Bias | MSE | Estimates | Bias | MSE | |
| α | 0.25 | 50 | 0.2549 | 0.0714 | 0.0080 | 0.2241 | 0.1162 | 0.0210 | 0.2188 | 0.1146 | 0.0204 |
| | | 100 | 0.2506 | 0.0504 | 0.0039 | 0.2340 | 0.0804 | 0.0103 | 0.2317 | 0.0799 | 0.0102 |
| | | 150 | 0.2493 | 0.0408 | 0.0026 | 0.2371 | 0.0666 | 0.0069 | 0.2355 | 0.0664 | 0.0069 |
| | | 200 | 0.2514 | 0.0348 | 0.0019 | 0.2426 | 0.0577 | 0.0051 | 0.2415 | 0.0577 | 0.0051 |
| | | 500 | 0.2501 | 0.0309 | 0.0015 | 0.2405 | 0.0512 | 0.0042 | 0.2395 | 0.0511 | 0.0042 |
| θ | 0.6 | 50 | 0.6202 | 0.0769 | 0.0099 | 0.6086 | 0.0968 | 0.0156 | 0.6071 | 0.0948 | 0.0148 |
| | | 100 | 0.6066 | 0.0501 | 0.0041 | 0.6001 | 0.0656 | 0.0069 | 0.5997 | 0.0648 | 0.0067 |
| | | 150 | 0.6070 | 0.0403 | 0.0026 | 0.6022 | 0.0546 | 0.0048 | 0.6022 | 0.0543 | 0.0048 |
| | | 200 | 0.6045 | 0.0352 | 0.0020 | 0.6010 | 0.0464 | 0.0035 | 0.6009 | 0.0463 | 0.0035 |
| | | 500 | 0.6007 | 0.0305 | 0.0015 | 0.5966 | 0.0414 | 0.0027 | 0.5963 | 0.0410 | 0.0026 |
| α | 0.5 | 50 | 0.4909 | 0.0790 | 0.0100 | 0.4487 | 0.1158 | 0.0214 | 0.4392 | 0.1159 | 0.0216 |
| | | 100 | 0.4949 | 0.0546 | 0.0048 | 0.4719 | 0.0793 | 0.0104 | 0.4670 | 0.0803 | 0.0106 |
| | | 150 | 0.4954 | 0.0438 | 0.0031 | 0.4770 | 0.0650 | 0.0068 | 0.4738 | 0.0655 | 0.0069 |
| | | 200 | 0.4963 | 0.0396 | 0.0025 | 0.4823 | 0.0571 | 0.0052 | 0.4798 | 0.0573 | 0.0052 |
| | | 500 | 0.4996 | 0.0234 | 0.0009 | 0.4959 | 0.0343 | 0.0019 | 0.4948 | 0.0343 | 0.0019 |
| θ | 1.3 | 50 | 1.3381 | 0.1667 | 0.0474 | 1.2853 | 0.2016 | 0.0651 | 1.2859 | 0.1982 | 0.0629 |
| | | 100 | 1.3200 | 0.1180 | 0.0223 | 1.2916 | 0.1482 | 0.0351 | 1.2915 | 0.1481 | 0.0350 |
| | | 150 | 1.3164 | 0.0976 | 0.0152 | 1.2926 | 0.1228 | 0.0239 | 1.2923 | 0.1225 | 0.0238 |
| | | 200 | 1.3066 | 0.0826 | 0.0107 | 1.2882 | 0.1025 | 0.0165 | 1.2880 | 0.1025 | 0.0165 |
| | | 500 | 1.3035 | 0.0503 | 0.0040 | 1.2994 | 0.0652 | 0.0067 | 1.2994 | 0.0654 | 0.0067 |
| α | 0.75 | 50 | 0.7242 | 0.0728 | 0.0105 | 0.6750 | 0.1091 | 0.0214 | 0.6580 | 0.1157 | 0.0236 |
| | | 100 | 0.7400 | 0.0450 | 0.0037 | 0.7146 | 0.0689 | 0.0082 | 0.7071 | 0.0712 | 0.0087 |
| | | 150 | 0.7451 | 0.0363 | 0.0022 | 0.7262 | 0.0549 | 0.0051 | 0.7214 | 0.0559 | 0.0053 |
| | | 200 | 0.7433 | 0.0309 | 0.0016 | 0.7294 | 0.0460 | 0.0035 | 0.7257 | 0.0472 | 0.0037 |
| | | 500 | 0.7491 | 0.0192 | 0.0006 | 0.7415 | 0.0288 | 0.0013 | 0.7399 | 0.0291 | 0.0014 |
| θ | 2.5 | 50 | 2.5662 | 0.3358 | 0.1945 | 2.4161 | 0.4191 | 0.2732 | 2.4129 | 0.4128 | 0.2688 |
| | | 100 | 2.5244 | 0.2187 | 0.0787 | 2.4450 | 0.2875 | 0.1306 | 2.4449 | 0.2859 | 0.1289 |
| | | 150 | 2.5296 | 0.1890 | 0.0578 | 2.4688 | 0.2456 | 0.0938 | 2.4690 | 0.2454 | 0.0933 |
| | | 200 | 2.4993 | 0.1494 | 0.0350 | 2.4541 | 0.2061 | 0.0663 | 2.4540 | 0.2044 | 0.0656 |
| | | 500 | 2.5060 | 0.0925 | 0.0140 | 2.4837 | 0.1333 | 0.0280 | 2.4837 | 0.1333 | 0.0281 |

5. Data analysis

In this section, the DNXL distribution is used to model two real data sets and the resulting fits are compared to the ones provided by other discrete distributions. Similarly, a third data set is modelled by an INAR(1)DNXL process and the results are compared to those obtained with different competitive INAR(1) processes.

5.1. Methodology

To be more precise, the two data sets are described in Subsections 5.2 and 5.3 and both sets were fitted with the DNXL distribution together with the following distributions: discrete Rayleigh (DR) introduced by Roy (2004), discrete Pareto (DP) introduced by Krishna and Pundir (2009), discrete Burr (DB) introduced by Krishna and Pundir (2009), discrete inverse Weibull (DIW) introduced by Jazi et al. (2010), discrete Lomax (DL) introduced by Para and Jan (2016), discrete Burr type II (DB-XII) introduced by Para

and Jan (2016), discrete Teissier (DT) introduced by Irshad et al. (2023) and Poisson (P). For the sake of completeness, Table 5 shows the pmf of the above distributions.

Table 5: Discrete distributions fitted to the data.

| Model | Pmf | Support | Parameters |
|--------|---|---------------------------------|-------------------------------------|
| DR | $p(x; \alpha) = \alpha^{x^2} - \alpha^{(x+1)^2}$ | $x = 0, 1, 2, \dots$ | $0 < \alpha < 1$ |
| DP | $p(x; \alpha) = \alpha^{\log(1+x)} - \alpha^{\log(2+x)}$ | $x = 0, 1, 2, \dots$ | $0 < \alpha < 1$ |
| DB | $p(x; \alpha, \beta) = \beta^{\log(1+x^\alpha)} - \beta^{\log(1+(x+1)^\alpha)}$ | $x = 0, 1, 2, \dots$ | $\alpha > 0, 0 < \beta < 1$ |
| DIW | $p(x; \alpha, \beta) = \begin{cases} \beta & x = 1 \\ \beta^{x-\alpha} - \beta^{(x-1)-\alpha} & x = 2, 3, 4, \dots \end{cases}$ | $x = 1$ $x = 2, 3, 4, \dots$ | $\alpha > 0, 0 < \beta < 1$ |
| DL | $p(x; \alpha, \beta) = \beta^{\log(1+\frac{x}{\alpha})} - \beta^{\log(1+\frac{x+1}{\alpha})}$ | $x = 0, 1, 2, \dots$ | $\alpha > 0, 0 < \beta < 1$ |
| DB-XII | $p(x; \alpha, \gamma, \beta) = \beta^{\log(1+(\frac{x}{\alpha})^\gamma)} - \beta^{\log(1+(\frac{x+1}{\alpha})^\gamma)}$ | $x = 0, 1, 2, \dots$ | $\alpha, \gamma > 0, 0 < \beta < 1$ |
| DT | $p(x; \alpha) = \alpha^{-x} \left(e^{1-\alpha^{-x}} - \frac{1}{\alpha} e^{1-\alpha^{-(x+1)}} \right)$ | $x = 0, 1, 2, \dots$ | $0 < \alpha < 1$ |
| P | $p(x; \alpha) = \alpha^x e^{-\alpha} / x!$ | $x = 0, 1, 2, \dots$ | $\alpha > 0$ |

The fits provided by the DNXL distribution and the distributions in Table 5 were compared using the standard criteria of the lowest values of the Akaike information criterion (AIC) and Bayesian information criterion (BIC). Moreover, a Kolmogorov–Smirnov (KS) test for discrete distributions was applied to determine the goodness-of-fit of the above distributions. To further assess the goodness-of-fit of the DNXL distribution to the data sets, the following additional tests were also applied: Cramér von Mises (statistic W^2), Watson (statistic U^2), Anderson–Darling (statistic A^2) and Kuiper (statistic V). The p -values corresponding to these tests were calculated by a parametric bootstrap (cf. Babu and Rao (2004)) generating 10^4 bootstrap replicates. Specifically, the corresponding calculations were performed in R language. For the sake of completeness, the mathematical expressions of the above measures and statistics are given below:

$$\begin{aligned}
 AIC &= -2 \log L + 2r, \\
 BIC &= -2 \log L + r \log n, \\
 W^2 &= \frac{1}{12n} + \sum_{i=1}^n \left[\frac{2i-1}{2n} - F(x_{i:n}) \right]^2, \\
 U^2 &= W^2 n \left(\bar{F} - \frac{1}{2} \right)^2, \\
 A^2 &= \sum_{i=1}^n (2i-1) \left(\ln F(x_{i:n}) + \ln \left[1 - F(x_{n+1-i:n}) \right] \right), \\
 V &= D^+ + D^-,
 \end{aligned}$$

where n denotes the sample size, r the number of parameters, L denotes the maximized

value of the likelihood function, $\bar{F} = \frac{1}{n} \sum_{i=1}^n F(x_i)$, F is the cdf of the involved distribution, $\{x_i\}_{i=1}^n$ are the observed data, $x_{1:n} \leq x_{2:n} \leq \dots \leq x_{n:n}$ are the ordered observed data, $D^+ = \sup\{F_n(x) - F(x)\}$ and $D^- = \sup\{F(x) - F_n(x)\}$ where $F_n(x) = i/n$ for $x_{i:n} \leq x < x_{i+1:n}$ and $i = 1, \dots, n-1$, $F_n(x) = 0$ for $x < x_{1:n}$ and $F_n(x) = 1$ for $x_{n:n} \leq x$.

