

Latent Gaussian Count Time Series

Yisu Jia

University of North Florida

Stefanos Kechagias

SAS Institute

James Livsey

United States Census Bureau

Robert Lund

University of California - Santa Cruz

Vladas Pipiras

University of North Carolina - Chapel Hill

June 7, 2021

Abstract

Keywords: Count Distributions; Hermite Expansions; Likelihood Estimation; Particle Filtering; Sequential Monte Carlo; State Space Models

*

†

1 Introduction

This paper develops the theory and methods for modeling a stationary discrete-valued time series by transforming a Gaussian process. Since the majority of discrete-valued time series involve integer counts supported on some subset of $\{0, 1, \dots\}$, we isolate on this support set. Our methods are based on a copula-style transformation of a latent Gaussian stationary series and are able to produce any desired count marginal distribution. It is shown that the proposed model class produces the most flexible pairwise correlation structures possible, including negatively dependent series. Model parameters are estimated via 1) a Gaussian pseudo-likelihood approach, developed from some new Hermite expansion techniques, which use only the mean and the autocovariance of the series, 2) an implied Yule-Walker moment estimation approach when the latent Gaussian process is an autoregression, and 3) a particle filtering (PF) / sequential Monte Carlo (SMC) approach that uses a state space model (SSM) representation of the transformation to approximate the true likelihood. Extensions to non-stationary settings, particularly those with covariates, are discussed.

The theory of stationary Gaussian time series is by now well developed. A central result is that a stationary Gaussian series $\{X_t\}_{t \in \mathbb{Z}}$ having the lag- h autocovariance $\gamma(h) = \text{Cov}(X_t, X_{t+h})$ exist if and only if γ is symmetric about lag zero and non-negative definite (see Theorem 1.5.1 in [6]). However, such a result does not hold for stationary count series having a certain prescribed marginal distribution (e.g, Poisson). In principle, distributional existence issues are checked with Kolmogorov's consistency criterion (see Theorem 1.2.1 in [6]); in practice, one needs a specified joint distribution to check for consistency. Phrased another way, Kolmogorov's consistency criterion is not a constructive result and does not illuminate how to build stationary time series having a particular marginal distribution and correlation structure. Perhaps owing to this, count time series have been constructed

from a plethora of approaches over the years, as is next reviewed.

Drawing from the success of autoregressive moving-average (ARMA) models in describing stationary Gaussian series, early count authors constructed correlated count series from discrete ARMA (DARMA) and integer ARMA (INARMA) difference equation methods. Focusing on the first order autoregressive case for simplicity, a DAR(1) series X_t with specified marginal distribution $f(x)$ is obtained by generating X_1 from $f(x)$ and then at each subsequent time, either keeping the previous count value with probability p or generating an independent copy of $f(x)$ with probability $1-p$. INAR(1) series are built via the thinned AR(1) equation $X_t = \alpha X_{t-1} + Y_t$, where Y_t is an IID count-valued random sequence and α is a thinning operator defined by $\alpha X = \sum_{i=1}^X B_i$ for a binomial distribution $B(n, p)$ with n trials and success probability p . DARMA methods were initially explored in [23], but were subsequently discarded by practitioners because their sample paths often remained constant for long periods, especially in highly correlated cases; INARMA series are still used today. In contrast to their Gaussian ARMA brethren, DARMA and INARMA models, and their extensions in [24], cannot produce negative autocorrelations.

The works [5] and [9] take a different approach, producing the desired count marginal distribution by combining IID copies of a correlated Bernoulli series Z_t built from a stationary renewal sequence. Explicit autocovariance functions when Z_t is made by binning (clipping) a stationary Gaussian sequence into zero-one categories are derived in [33]. While these models can have negative correlations, they do not necessarily produce the most negatively correlated count structures possible. Also, some important count marginal distributions, including generalized Poisson, are not easily built from these methods. The results here easily generate any desired count marginal distribution. Other count model classes studied include Gaussian processes rounded to their nearest integer [26], hierarchical

Bayesian count model approaches [2], and others (see [18] and [11] for recent reviews). Each approach has some drawbacks.

The models here impose a fixed marginal distribution for the counts. This is in contrast to generalized ARMA methods (GLARMA), which typically posit conditional distributions in lieu of marginal distributions, with model parameters typically being random. As [1] shows in the Poisson case, once the randomness of the parameters is taken into account, the true marginal distribution of the series can be far from the posited conditional distribution. This said, the literature on GLARMA and other conditional models is extensive [3, 43]. See [16] for a recent review of GLARMA models.

A time series analyst generally needs four features in a count model: 1) general marginal distributions; 2) the most general correlation structures possible, both positive and negative; 3) the straight-forward accomodation of covariates; and 4) a well performing and computationally feasible likelihood inference approach. All previous count classes fail to accommodate one or more of these tenets. This paper s purpose is to introduce and study a count model class that, for the first time, simultaneously achieves all four features. Our model employs a latent Gaussian process and a copula-style transformation. This type of construction has recently shown promise in spatial statistics [12, 21], multivariate modeling [39, 40], and regression [35], but the theory has yet to be developed for count series ([35, 30] provide some partial results). Our objectives here are several-fold. On a methodological level, it is shown, through some newly derived Hermite polynomial expansions, that accurate and efficient numerical quantification of the correlation structure of this count model class is feasible. Based on a result in [42], the class is shown to produce the most flexible pairwise correlation structures possible, positive or negative (see Remark 2.2 below). Connections to both importance sampling schemes, where the popular GHK sampler in [35] is

adapted to our needs, and to the SSM and SMC literature, which allow natural extensions of the GHK sampler and likelihood evaluation, are made. The methods are tested on both synthetic and real data.

The works [35, 30] are perhaps the closest papers to this study. While the general latent Gaussian construct adopted is the same, our work differs in that explicit autocovariance relations are developed via Hermite expansions, flexibility and optimality issues of the model class are addressed, Gaussian pseudo-likelihood and implied least-squares parameter estimation approaches are developed, and both the importance sampling and SSM connections are explored in detail. Additional connections to [35, 30] and to the spatial count modeling papers [21, 22] are later made.

The rest of this paper proceeds as follows. The next section and Appendix A introduce our Gaussian transformation count model and establish its basic mathematical and statistical properties. Section 3 and Appendix B move to estimation, developing three techniques: a Gaussian pseudo-likelihood approach, implied Yule-Walker estimation, and PF/SMC methods. Section 4 and Appendix C present simulation results. Section 5 and Appendix D analyze soft drink sales counts at one location of the now defunct Dominicks Finer Foods retail chain. This series exhibits overdispersion, negative lag one autocorrelation, and dependence on a price reduction (sales) covariate, which illustrates the flexibility of our approach. Section 6 concludes with comments and suggestions for future research.

2 Theory

We seek to construct a strictly stationary time series $\{Y_t\}_{t \in \mathbb{Z}}$ having marginal distributions from any family of count distributions supported in $[0, 1]$, including the Binomial,

Poisson, mixture Poisson, negative binomial, generalized Poisson, and Conway-Maxwell-Poisson distributions. The later three distributions are over-dispersed (their variances are larger than their respective means), which is the case for many observed count time series.

Let $\{Y_t\}_{t \in \mathbb{Z}}$ be the stationary count time series of interest. Suppose that one wants the marginal cumulative distribution function (CDF) of Y_t for each t of interest to be $F_Y(y) = \mathbb{P}[Y_t \leq y]$, depending on a vector $\boldsymbol{\theta}$ containing all CDF model parameters. The series $\{Y_t\}_{t \in \mathbb{Z}}$ will be modeled through

$$Y_t = \lfloor \Phi^{-1}(F_Y(Y_t)) \rfloor \quad \text{where} \quad \Phi^{-1}(u) = \inf\{x \in \mathbb{R} : \Phi(x) \geq u\} \quad (1)$$

and $\Phi(\cdot)$ is the CDF of a standard normal variable and $\Phi^{-1}(u) = \inf\{x \in \mathbb{R} : \Phi(x) \geq u\}$, $(0, 1)$, is the generalized inverse (quantile function) of the CDF $F_Y(\cdot)$. The process $\{Z_t\}_{t \in \mathbb{Z}}$ is standard Gaussian for each fixed t , but possibly correlated in time:

$$\mathbb{E}[Z_t] = 0 \quad \mathbb{E}[Z_t^2] = 1 \quad \rho_{t,s} =: \text{Corr}(Z_t, Z_s) = \mathbb{E}[Z_t Z_s] \quad (2)$$

This approach has been used in [39, 35, 21, 30] with good results. The autocovariance function (ACVF) of $\{Z_t\}_{t \in \mathbb{Z}}$, denoted by $\gamma_{Z_t}(h)$, is the same as the autocorrelation function (ACF) due to standard normality and depends on another vector $\boldsymbol{\eta}$ of ACVF parameters.