The third data set is considered in Subsection 5.4. The data were fitted with an INAR(1)DNXL process and with the following processes: INAR(1)P with Poisson innovations (Al-Osh and Alzaid (1987)), INAR(1)G with geometric innovations (Aghababaei Jazi et al. (2022)), INAR(1)DPL process with Poisson–Lindley innovations (Lívio et al. (2018)), INAR(1)NPWE process with new Poisson-weighted exponential innovations (Altun (2020)) and INAR(1)DT process with discrete Teissier innovations (Irshad et al. (2023)).

To further assess the model accuracy of the INAR(1)DNXL process for a data set, the standardized Pearson residuals were used, which are defined as follows

$$r_t = \frac{X_t - E(X_t|X_{t-1} = x_{t-1})}{\sqrt{\text{Var}(X_t|X_{t-1} = x_{t-1})}}, \quad t = 2, \dots, T,$$

where $E(X_t|X_{t-1})$ and $\text{Var}(X_t|X_{t-1})$ are given in Proposition 4.1. A cumulative periodogram (cpgram) of the standardized Pearson residuals will be plotted to determine whether the fitted INAR(1)DNXL process is random for a data set and the process will be considered statistically valid if these residuals are uncorrelated, have zero mean and unit variance (see Harvey and Fernandes (1989) for more details). In this regard, the ACF of the residuals will be represented graphically to determine the existence of correlation.

5.2. Failure times data set

The data represent the failure times in minutes of 15 electronic components in an acceleration lifetime test. These data can be found in Lawless (2011), page 204, and the sorted values (discretized) are the following: 1, 5, 6, 11, 12, 19, 20, 22, 23, 31, 37, 46, 54, 60, 66. Figure 3 displays the TTT (total time on test) plot for the data set and reveals an increasing hrf. Table 6 shows the results of the fitted models together with the KS test and Table 7 displays the bootstrap p -values of the additional goodness-of-fit tests for the DNXL distribution.

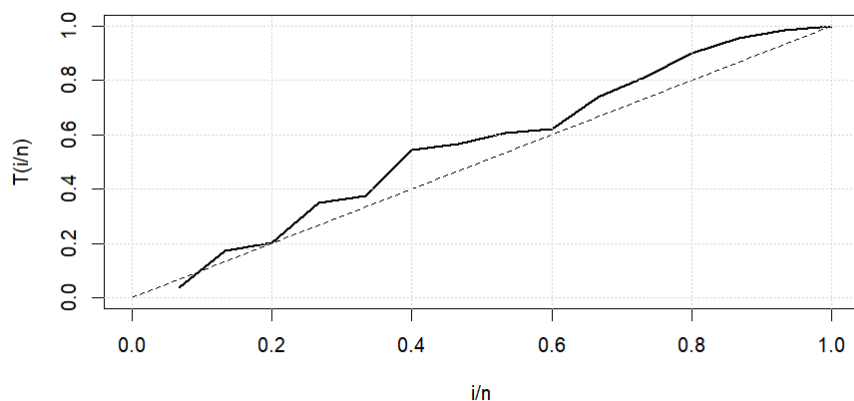


Figure 3: TTT plot for the failure times data.

Table 6: Failure times. Model, ML estimate, information criteria, KS statistic and p -value.

| Model | Estimates | $-\log L$ | AIC | BIC | KS stat. | p -value |
|--------|---------------------------|-----------|----------|----------|----------|------------|
| DNXL | $\hat{\theta} = 0.0543$ | 64.4653 | 130.9305 | 131.6386 | 0.1456 | 0.8639 |
| | $\hat{\alpha} = 9.1186$ | | | | | |
| DB-XII | $\hat{\gamma} = 0.0502$ | 65.6388 | 137.2777 | 139.4018 | 0.1505 | 0.8376 |
| | $\hat{\beta} = 0.0018$ | | | | | |
| DL | $\hat{\alpha} = 0.0128$ | 65.7182 | 135.4365 | 136.8526 | 0.2000 | 0.5219 |
| | $\hat{\beta} = 0.0051$ | | | | | |
| DR | $\hat{\alpha} = 0.9992$ | 66.3943 | 134.7886 | 135.4966 | 0.2401 | 0.3022 |
| DT | $\hat{\alpha} = 0.9710$ | 68.4277 | 138.8555 | 139.5635 | 0.3255 | 0.0647 |
| DB | $\hat{\alpha} = 169.4125$ | 70.3262 | 144.6524 | 146.0685 | 0.3318 | 0.0567 |
| | $\hat{\beta} = 0.9981$ | | | | | |
| DIW | $\hat{\alpha} = 0.7111$ | 70.4214 | 144.8427 | 146.2588 | 0.2195 | 0.4069 |
| | $\hat{\beta} = 0.0077$ | | | | | |
| DP | $\hat{\alpha} = 0.7201$ | 77.4023 | 156.8047 | 157.5127 | 0.3781 | 0.0195 |
| P | $\hat{\alpha} = 27.5333$ | 151.2064 | 304.4129 | 305.1209 | 0.3815 | 0.0180 |

Table 7: Failure times. Goodness-of-fit tests for the DNXL distribution.

| | W^2 | U^2 | A^2 | V |
|----------------------|--------|--------|--------|--------|
| Bootstrap p -value | 0.8184 | 0.9598 | 0.8533 | 0.8718 |

From the results in Tables 6 and 7, it may be concluded that the DNXL distribution provides a better fit than the other competing distributions, because it yields the lowest AIC and BIC values, the largest p -value, the lowest KS statistic value and the additional goodness-of-fit tests also show a suitable fit of the data. The Hessian matrix associated with the failure times data was calculated by using the R package `matrixcalc`, which provides $H(\hat{\theta}) = -6687.389$. This result implies that the Hessian matrix is negative definite confirming the uniqueness of the ML estimate.

5.3. Leukemia remission times data set

The data represent the number of weeks that 20 leukaemia patients spent in remission in a particular treatment. These data can be found in Lawless (2011), page 346, and the sorted values are the following: 1, 3, 3, 6, 7, 7, 10, 12, 14, 15, 18, 19, 22, 26, 28, 29, 34, 40, 48, 49. The TTT plot depicting the data set is showcased in Figure 4 and reveals an increasing hrf. Table 8 reports the results of the fitted models together with the KS test and Table 9 displays the bootstrap p -values of the additional goodness-of-fit tests for the DNXL distribution.

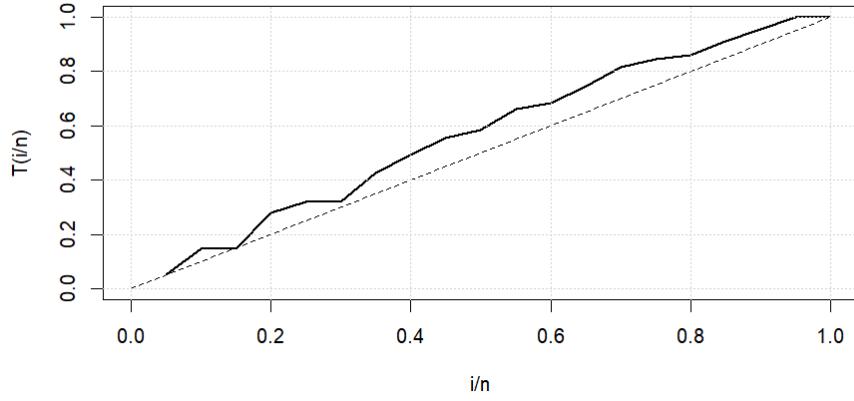


Figure 4: TTT plot for the remission times data.

Table 8: Leukemia remission times. Model, ML estimate, information criteria, KS statistic and p -value.

| Model | Estimates | $-\log L$ | AIC | BIC | KS stat. | p -value |
|--------|---|-----------|----------|----------|----------|------------|
| DNXL | $\hat{\theta}=0.0760$ | 79.2063 | 160.4125 | 161.4083 | 0.1064 | 0.9773 |
| DB-XII | $\hat{\alpha} = 1.9957$ $\hat{\gamma} = 1.6581 \cdot 10^{-3}$ $\hat{\beta} = 5.0288 \cdot 10^{-27}$ | 80.0435 | 166.0869 | 169.0741 | 0.1461 | 0.7866 |
| DL | $\hat{\alpha}=93.9177$ $\hat{\beta}=0.0046$ | 80.9117 | 165.8233 | 167.8148 | 0.1708 | 0.6041 |
| DR | $\hat{\alpha}=0.9984$ | 81.1750 | 164.3500 | 165.3458 | 0.2226 | 0.2748 |
| DT | $\hat{\alpha}=0.9599$ | 84.1298 | 170.2596 | 171.2554 | 0.2637 | 0.1239 |
| DIW | $\hat{\alpha} = 0.8166$ $\hat{\beta} = 0.0070$ | 84.9098 | 173.8197 | 175.8111 | 0.1867 | 0.4887 |
| DB | $\hat{\alpha}=182.3745$ $\hat{\beta}=0.9980$ | 87.6893 | 179.3787 | 181.3701 | 0.3238 | 0.0302 |
| DP | $\hat{\alpha}=0.6958$ | 95.4480 | 192.8959 | 193.8917 | 0.3563 | 0.0125 |
| P | $\hat{\alpha}=19.5500$ | 152.7180 | 307.4360 | 308.4317 | 0.3523 | 0.0140 |

Table 9: Leukemia remission times. Goodness-of-fit tests for the DNXL distribution.

| | W^2 | U^2 | A^2 | V |
|----------------------|--------|--------|--------|--------|
| Bootstrap p -value | 0.7118 | 0.9735 | 0.7058 | 0.9573 |

Based on the results in Tables 8 and 9, it can be concluded that the DNXL distribution provides a better fit than the other competing distributions for the data set under study. Additionally, the uniqueness of the ML estimate can be confirmed through the Hessian matrix since $H(\hat{\theta}) = -4549.184$.