As expanded on in Section 2.3, (1) can be viewed as a SSM:

State equation : $\{Z_t\}_{t \in \mathbb{Z}}$ governing latent Gaussian dynamics;

Observation equation : $\mathbb{P}(Y_t = y) = \int \mathbb{P}(Y_t = y | Z_t = z) \phi(z) dz$ with the set \mathcal{Y} defined below.

Here, $\phi(\cdot)$ is notation for an arbitrary conditional distribution.

This model has alternative names in other literature. For example, [7] call this setup the normal to anything (NORTA) procedure in operations research, whereas [20] calls this

a translational model in mechanical engineering. Our goal is to give a reasonably complete analysis of the probabilistic and statistical properties of these models.

The construction in (1) ensures that the marginal CDF of X_t is indeed $F(x)$. Elaborating, the probability integral transformation theorem shows that X_t has a uniform distribution over $(0, 1)$ for each t ; a second application of the result justifies that X_t has the marginal distribution $F(x)$ for each t . Moreover, temporal dependence in X_t will induce temporal dependence in $\gamma(t, s)$ as quantified below. For notation, let $\gamma(t, s) = \mathbb{E}[X_{t+s} - X_t] = \mathbb{E}[X_{t+s}] - \mathbb{E}[X_t]$ denote the ACVF of X_t .

2.1 Relationship between autocovariances

The autocovariance functions of X_t and Y_t can be related using Hermite expansions (see Chapter 5 of [37]). In particular, using the Hermite polynomials $H_k(x) = (-1)^k (1-x^2)^{k/2} \frac{d^k}{dx^k} (1-x^2)^{-k/2}$, \mathbb{R} we can expand the $\gamma(t, s)$ function as

$$\gamma(t, s) = \mathbb{E}[X_{t+s} - X_t] + \sum_{k=1}^{\infty} \frac{1}{k!} H_k(X_t) H_k(X_{t+s}) \quad (3)$$

where the *Hermite coefficients* are given by

$$c_k = \frac{1}{k!} \int_{-\infty}^{\infty} H_k(x) H_k(x) \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx = \frac{1}{k!} \mathbb{E}[H_k(X_0) H_k(X_0)] \quad (4)$$

for a standard normal variable X_0 . The relationship between $\gamma(t, s)$ and c_k is key and is extracted from Chapter 5 of [37]:

$$\gamma(t, s) = \sum_{k=1}^{\infty} \frac{1}{k!} c_k \gamma_k(t, s) =: \gamma(t, s) \quad (5)$$

where $\psi(\mathbf{x}) = \sum_{i=1}^{\infty} \frac{1}{i!} \mathbf{x}^i$. For $\mathbf{x} = \mathbf{0}$, (5) yields $\text{Var}(\mathbf{X}) = \psi(\mathbf{0}) = \sum_{i=1}^{\infty} \frac{1}{i!} \mathbf{x}^i$, which depends only on the marginal parameters in $\boldsymbol{\theta}$. Moreover, the ACF of \mathbf{X} is

$$\psi(\mathbf{x}) = \sum_{i=1}^{\infty} \frac{1}{i!} \mathbf{x}^i \quad \psi(\mathbf{x}) =: \psi(\mathbf{x}) \quad (6)$$

where

$$\psi(\mathbf{x}) = \sum_{i=1}^{\infty} \frac{1}{i!} \mathbf{x}^i \quad =: \sum_{i=1}^{\infty} \quad (7)$$

and $\psi(\mathbf{0}) = \sum_{i=1}^{\infty} \frac{1}{i!} \mathbf{x}^i$. The function ψ maps $[0, 1]^d$ into (but not necessarily onto) $[0, 1]^d$. For future reference, note that $\psi(\mathbf{0}) = 0$ and $\psi(\mathbf{1}) = \sum_{i=1}^{\infty} \frac{1}{i!} \mathbf{1}^i = 1$. Using (3) and $\mathbb{E}[\psi(\mathbf{X}_0) \psi(\mathbf{X}_1)] = \psi(\mathbf{1}) \mathbb{1}_{[\mathbf{X}_0 = \mathbf{X}_1]}$ gives $\psi(\mathbf{1}) = \text{Corr}(\psi(\mathbf{X}_0), \psi(\mathbf{X}_1))$; however, $\psi(\mathbf{1})$ is not necessarily 1 in general. As such, ψ starts at $(\mathbf{0}, \psi(\mathbf{0}))$, passes through $(\mathbf{0}, 0)$, and connects to $(\mathbf{1}, 1)$. Examples are given in Figure 2 of Appendix A.

We call the quantity $\psi(\mathbf{x})$ a *link function*, and the coefficients ψ_i , $i \geq 1$, *link coefficients*. (Sometimes, slightly abusing terminology, we also use these terms for $\psi(\mathbf{x})$ and ψ_i , respectively.) A key feature in (5) is that the effects of the marginal CDF $F_i(x_i)$ and the ACVF $\psi(\mathbf{x})$ are decoupled in the sense that the correlation parameters in $\boldsymbol{\theta}$ do not influence the coefficients in (5) — this is useful later in estimation.

Further properties and the numerical calculation of the link function and the Hermite coefficients are discussed in Appendix A. The computation of the Hermite coefficients, in particular, is feasible due to the following lemma, which is proved in Appendix A.

Lemma 2.1. *If $\mathbb{E}[\psi(\mathbf{X})] = \psi(\mathbf{1})$ for some $\boldsymbol{\theta}$, then the coefficients ψ_i satisfy*

$$\psi_i = \frac{1}{i!} \sum_{n=0}^{\infty} \frac{1}{2^n} \psi_{i+n} \psi_{i-n} \quad (8)$$

where $\psi_i = \mathbb{P}[\psi(\mathbf{X}) = i]$. (When $\psi_{i-n} = 0$ (that is, $i-n < 0$ or $i-n > 1$), the summand $\frac{1}{2^n} \psi_{i+n} \psi_{i-n}$ is interpreted as zero.)

Returning to the relationship between $\gamma(h)$ and $\rho(h)$, from (6), one can see that

$$\gamma(h) = \rho(h) \quad (9)$$

which implies that a positive $\rho(h)$ leads to a positive $\gamma(h)$. A negative $\rho(h)$ produces a negative $\gamma(h)$ since $\rho(h)$ is, in fact, monotone increasing (see Proposition A.1 in Appendix A) and crosses zero at $h = 0$ (the negativeness of $\rho(h)$ when $h > 0$ can also be deduced from the nondecreasing nature of $\gamma(h)$ via an inequality on page 20 of [41] for Gaussian variables).

Remark 2.1. The short- and long-range dependence properties of $\gamma(h)$ can be extracted from those of $\rho(h)$. Recall that a time series $\{X_t\}$ is short-range dependent (SRD) if $\sum_{h=-\infty}^{\infty} \gamma(h) < \infty$. According to one definition, a series $\{X_t\}$ is long-range dependent (LRD) if $\gamma(h) = O(h^{-2+\beta})$, where $\beta \in (0, 1/2)$ is the LRD parameter and $L(h)$ is a slowly varying function at infinity [37]. The ACVF of such LRD series satisfies $\sum_{h=-\infty}^{\infty} \gamma(h) = \infty$. If $\rho(h)$ is SRD, then so is $\gamma(h)$ by (9). On the other hand, if $\rho(h)$ is LRD with parameter β , then $\gamma(h)$ can be either LRD or SRD. The conclusion depends, in part, on the Hermite rank of $\rho(h)$, which is defined as $\beta = \min\{1 : \rho(h) = O(h^{-\beta})\}$. Specifically, if $\beta \in (0, 1/2)$, then $\gamma(h)$ is SRD; if $\beta \in (1/2, 1)$, then $\gamma(h)$ is LRD with parameter $(\beta - 1/2) + 1/2$ (see [37], Proposition 5.2.4).