5.4. Robbery data set

The data set was collected by the 54th police vehicle in Pittsburgh and the data represent the monthly counts of robberies spanning the period of January 1990 to December 2001. The corresponding time series data can be found in the website Forecasting Principles (www.forecastingprinciples.com) and the data are reported in Table 10 below. The mean, variance and DI of the resulting 144 observations are 2.1528, 3.2772 and 1.5223, respectively, and clearly the data exhibit overdispersion.

Table 10: Monthly counts of robberies reported by the 54th police vehicle.

| | Jan | Feb | Mar | Apr | May | Jun | Jul | Aug | Sep | Oct | Nov | Dec |
|------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 1990 | 4 | 0 | 2 | 3 | 6 | 2 | 0 | 2 | 3 | 4 | 3 | 5 |
| 1991 | 3 | 3 | 4 | 4 | 1 | 5 | 1 | 5 | 1 | 2 | 1 | 1 |
| 1992 | 0 | 1 | 1 | 1 | 1 | 1 | 5 | 4 | 2 | 2 | 3 | 6 |
| 1993 | 5 | 4 | 2 | 0 | 0 | 4 | 3 | 2 | 3 | 1 | 1 | 4 |
| 1994 | 2 | 0 | 2 | 0 | 0 | 0 | 3 | 4 | 2 | 3 | 1 | 4 |
| 1995 | 5 | 2 | 0 | 3 | 3 | 2 | 1 | 3 | 2 | 0 | 3 | 5 |
| 1996 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 4 | 0 | 0 | 2 | 1 |
| 1997 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 2 | 2 | 4 |
| 1998 | 2 | 1 | 3 | 1 | 1 | 1 | 2 | 5 | 3 | 0 | 2 | 5 |
| 1999 | 7 | 3 | 2 | 0 | 4 | 0 | 1 | 5 | 2 | 7 | 5 | 0 |
| 2000 | 9 | 5 | 3 | 2 | 1 | 1 | 2 | 2 | 3 | 2 | 2 | 4 |
| 2001 | 0 | 0 | 1 | 1 | 0 | 4 | 2 | 0 | 0 | 0 | 2 | 1 |

Figure 5 displays the plots corresponding to the ACF, partial autocorrelation function (PACF), time series and barplot. It can be observed that the INAR(1)DNXL process is an appropriate candidate model for fitting the data set because only the first lag is noteworthy in the PACF plot.

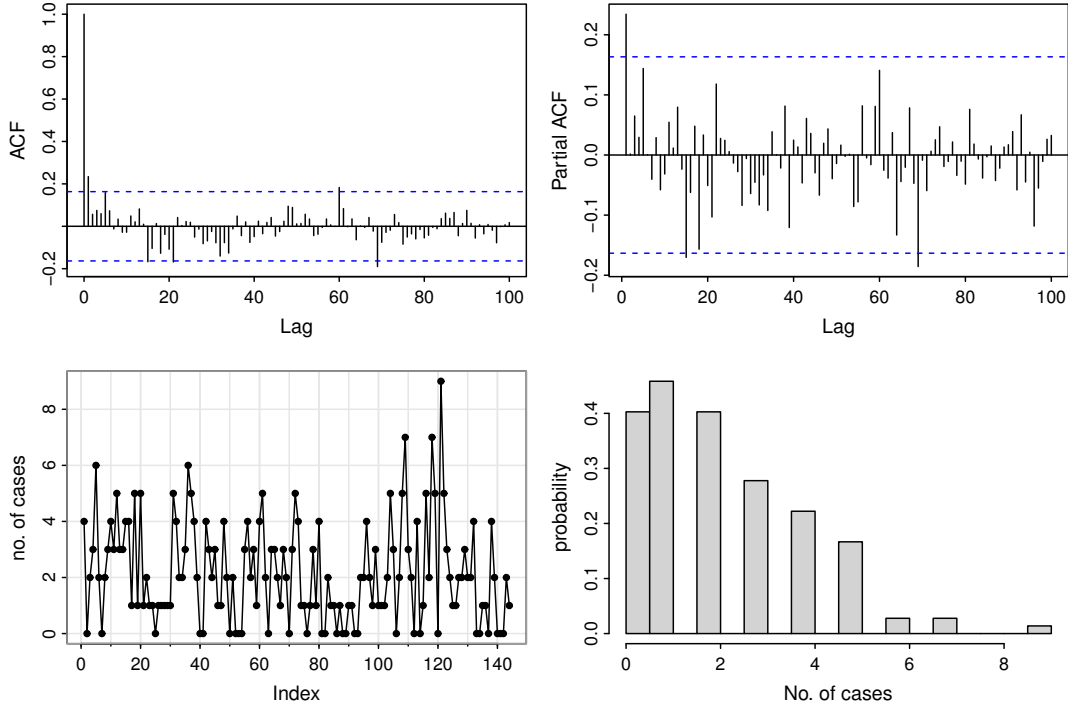


Figure 5: Monthly robberies. ACF, PACF, time series and barplot.

Table 11 reports the results of the different INAR(1) processes fitted to the data set together with some model adequacy measures. In the table, s.e. denotes the standard error of the estimates; μ , σ^2 and DI denote the mean, variance and DI of each INAR(1) process. The parameters of the processes were estimated by the CML method. As can be seen, the INAR(1)DNXL process yielded the lowest AIC and BIC values and, therefore, the new model offers a better fit than the other INAR(1) processes under consideration.

Table 11: Monthly robberies. CML estimates and adequacy measures of the fitted processes.

| Process | Param. | CML est. | s.e. | $-\log L$ | AIC | BIC | μ | σ^2 | DI |
|-------------|----------|----------|---------|-----------|----------|----------|--------|------------|--------|
| INAR(1)DNXL | α | 0.2793 | 0.0571 | -265.2417 | 534.4834 | 540.4231 | 2.1216 | 3.7815 | 1.7824 |
| | θ | 0.7510 | 0.0681 | | | | | | |
| INAR(1)DPL | α | 0.2947 | 0.0574 | -265.7662 | 535.5325 | 541.4721 | 2.1375 | 4.0717 | 1.9049 |
| | β | 0.9957 | 0.1069 | | | | | | |
| INAR(1)NPWE | α | 0.3206 | 0.0534 | -266.9818 | 539.9636 | 548.8730 | 2.1300 | 4.4637 | 2.0957 |
| | β | 0.6775 | 26.8266 | | | | | | |
| | γ | 0.0200 | 40.3901 | | | | | | |
| INAR(1)G | α | 0.3424 | 0.0733 | -267.9794 | 539.9589 | 545.8985 | 2.1289 | 4.9146 | 2.3085 |
| | β | 0.4167 | 0.0354 | | | | | | |
| INAR(1)P | α | 0.1847 | 0.1646 | -272.6633 | 549.3267 | 555.2663 | 2.1351 | 2.1351 | 1.0000 |
| | β | 1.7408 | 0.0607 | | | | | | |
| INAR(1)DT | α | 0.0277 | 0.0430 | -300.9203 | 605.8407 | 611.7803 | 2.6875 | 2.1744 | 0.8091 |
| | β | 0.7250 | 0.0111 | | | | | | |

The fitted INAR(1)DNXL process for the robbery data set is given by

$$X_t = 0.2793 X_{t-1} + \epsilon_t, \quad t = 2, 3, \dots, T,$$

where ϵ_t has a DNXL distribution with parameter $\theta = 0.7510$. Furthermore, by virtue of Proposition 4.1, the predicted values obtained from this process are the following:

$$\begin{aligned}\hat{X}_1 &= E(X_1) = 2.1216, \\ \hat{X}_t &= E(X_t|X_{t-1}) = 0.2793 X_{t-1} + 1.5289, \quad t = 2, 3, \dots, T.\end{aligned}$$

As said before, the standardized Pearson residuals were used to determine the model accuracy of the fitted INAR(1)DNXL process. Figure 6 displays the ACF of the residuals and shows graphically that there is no correlation. Additionally, the Ljung–Box test was performed with 10 degrees of freedom and the resulting p -value was 0.8146, which clearly indicates that the residuals are uncorrelated. Moreover, their mean and variance are 0.006 and 0.9036, respectively, values close to zero and one as desired. Figure 7 also displays the corresponding cpggram which shows that the INAR(1)DNXL process is random for the data set. The functions `Box.test` and `cpgram` available in the R programming language can be used, respectively, to perform the Ljung–Box test and to represent the cpggram. From the overall results, it can be concluded that the INAR(1)DNXL process provides a suitable fit for the data set under consideration.

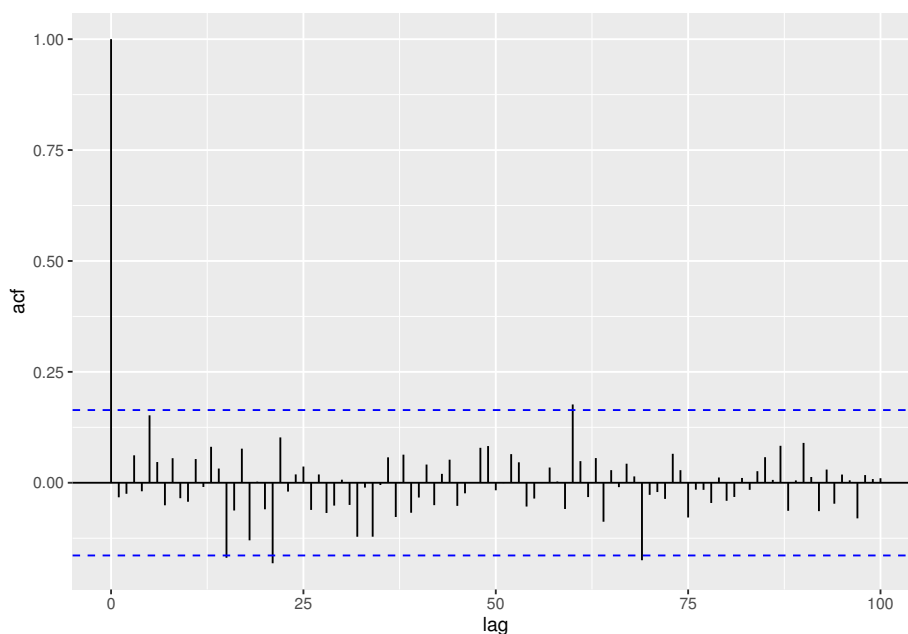


Figure 6: Monthly robberies. ACF plot of the Pearson residuals.