The model in (1) admits the following structure: if $\{X_t\}$ and $\{Y_t\}$ are independent, then so are $\{Z_t\}$ and $\{W_t\}$. It follows that if $\{Z_t\}$ is stationary and β -dependent, then both $\{X_t\}$ and $\{Y_t\}$ must be β th order moving-average time series. Unfortunately, no analogous autoregressive structure holds; in fact, if $\{Z_t\}$ is a first order autoregression, then $\{X_t\}$ may not be an autoregression of any order (this can be inferred from [28]).

Remark 2.2. The construction in (1) yields models with the most flexible correlations possible for $\text{Corr}(X_1, X_2)$ for two variables X_1 and X_2 with the same marginal distribution F . Indeed, let $\rho_- = \min \text{Corr}(X_1, X_2) : X_1, X_2 \sim F$ and define ρ_+ similarly with min replaced by max. Then, as shown in Theorem 2.5 of [42],

$$\rho_+ = \text{Corr}(F^{-1}(U), F^{-1}(U)) = 1 \quad \rho_- = \text{Corr}(F^{-1}(U), F^{-1}(1-U))$$

where U is a uniform random variable over $(0, 1)$. Since $\frac{D}{D} F^{-1}(U) \stackrel{D}{=} Z$ and $1 - \frac{D}{D} F^{-1}(U) \stackrel{D}{=} (-Z)$ for a standard normal random variable Z , the maximum and minimum correlations ρ_+ and ρ_- are indeed achieved with (1) when $X_1 = X_2$ and $X_1 = -X_2$, respectively. The preceding statements are non-trivial for ρ_- only since $\rho_+ = 1$ is attained whenever $X_1 = X_2$. It is worthwhile to compare this to the discussion following (7). Finally, all correlations in $(-\rho_-, \rho_+) = (-1, 1)$ are achievable since $F^{-1}(U)$ in (7) is continuous in U . The flexibility of correlations for Gaussian copula models in the spatial context was also noted and studied in [21], especially in comparison to a class of hierarchical, e.g. Poisson, models.

The preceding remark settles autocovariance flexibility issues for stationary count series. Flexibility is a concern when the series is negatively correlated, an issue arising, for example, with hurricane counts in [33] and chemical process counts in [26]. Since any general count marginal distribution can also be achieved, the model class is quite general.

2.2 Covariates

There are situations where stationarity is not desired. Such scenarios can often be accommodated by simple variants of the above setup. For concreteness, consider a situation where a vector \mathbf{M} of p non-random covariates is available to explain the series at time t . If one wants X_t to have the marginal distribution $\theta_t(\cdot)$, where $\boldsymbol{\theta}(\cdot)$ is a vector-valued

function of $\boldsymbol{\theta}$ containing marginal distribution parameters, then simply set

$$= \boldsymbol{\theta}^{-1}(\boldsymbol{\theta}(\boldsymbol{\theta})) \quad (10)$$

and reason as before. We do not recommend modifying $\boldsymbol{\theta}$ for the covariates as this may bring process existence issues into play.

Generalized linear models link functions (not to be confused with $\boldsymbol{\theta}$ in (6)–(7)) can be used when parametric support set bounds are encountered. For example, a Poisson regression with correlated errors can be formulated via a parameter vector $\boldsymbol{\beta}$ of regression coefficients with $\boldsymbol{\theta}(\boldsymbol{\theta}) = \mathbb{E}[\boldsymbol{\theta}] = \exp(\boldsymbol{\beta}'\boldsymbol{M})$. Here, the exponential link guarantees that the Poisson parameter is positive. The above construct requires the covariates to be non-random; should covariates be random, marginal distributions may change from $\boldsymbol{\theta}(\boldsymbol{\theta})$.