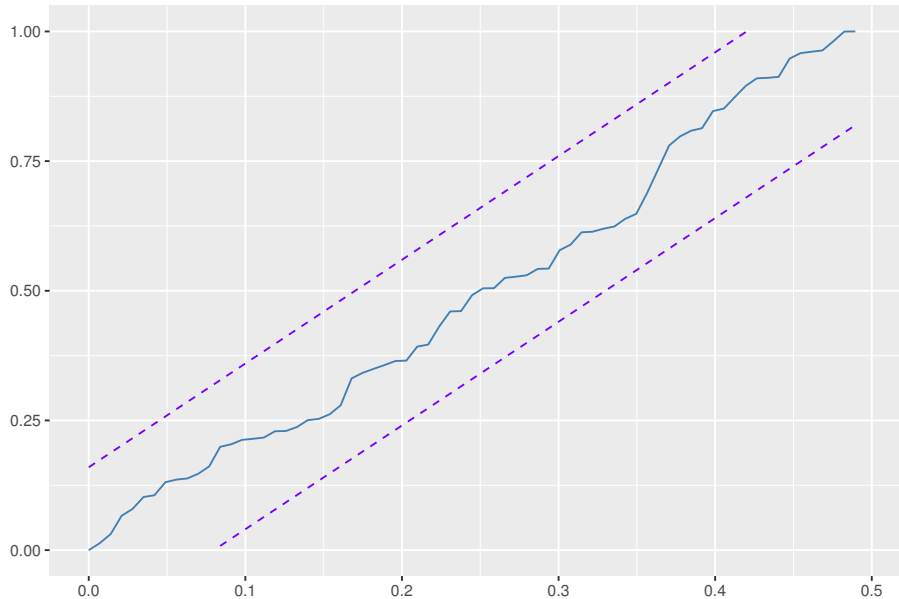


Figure 7: Monthly robberies. Cpgram of the standardized Pearson residuals.

6. CUSUM chart to monitor a mean increase of the INAR(1)DNXL process

In this section, a control chart for detection of upward shifts in the mean is studied. Specifically, it is used the CUSUM control chart suggested by [Page \(1961\)](#) due to their high sensitivity for detecting small shifts in the process mean. It relies on the theory of the sequential probability ratio test and the objective is to accumulate information from sample data and thus magnifying the impact of minor process deviations. Moreover, CUSUM charts exhibit established optimality properties when identifying a persistent shift from a known in-control distribution to a predefined out-of-control distribution (see [Weiss and Testik \(2009\)](#)). The upper-sided control chart is constructed by exclusively focusing on detecting upward shifts in the mean μ , as they are directly linked to process deterioration. The primary purpose is to detect quickly and accurately a change in μ . The CUSUM statistic is the following

$$C_t = \max\{0, X_t - k + C_{t-1}\}, \quad t \in \mathbb{N},$$

where $C_0 = c_0$. More specifically, c_0 denotes an initial non-negative value typically chosen to be zero, $k \geq \mu$ is the reference value and, for $\mu < k < h$, the monitoring statistics $\{C_t\}_{t \in \mathbb{N}}$ are plotted on a CUSUM chart with control region $[0, h]$, where h is the upper control limit. The INAR(1)DNXL process is considered as being in control unless $C_t > h$. The inclusion of the reference value k prevents the CUSUM statistic from drifting towards h , enhancing its sensitivity to respond to a mean increase.

6.1. Performance of the CUSUM control chart

A widely used metric for evaluating the efficiency of a control chart within the statistical process control framework is the average run length (ARL). The ARL represents the anticipated number of data points plotted on a control chart before it signals an out-of-control alarm. The two types of ARL performance that are of interest involve ARL_0 (in-control ARL), which assesses the number of points on a chart from the monitoring

commencement until a false alarm is triggered, together with ARL_1 (out-of-control ARL), which assesses the number of points on a chart from the initiation of a process shift until the chart detects the shift.

It should be noted that the bivariate process $(X_t, C_t)_{t \in \mathbb{N}}$ has the Markov property, characterized by the following transition probabilities:

$$\begin{aligned} P(x, y|l, z) &= P(X_t = x, C_t = y | X_{t-1} = l, C_{t-1} = z) \\ &= P(X_t = x | X_{t-1} = l) \mathbf{I}_y(\max\{0, x - k + z\}), \end{aligned} \quad (6.1)$$

$$\begin{aligned} P_1(x, y|l) &= P(X_1 = x, C_1 = y | C_0 = l) \\ &= P_x \mathbf{I}_y \max\{(0, x - k + l)\}, \end{aligned} \quad (6.2)$$

where $\mathbf{I}_y(\cdot)$ denotes the indicator function, and $P(X_t = x | X_{t-1} = l)$ and P_x are given by (4.1) and (4.2), respectively.

The effectiveness of the CUSUM chart in detecting changes in the process mean is evaluated through a simulation study. To determine the ARLs for the CUSUM chart design, it is used a Brook's Markov chain approach to perform the calculations (see Brook and Evans (1972)). The empirical acquisition of marginal probabilities is achieved through the simulation of a process of a specified size under given values of (α, θ) . In addition to ARL as a performance metric for the control charts, the typical relative deviation (expressed as a percentage) is also integrated in the ARL, $dev_{ARL} = 100\% \times (ARL - ARL_0)/ARL_0$, which evaluates and assesses the effectiveness and performance of the chart.

In particular, we set θ_0 to be 1.5 and 1.0, with θ_0 decreasing by 0.1, resulting in an increase in the mean under α_0 . Under the CUSUM chart, for a given $C_0 = c_0$, the focus revolves around possible integers h and k such that the desired in control ARL is close to 200. Tables 12 and 13 report an analysis of the performance of ARL within the CUSUM chart with $\theta_0 = (1.5, 1.0)$. The findings indicate the notable effectiveness of the CUSUM chart in identifying upward shifts in the mean μ . Moreover, there is a substantial reduction in the ARL as μ increases. It is essential to highlight that the decline in ARL is considerably more pronounced than the corresponding increase in the magnitude of the shift in μ . Therefore, it can be inferred that the CUSUM control chart exhibits effective performance when applied to the INAR(1)DNXL process. Its proficiency in efficiently identifying upward shifts in the mean renders it in a suitable and advantageous tool for monitoring and controlling the INAR(1)DNXL process.

Table 12: ARL with deviation (in parenthesis) for the CUSUM chart with $\alpha_0 = 0.6$ and $\theta_0 = 1.5$.

| α_0 | k | h | c_0 | α | θ | | | | | |
|------------|------|-----|-------|----------|---------------------|---------------------|---------------------|--------------------|--------------------|--------------------|
| | | | | | 1.5 | 1.4 | 1.3 | 1.2 | 1.1 | |
| 0.6 | 3 | 5 | 0 | 0.6 | 204.63 | 129.24 (-36.84%) | 82.71 (-59.58%) | 53.9 (-73.66%) | 35.86 (-82.48%) | |
| | | | | 0.65 | 120.07 (-41.32%) | 79.76 (-61.02%) | 53.87 (-73.67%) | 37.06 (-81.89%) | 26 (-87.29%) | |
| | | | | 0.7 | 71.58 (-65.02%) | 50.25 (-75.44%) | 35.87 (-82.47%) | 26.05 (-87.27%) | 19.26 (-90.59%) | |
| | 3 | 5 | 2 | 0.6 | 200.76 | 126.59 (-36.94%) | 80.63 (-59.84%) | 52.25 (-73.97%) | 34.53 (-82.80%) | |
| | | | | 0.65 | 117.48 (-41.48%) | 77.73 (-61.28%) | 52.25 (-73.97%) | 35.75 (-82.19%) | 24.92 (-87.59%) | |
| | | | | 0.7 | 69.69 (-65.29%) | 48.73 (-75.73%) | 34.63 (-82.75%) | 25.02 (-87.54%) | 18.39 (-90.84%) | |
| | 3 | 6 | 5 | 0.6 | 255.39 | 154.72 (-39.42%) | 94.07 (-63.17%) | 58.15 (-77.23%) | 36.65 (-85.65%) | |
| | | | | 0.65 | 139.82 (-45.25%) | 88.94 (-65.18%) | 57.47 (-77.50%) | 37.79 (-85.20%) | 25.35 (-90.07%) | |
| | | | | 0.7 | 77.73 (-69.56%) | 52.62 (-79.40%) | 36.2 (-85.83%) | 25.35 (-90.07%) | 18.06 (-92.93%) | |
| | 0.75 | 4 | 9 | 0 | 0.75 | 215.36 | 131.6 (-38.89%) | 82.64 (-61.63%) | 53.48 (-75.17%) | 35.73 (-83.41%) |
| | | | | | 0.8 | 89.03 (-58.66%) | 60.55 (-71.88%) | 42.27 (-80.37%) | 30.29 (-85.94%) | 22.27 (-89.66%) |
| | | | | | 0.85 | 40.75 (-81.08%) | 30.86 (-85.67%) | 23.83 (-88.93%) | 18.74 (-91.30%) | 14.98 (-93.04%) |
| 4 | | 9 | 2 | 0.75 | 213.53 | 130.24 (-39.01%) | 81.6 (-61.78%) | 52.67 (-75.33%) | 35.08 (-83.57%) | |
| | | | | 0.8 | 87.95 (-58.81%) | 59.71 (-72.04%) | 41.59 (-80.52%) | 29.74 (-86.07%) | 21.8 (-89.79%) | |
| | | | | 0.85 | 40.14 (-81.20%) | 30.35 (-85.79%) | 23.39 (-89.05%) | 18.35 (-91.41%) | 14.63 (-93.15%) | |
| 4 | | 9 | 5 | 0.75 | 207.84 | 126.09 (-39.33%) | 78.49 (-62.23%) | 50.28 (-75.81%) | 33.19 (-84.03%) | |
| | | | | 0.8 | 84.86 (-59.17%) | 57.31 (-72.43%) | 39.68 (-80.91%) | 28.17 (-86.45%) | 20.49 (-90.14%) | |
| | | | | 0.85 | 38.46 (-81.50%) | 28.94 (-86.08%) | 22.19 (-89.32%) | 17.3 (-91.68%) | 13.69 (-93.41%) | |
| 0.85 | | 6 | 12 | 0 | 0.85 | 207.7 | 123.96 (-40.32%) | 77.1 (-62.88%) | 50.04 (-75.91%) | 33.86 (-83.70%) |
| | | | | | 0.9 | 51.56 (-75.18%) | 37.87 (-81.77%) | 28.6 (-86.23%) | 22.13 (-89.35%) | 17.5 (-91.57%) |
| | | | | | 0.95 | 19.62 (-90.55%) | 16.89 (-91.87%) | 14.63 (-92.96%) | 12.74 (-93.87%) | 11.14 (-94.64%) |
| | 6 | 12 | 2 | 0.85 | 206.31 | 123.01 (-40.37%) | 76.42 (-62.96%) | 49.53 (-75.99%) | 33.47 (-83.78%) | |
| | | | | 0.9 | 51.09 (-75.24%) | 37.49 (-81.83%) | 28.28 (-86.29%) | 21.87 (-89.40%) | 17.26 (-91.63%) | |
| | | | | 0.95 | 19.42 (-90.59%) | 16.69 (-91.91%) | 14.44 (-93.00%) | 12.56 (-93.91%) | 10.97 (-94.68%) | |
| | 6 | 12 | 5 | 0.85 | 202.87 | 120.65 (-40.53%) | 74.74 (-63.16%) | 48.28 (-76.20%) | 32.49 (-83.99%) | |
| | | | | 0.9 | 49.96 (-75.37%) | 36.58 (-81.97%) | 27.52 (-86.43%) | 21.21 (-89.55%) | 16.68 (-91.78%) | |
| | | | | 0.95 | 18.91 (-90.68%) | 16.21 (-92.01%) | 13.98 (-93.11%) | 12.12 (-94.03%) | 10.53 (-94.81%) | |

Table 13: ARL with deviation (in parenthesis) for the CUSUM chart with $\alpha_0 = 0.6$ and $\theta_0 = 1.0$.

| α_0 | k | h | c_0 | α | θ | | | | | |
|------------|-----|-----|-------|----------|-----------|-----------|-----------|-----------|-----------|-----------|
| | | | | | 1 | 0.9 | 0.8 | 0.7 | 0.6 | |
| 0.6 | 4 | 13 | 0 | 0.6 | 201.02 | 95.08 | 48.14 | 26.41 | 15.69 | |
| | | | | | | (-52.70%) | (-76.05%) | (-86.86%) | (-92.19%) | |
| | | | | | 100.1 | 53.86 | 30.93 | 18.98 | 12.37 | |
| | | | | | | (-50.20%) | (-73.21%) | (-84.61%) | (-90.56%) | (-93.85%) |
| | 4 | 13 | 0 | 0.7 | 53.36 | 32.57 | 20.98 | 14.2 | 10.02 | |
| | | | | | | (-73.46%) | (-83.80%) | (-89.56%) | (-92.94%) | (-95.02%) |
| | | | | | 199.13 | 93.88 | 47.25 | 25.74 | 15.17 | |
| | | | | | | (-52.86%) | (-76.27%) | (-87.07%) | (-92.38%) | |
| | 4 | 13 | 2 | 0.65 | 98.83 | 52.92 | 30.22 | 18.43 | 11.93 | |
| | | | | | | (-50.37%) | (-73.42%) | (-84.82%) | (-90.74%) | (-94.01%) |
| | | | | | 52.41 | 31.85 | 20.41 | 13.74 | 9.65 | |
| | | | | | | (-73.68%) | (-84.01%) | (-89.75%) | (-93.10%) | (-95.15%) |
| 4 | 13 | 5 | 0.6 | 193.8 | 90.55 | 44.92 | 24.06 | 13.94 | | |
| | | | | | (-53.28%) | (-76.82%) | (-87.59%) | (-92.81%) | | |
| | | | | 95.41 | 50.48 | 28.43 | 17.09 | 10.92 | | |
| | | | | | (-50.77%) | (-73.95%) | (-85.33%) | (-91.18%) | (-94.37%) | |
| 4 | 13 | 5 | 0.7 | 50.04 | 30.08 | 19.05 | 12.68 | 8.81 | | |
| | | | | | (-74.18%) | (-84.48%) | (-90.17%) | (-93.46%) | (-95.45%) | |
| | | | | 212.31 | 99.37 | 50.12 | 27.38 | 16.17 | | |
| | | | | | (-53.20%) | (-76.39%) | (-87.10%) | (-92.38%) | | |
| 0.75 | 7 | 11 | 0 | 0.8 | 77.41 | 43.89 | 26.55 | 17.05 | 11.53 | |
| | | | | | | (-63.54%) | (-79.33%) | (-87.49%) | (-91.97%) | (-94.57%) |
| | | | | | 32.49 | 22.18 | 15.77 | 11.6 | 8.74 | |
| | | | | | | (-84.70%) | (-89.55%) | (-92.57%) | (-94.54%) | (-95.88%) |
| | 7 | 11 | 2 | 0.75 | 211.14 | 98.63 | 49.62 | 27.03 | 15.9 | |
| | | | | | | (-53.29%) | (-76.50%) | (-87.20%) | (-92.47%) | |
| | | | | | 76.78 | 43.45 | 26.22 | 16.79 | 11.31 | |
| | | | | | | (-63.63%) | (-79.42%) | (-87.58%) | (-92.05%) | (-94.64%) |
| | 7 | 11 | 2 | 0.85 | 32.15 | 21.92 | 15.55 | 11.4 | 8.56 | |
| | | | | | | (-84.77%) | (-89.62%) | (-92.64%) | (-94.60%) | (-95.95%) |
| | | | | | 207.88 | 96.61 | 48.29 | 26.09 | 15.2 | |
| | | | | | | (-0.00%) | (-53.53%) | (-76.77%) | (-87.45%) | (-92.69%) |
| 7 | 11 | 5 | 0.8 | 75.14 | 42.31 | 25.38 | 16.14 | 10.77 | | |
| | | | | | (-63.85%) | (-79.65%) | (-87.79%) | (-92.24%) | (-94.82%) | |
| | | | | 31.31 | 21.25 | 14.99 | 10.92 | 8.13 | | |
| | | | | | (-84.94%) | (-89.78%) | (-92.79%) | (-94.75%) | (-96.09%) | |
| 0.85 | 10 | 19 | 0 | 0.85 | 200.17 | 90.73 | 46.06 | 26.07 | 16.2 | |
| | | | | | | (-54.67%) | (-76.99%) | (-86.98%) | (-91.91%) | |
| | | | | | 39.95 | 26.58 | 18.71 | 13.75 | 10.42 | |
| | | | | | | (-80.04%) | (-86.72%) | (-90.65%) | (-93.13%) | (-94.79%) |
| | 10 | 19 | 0 | 0.95 | 15.41 | 12.9 | 10.87 | 9.19 | 7.75 | |
| | | | | | | (-92.30%) | (-93.56%) | (-94.57%) | (-95.41%) | (-96.13%) |
| | | | | | 199.3 | 90.24 | 45.75 | 25.85 | 16.03 | |
| | | | | | | (-0.00%) | (-54.72%) | (-77.05%) | (-87.03%) | (-91.96%) |
| | 10 | 19 | 2 | 0.9 | 39.7 | 26.39 | 18.55 | 13.61 | 10.29 | |
| | | | | | | (-80.08%) | (-86.76%) | (-90.69%) | (-93.17%) | (-94.84%) |
| | | | | | 15.28 | 12.78 | 10.76 | 9.08 | 7.65 | |
| | | | | | | (-92.33%) | (-93.59%) | (-94.60%) | (-95.44%) | (-96.16%) |
| 10 | 19 | 5 | 0.85 | 197.43 | 89.17 | 45.07 | 25.37 | 15.66 | | |
| | | | | | (-54.84%) | (-77.17%) | (-87.15%) | (-92.07%) | | |
| | | | | 39.14 | 25.97 | 18.2 | 13.3 | 10.01 | | |
| | | | | | (-80.18%) | (-86.85%) | (-90.78%) | (-93.26%) | (-94.93%) | |
| 10 | 19 | 5 | 0.95 | 15.01 | 12.52 | 10.51 | 8.83 | 7.41 | | |
| | | | | | (-92.40%) | (-93.66%) | (-94.68%) | (-95.53%) | (-96.25%) | |

6.2. Real data analysis

To illustrate the practical use of the CUSUM control chart, it is employed the disorderly conduct data series from the 44th police car beat in Pittsburgh. The data set comprises monthly observations spanning from 1990 to 2001. Initially, it is used the data spanning from 1990 to 1996 to model the INAR(1)DNXL process. Subsequently, the CUSUM control chart is implemented on the data from 1997 onwards to identify any substantial increase in the mean.

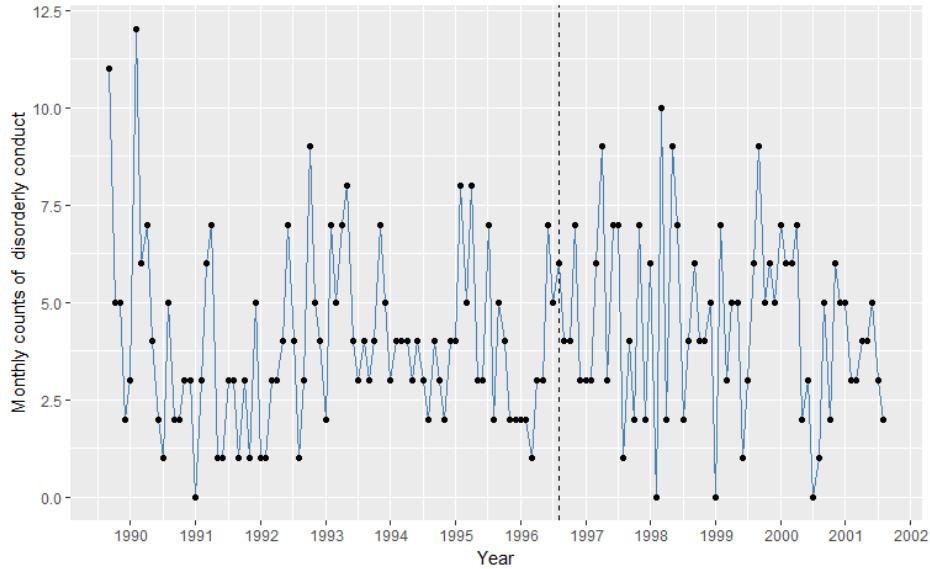


Figure 8: Time series plot of the disorderly conduct data set.

Figure 8 represents the time series plot of the disorderly conduct data set, with the dashed line indicating December 1996. From January 1990 to December 1996, the sample mean is 3.9643 and from January 1997 to May 2000 is 4.818. From Figure 8, it is noticeable that detecting a mean increase after the dashed line is not straightforward. Figure 9 displays the ACF and PACF plots for the Phase I data. From the latter figure, it is clear that only the initial lag is notable in the PACF plot and that the ACF plot exhibits exponential decay, making it suitable to model the data as an INAR(1) process. Accordingly, the Phase I data set is modelled using the INAR(1)DNXL model and other competing processes.

The considered competing models are the INAR(1) process with Poisson marginals denoted by INARP(1) (Al-Osh and Alzaid (1987)), INAR(1) process with geometric innovations denoted by INAR(1)G (Aghababaei Jazi et al. (2022)), INAR(1) process with geometric marginals denoted by INARG(1) (Alzaid and Al-Osh (1988)), INAR(1) process with negative binomial marginals denoted by INARNB(1) (McKenzie (1986)), negative binomial thinning-based INAR(1) process with geometric marginals denoted by NBINARG(1) (Ristić et al. (2009)), INAR(1) process with zero-inflated Poisson innovations denoted by INAR(1)ZP (Jazi et al. (2012)) and INAR(1) process with Katz family innovations denoted by INAR(1)KF (Kim and Lee (2017)).

Table 14 reports the CML estimates, $-\log L$, AIC, and BIC of the aforementioned processes. The INAR(1)DNXL process exhibits the highest log-likelihood value and the lowest AIC and BIC values among the competing models, indicating that it provides a better fit.

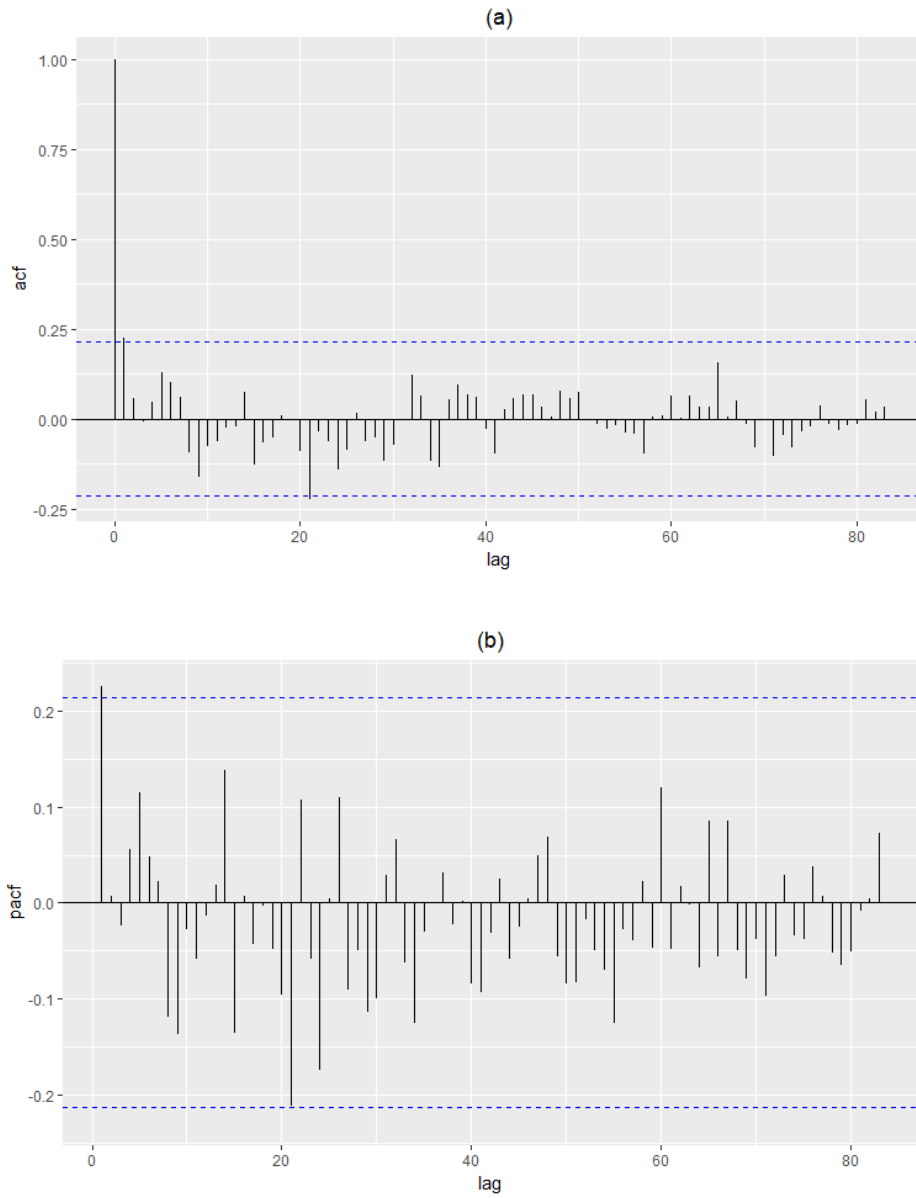


Figure 9: ACF (a) and PACF (b) plots for the Phase I data set.

Table 14: CML estimates, $-\log L$, AIC, and BIC of processes fitted to Phase I data set.

| Model | Parameters | CML est. | $-\log L$ | AIC | BIC |
|-------------|------------|----------------------|-----------|---------|---------|
| INAR(1)DNXL | α | 0.459 | -177.024 | 358.048 | 362.910 |
| | θ | 0.595 | | | |
| INAR(1)G | π | 0.324 | -183.025 | 370.050 | 374.910 |
| | α | 0.500 | | | |
| INARG(1) | p | 0.257 | -196.050 | 396.100 | 400.960 |
| | ρ | 0.518 | | | |
| INARP(1) | λ | 3.133 | -181.365 | 366.730 | 371.590 |
| | α | 0.215 | | | |
| INARNB(1) | n | 11.639 | -179.330 | 364.660 | 371.950 |
| | p | 0.744 | | | |
| | ρ | 0.280 | | | |
| NBINARG(1) | μ | 4.307 | -192.195 | 388.390 | 393.250 |
| | α | 0.799 | | | |
| INAR(1)ZP | λ | 3.188 | -181.150 | 368.300 | 375.590 |
| | ρ | 0.6×10^{-7} | | | |
| | α | 0.211 | | | |
| INAR(1)KF | θ_1 | 2.208 | -179.145 | 364.290 | 371.580 |
| | θ_2 | 0.254 | | | |
| | α | 0.251 | | | |

To assess the statistical accuracy of the fitted INAR(1)DNXL process, residual analysis is conducted using Pearson residuals. Figure 10 exhibits the ACF of the Pearson residuals, indicating the absence of autocorrelation. Additionally, the Ljung–Box test for the existence of autocorrelation is conducted with 10 degrees of freedom, yielding a p -value of 0.8625. This result suggests that the residuals are uncorrelated. Furthermore, the mean and variance of the Pearson residuals are 0.02024 and 0.824, respectively, which are in proximity to the desired values 0 and 1. Accordingly, the INAR(1)DNXL process adequately captures the characteristics of the data set under consideration. Figure 11 represents the cpggram of Pearson residuals of the Phase I data set and the residuals clearly display a random behaviour. Figure 12 displays the time series plot of the original versus predicted values of the data.

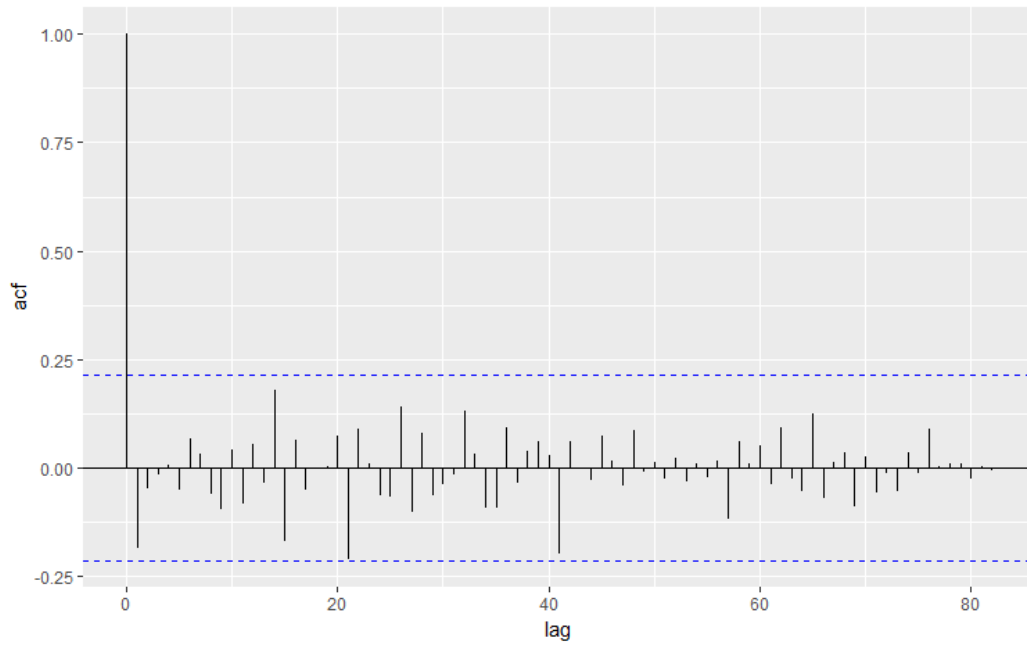


Figure 10: ACF plot of the Pearson residuals for the Phase I data set.

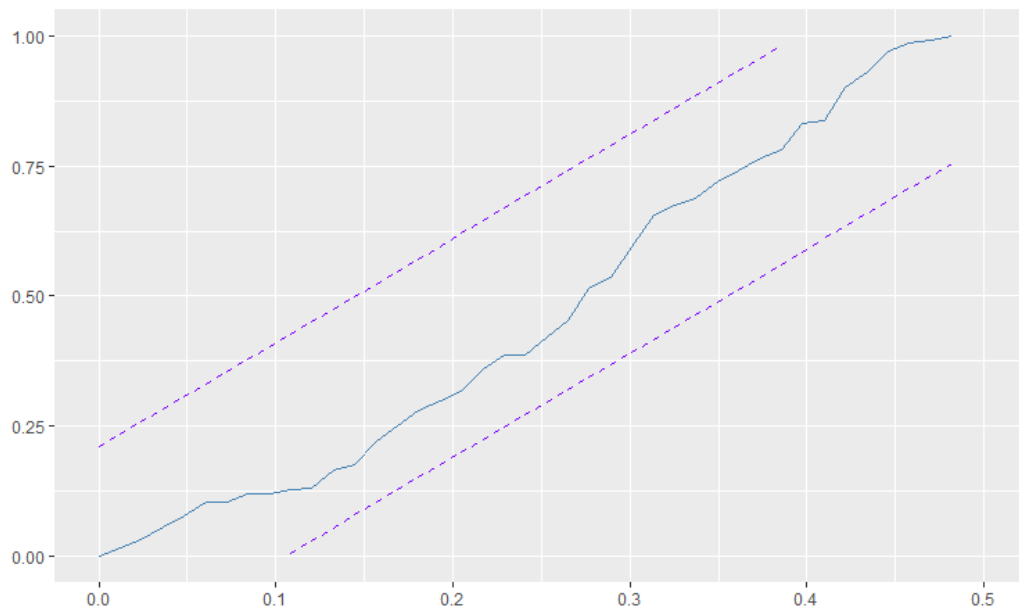


Figure 11: Cpgram of Pearson residuals for the Phase I data set.

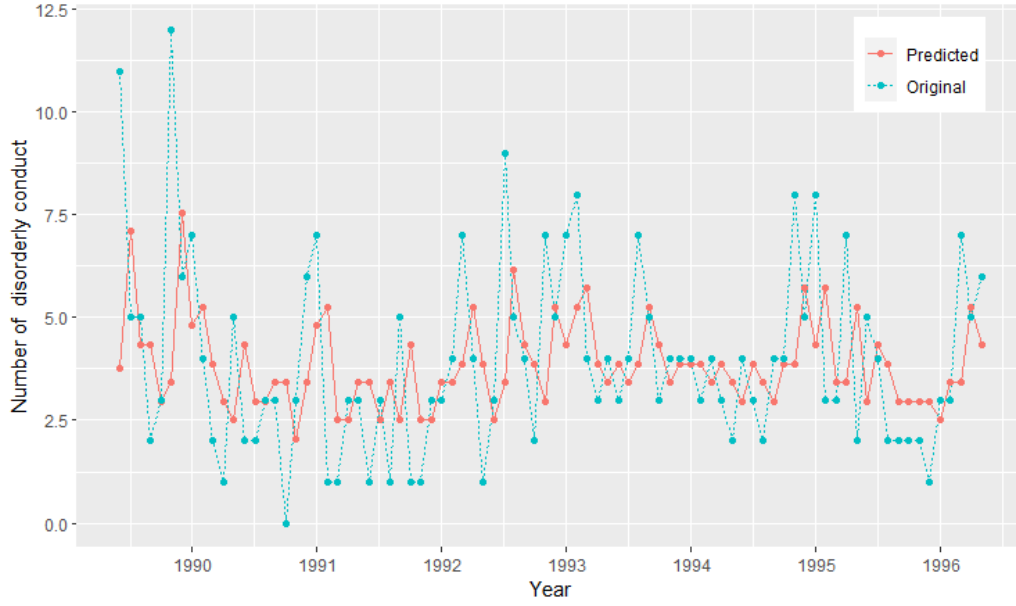


Figure 12: Predicted versus original values of the monthly counts of disorderly conduct of the Phase I data set.

The INAR(1)DNXL model of the Phase I data set is given by

$$X_t = 0.459 X_{t-1} + \epsilon_t, \quad t = 2, 3, \dots, T,$$

where $\epsilon_t \sim \text{DNXL}(0.595)$. The predicted values can be obtained by

$$\begin{aligned} \hat{X}_1 &= E(X_1) = 3.7820, \\ \hat{X}_t &= E(X_t|X_{t-1}) = 0.459 X_{t-1} + 2.0454, \quad t = 2, 3, \dots, T. \end{aligned}$$

Now, it is considered the CUSUM chart with reference value $k = 4$ (see [Kim and Lee \(2017\)](#) for more details) based on the INAR(1)DNXL process with parameter values ($\alpha = 0.459, \theta = 0.595$) and a suitable limit value h is chosen. Table 15 shows the ARL values corresponding to the INAR(1)DNXL, INAR(1)KF, and INARP(1) processes setting $k = 4$ and $c_0 = 0$. The INAR(1)DNXL process demonstrates a superior performance compared to the others, exhibiting lower ARL values. By setting $ARL_0 = 100$, the upper control limit for the CUSUM chart under the INAR(1)DNXL process is determined as $h = 28$.

Table 15: ARL values of CUSUM chart for disorderly conduct data with $k = 4$.

| h | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 |
|-------------|-------|-------|------|-------|-------|-------|-------|-------|-------|
| INAR(1)DNXL | 78.0 | 83.1 | 88.4 | 93.9 | 99.6 | 105.5 | 111.6 | 118.0 | 124.5 |
| INAR(1)KF | 108.4 | 116.6 | 125 | 133.9 | 143 | 152.6 | 162.4 | 172.6 | 183.2 |
| INARP(1) | 126 | 135.3 | 145 | 155.1 | 165.5 | 176.2 | 187.3 | 198.8 | 210.5 |

The CUSUM chart for disorderly conduct data under INAR(1)DNXL process is given in Figure 13, where the horizontal and vertical dashed lines stand for $h = 28$ and December 1996, respectively.

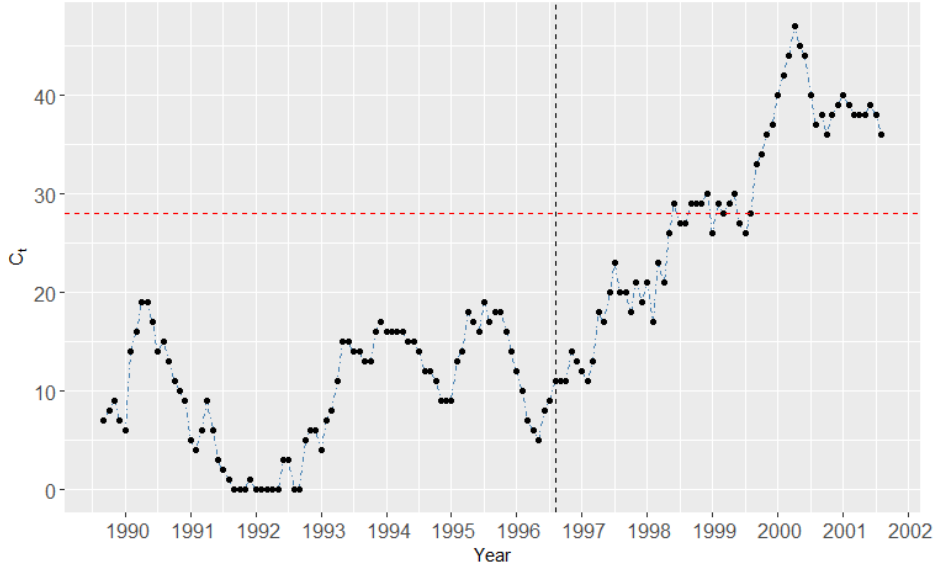


Figure 13: CUSUM plot of the monthly counts of disorderly conduct data.

From the above figure, it can be seen that $C_t \geq h$ occurred in October 1998. As said before, the sample mean of the data from January 1997 to May 2000 is 4.818, which is greater than that of the past observations. Clearly, it is difficult to identify the change in the mean unless the CUSUM control scheme has been implemented.

7. Conclusions and future research

This paper introduces a discrete analogue of the continuous new XLindley (NXL) distribution proposed by [Nawel et al. \(2023\)](#). The new model has been obtained by the survival discretization method and it is referred to as the discrete new XLindley distribution (DNXL). One remarkable property of the DNXL distribution is that the pmf and cdf are very simple, which implies that its most relevant statistical properties can be given in closed form. Moreover, it is suitable to describe both overdispersed and underdispersed count data and it is characterised by an increasing failure rate. Additionally, a new INAR(1) process with DNXL innovations is also constructed. The parameter estimation of the DNXL distribution and the INAR(1)DNXL process can be easily performed by the methods of maximum likelihood and conditional maximum likelihood, respectively. Several real data sets illustrate the usefulness of the novel discrete distribution and its associated INAR(1) process for modelling count data and counts of time series. Additionally, a CUSUM control chart is developed to detect the increase in mean of autocorrelated count processes and its usefulness is illustrated by a real data example. As future work, it could be interesting to examine other models obtained from the NXL distribution by means of different discretization techniques. Furthermore, a bivariate discrete DNXL distribution together with its associated BINAR(1) process may also be studied.

Acknowledgments

The authors would like to thank the editor and referees for their careful reading and comments which greatly improved the paper.

References

- M. Aghababaei Jazi, G. Jones, and C.-D. Lai. Integer valued AR (1) with geometric innovations. *Journal of the Iranian Statistical Society*, 11(2):173–190, 2022.
- A. A. Al-Babtain, A. M. Gemeay, and A. Z. Afify. Estimation methods for the discrete Poisson-Lindley and discrete Lindley distributions with actuarial measures and applications in medicine. *Journal of King Saud University-Science*, 33(2):101224, 2021.
- M. A. Al-Osh and A. A. Alzaid. First-order integer-valued autoregressive (INAR (1)) process. *Journal of Time Series Analysis*, 8(3):261–275, 1987.
- H. M. Aljohani, M. Ahsan-ul Haq, J. Zafar, E. M. Almetwally, A. S. Alghamdi, E. Husam, and A. H. Muse. Analysis of covid-19 data using discrete Marshall–Olkin length biased exponential: Bayesian and frequentist approach. *Scientific Reports*, 13(1):12243, 2023.
- A. D. H. E. J. T. S.-A. N. . A. H. Almetwally, E.M. The new discrete distribution with application to COVID-19 data. *Results in Physics*, 32:104987, 2022.
- R. Alotaibi, E. M. Almetwally, and H. Rezk. Optimal test plan of discrete alpha power inverse Weibull distribution under censored data. *Journal of Radiation Research and Applied Sciences*, 16(2):100573, 2023.
- E. Altun. A new generalization of geometric distribution with properties and applications. *Communications in Statistics-Simulation and Computation*, 49(3):793–807, 2020.
- E. Altun, D. Bhati, and N. M. Khan. A new approach to model the counts of earthquakes: INAR_{PQX} (1) process. *SN Applied Sciences*, 3:1–17, 2021.
- E. Altun, M. El-Morshedy, and M. Eliwa. A study on discrete Bilal distribution with properties and applications on integervalued autoregressive process. *REVSTAT-Statistical Journal*, 20(4):501–528, 2022.
- A. Alzaid and M. Al-Osh. First-order integer-valued autoregressive (INAR (1)) process: distributional and regression properties. *Statistica Neerlandica*, 42(1):53–61, 1988.
- G. J. Babu and C. R. Rao. Goodness-of-fit tests when parameters are estimated. *Sankhyā: The Indian Journal of Statistics*, pages 63–74, 2004.
- R. E. Barlow and F. Proschan. *Statistical theory of reliability and life testing: probability models*, volume 1. Holt, Rinehart and Winston New York (1975), 1975.
- A. Beghriche, H. Zeghdoudi, V. Raman, and S. Chouia. New polynomial exponential distribution: properties and applications. *Statistics in Transition new series*, 23(3): 95–112, 2022.
- A. Beghriche, Y. A. Tashkandy, M. Bakr, Z. Halim, A. M. Gemeay, M. M. Hossain, and A. H. Muse. The inverse XLindley distribution: Properties and application. *IEEE Access*, 11:47272–47281, 2023.
- D. Brook and D. Evans. An approach to the probability distribution of CUSUM run length. *Biometrika*, 59(3):539–549, 1972.

- S. Chakraborty. Generating discrete analogues of continuous probability distributions—a survey of methods and constructions. *Journal of Statistical Distributions and Applications*, 2:1–30, 2015.
- H. Chouia, S. & Zeghdoudi. The XLindley distribution: properties and application. *Journal of Statistical Theory and Applications*, 20:318–327, 2021.
- R. M. Corless, G. H. Gonnet, D. E. Hare, D. J. Jeffrey, and D. E. Knuth. On the Lambert W function. *Advances in Computational mathematics*, 5:329–359, 1996.
- M. El-Morshedy, M. S. Eliwa, and E. Altun. Discrete Burr-Hatke distribution with properties, estimation methods and regression model. *IEEE access*, 8:74359–74370, 2020.
- A. S. Eldeeb, M. Ahsan-ul Haq, and A. Babar. A new discrete XLindley distribution: Theory, actuarial measures, inference, and applications. *International Journal of Data Science and Analytics*, pages 1–11, 2023.
- T. Ghosh, D. Roy, and N. K. Chandra. Reliability approximation through the discretization of random variables using reversed hazard rate function. *International Journal of Mathematical, Computational, Statistical, Natural and Physical Engineering*, 7(4): 96–100, 2013.
- E. Gómez-Déniz and E. Calderín-Ojeda. The discrete Lindley distribution: properties and applications. *Journal of statistical computation and simulation*, 81(11):1405–1416, 2011.
- I. J. Good. The population frequencies of species and the estimation of population parameters. *Biometrika*, 40(3-4):237–264, 1953.
- H. Haj Ahmad and E. M. Almetwally. Generating optimal discrete analogue of the generalized Pareto distribution under bayesian inference with applications. *Symmetry*, 14(7):1457, 2022.
- A. Harvey and C. Fernandes. Time series models for count or qualitative observations: Reply. *Journal of Business & Economic Statistics*, 7(4), 1989.
- J. Huang and F. Zhu. A new first-order integer-valued autoregressive model with Bell innovations. *Entropy*, 23(6):713, 2021.
- M. Irshad, P. Jodrá, A. Krishna, and R. Maya. On the discrete analogue of the Teissier distribution and its associated inar (1) process. *Mathematics and Computers in Simulation*, 214:227–245, 2023.
- M. R. Irshad, C. Chesneau, V. D’cruz, and R. Maya. Discrete pseudo Lindley distribution: Properties, estimation and application on inar (1) process. *Mathematical and Computational Applications*, 26(4):76, 2021.
- M. A. Jazi, C.-D. Lai, and M. H. Alamatsaz. A discrete inverse Weibull distribution and estimation of its parameters. *Statistical Methodology*, 7(2):121–132, 2010.
- M. A. Jazi, G. Jones, and C.-D. Lai. First-order integer valued AR processes with zero inflated Poisson innovations. *Journal of Time Series Analysis*, 33(6):954–963, 2012.

- P. Jodrá. Computer generation of random variables with Lindley or Poisson–Lindley distribution via the Lambert W function. *Mathematics and Computers in Simulation*, 81(4):851–859, 2010.
- J. Keilson and H. Gerber. Some results for discrete unimodality. *Journal of the American Statistical Association*, 66(334):386–389, 1971.
- H. Kim and S. Lee. On first-order integer-valued autoregressive process with Katz family innovations. *Journal of Statistical Computation and Simulation*, 87(3):546–562, 2017.
- L. A. Klimko and P. I. Nelson. On conditional least squares estimation for stochastic processes. *The Annals of statistics*, pages 629–642, 1978.
- H. Krishna and P. S. Pundir. Discrete Burr and discrete Pareto distributions. *Statistical methodology*, 6(2):177–188, 2009.
- K. Kulasekera and D. W. Tonkyn. A new discrete distribution, with applications to survival, dispersal and dispersion. *Communications in Statistics-Simulation and Computation*, 21(2):499–518, 1992.
- C. D. Lai and M. Xie. *Stochastic ageing and dependence for reliability*. Springer Science and Business Media, 2006.
- J. F. Lawless. *Statistical models and methods for lifetime data*. John Wiley & Sons, 2011.
- C. Li, H. Zhang, and D. Wang. Modelling and monitoring of INAR (1) process with geometrically inflated Poisson innovations. *Journal of Applied Statistics*, 49(7):1821–1847, 2022.
- T. Lívio, N. M. Khan, M. Bourguignon, and H. S. Bakouch. An INAR (1) model with Poisson–Lindley innovations. *Econ Bull*, 38(3):1505–1513, 2018.
- T. Mäkeläinen, K. Schmidt, and G. P. Styan. On the existence and uniqueness of the maximum likelihood estimate of a vector-valued parameter in fixed-size samples. *The Annals of Statistics*, pages 758–767, 1981.
- E. McKenzie. Some simple models for discrete variate time series 1. *JAWRA Journal of the American Water Resources Association*, 21(4):645–650, 1985.
- E. McKenzie. Autoregressive moving-average processes with negative-binomial and geometric marginal distributions. *Advances in Applied Probability*, 18(3):679–705, 1986.
- T. Nakagawa and S. Osaki. The discrete Weibull distribution. *IEEE transactions on reliability*, 24(5):300–301, 1975.
- K. Nawel, A. M. Gemeay, H. Zeghdoudi, K. Karakaya, A. M. Alshangiti, M. Bakr, O. S. Balogun, A. H. Muse, and E. Hussam. Modeling voltage real data set by a new version of Lindley distribution. *IEEE Access*, 11:67220–67229, 2023.
- E. Page. Cumulative sum charts. *Technometrics*, 3(1):1–9, 1961.
- B. Para and T. Jan. On discrete three-parameter Burr type XII and discrete Lomax distributions and their applications to model count data from medical science. *Biometrics and Biostatistics International Journal*, 4(2):1–15, 2016.

- X. Qi, Q. Li, and F. Zhu. Modeling time series of count with excess zeros and ones based on INAR (1) model with zero-and-one inflated Poisson innovations. *Journal of Computational and Applied Mathematics*, 346:572–590, 2019.
- A. C. Rakitzis, C. H. Weiß, and P. Castagliola. Control charts for monitoring correlated Poisson counts with an excessive number of zeros. *Quality and Reliability Engineering International*, 33(2):413–430, 2017.
- M. M. Ristić, H. S. Bakouch, and A. S. Nastić. A new geometric first-order integer-valued autoregressive (NGINAR (1)) process. *Journal of Statistical Planning and Inference*, 139(7):2218–2226, 2009.
- D. Roy. Discrete Rayleigh distribution. *IEEE transactions on reliability*, 53(2):255–260, 2004.
- H. Sato, M. Ikota, A. Sugimoto, and H. Masuda. A new defect distribution metrology with a consistent discrete exponential formula and its applications. *IEEE Transactions on Semiconductor Manufacturing*, 12(4):409–418, 1999.
- D. A. Sprott. Estimating the parameters of a convolution by maximum likelihood. *Journal of the American Statistical Association*, 78(382):457–460, 1983.
- E. A. . G. A. Teamah, A.M. Discrete half-logistic distribution: Statistical properties, estimation, and application. *Journal of Statistics Applications and Probability*, 13: 273–284, 2024.
- C. H. Weiß. *An introduction to discrete-valued time series*. John Wiley & Sons.
- C. H. Weiss and M. C. Testik. CUSUM monitoring of first-order integer-valued autoregressive processes of Poisson counts. *Journal of quality technology*, 41(4):389–400, 2009.