2.3 Particle filtering and state space model connections

This subsection studies the implications of the latent structure of our model, especially as it relates to SSMs and importance sampling approaches. This will be used to construct PF/SMC approximations of various quantities, and in goodness-of-fit assessments. Our main reference is [14]. As in that monograph, let $\boldsymbol{\theta}_0 = \boldsymbol{\theta}_0 = \boldsymbol{\theta}_0 = \boldsymbol{\theta}_0$, $\boldsymbol{\theta}_0 = \boldsymbol{\theta}_0 = \boldsymbol{\theta}_0 = \boldsymbol{\theta}_0$, and $\boldsymbol{\theta}(\boldsymbol{\theta})$ and $\boldsymbol{\theta}(\boldsymbol{\theta})$ denote joint and conditional probabilities (or their densities, depending on the context). For example, $\boldsymbol{\theta}(\boldsymbol{\theta}_0 : \boldsymbol{\theta}_0)$ denotes the conditional density of $\boldsymbol{\theta}_0$ given $\boldsymbol{\theta}_0$. Similarly, let $\mathbb{E}[\boldsymbol{\theta}_0 : \boldsymbol{\theta}_0]$ denote conditional expectation given $\boldsymbol{\theta}_0$. The SSM formulation starts by specifying $\boldsymbol{\theta}(\boldsymbol{\theta}_{+1} : \boldsymbol{\theta}_0)$ and $\boldsymbol{\theta}(\boldsymbol{\theta}_0)$. While $\boldsymbol{\theta}$ is often first order Markov, implying that $\boldsymbol{\theta}(\boldsymbol{\theta}_{+1} : \boldsymbol{\theta}_0) = \boldsymbol{\theta}(\boldsymbol{\theta}_{+1})$, this is not necessary.

To specify $\boldsymbol{\theta}(\boldsymbol{\theta}_{+1} : \boldsymbol{\theta}_0)$ in our stationary Gaussian case, we compute the best one-step-ahead linear prediction of $\boldsymbol{\theta}_{+1}$ from $\boldsymbol{\theta}_0$ given by $\boldsymbol{\theta}_{+1} = \boldsymbol{\theta}_0 + \boldsymbol{\theta}_0 + \boldsymbol{\theta}_0$. The coefficients

, σ_0^2 , can be computed recursively in \mathcal{D} from the ACF of \mathcal{D} via the classical Durbin-Levinson (DL) or the Innovations algorithm, for examples. As a convention, we take $\sigma_0^2 = 0$. Let $\sigma^2 = \mathbb{E}[(\epsilon_t)^2]$ be the corresponding unconditional mean squared prediction error. With this notation,

$$(\epsilon_{t+1} | \mathcal{D}_t) \stackrel{\mathcal{D}}{=} \mathcal{N}(\epsilon_{t+1} | \sigma_{t+1}^2) \quad (11)$$

where $\sigma_{t+1}^2 = \sigma_t^2 + \sigma^2$. Again, \mathcal{D} does not have to be Markovian (of any order). On the other hand, with (1),

$$(\epsilon_t | \mathcal{D}_t) = (\epsilon_t | \mathcal{D}_{t-1}) = \begin{cases} 1 & \text{if } \epsilon_t = (\epsilon_t) \\ 0 & \text{otherwise} \end{cases} \quad (12)$$

where (ϵ_t) is a unit point mass at ϵ_t . The equations in (11) and (12) constitute the SSM representation of (1).

In inference and related tasks for SSMs, the basic goal is to compute the conditional expectation $\mathbb{E}[(\epsilon_{0:t} | \mathcal{D}_{0:t})]$ for some function \cdot . This is often carried out through an importance sampling algorithm such as sequential importance sampling (SIS), which generates N independent particle trajectories $\mathcal{D}_{0:t}^i$, $i = 1, \dots, N$, from a proposal distribution $(\epsilon_{0:t} | \mathcal{D}_{0:t}^i)$ and approximates the conditional expectation as

$$\mathbb{E}[(\epsilon_{0:t} | \mathcal{D}_{0:t})] \approx \sum_{i=1}^N (\epsilon_{0:t} | \mathcal{D}_{0:t}^i) \tilde{w}_i =: \mathbb{E}[(\epsilon_{0:t} | \mathcal{D}_{0:t})] \quad (13)$$

where

$$\tilde{w}_i = \frac{(\epsilon_{0:t} | \mathcal{D}_{0:t}^i)}{\sum_{i=1}^N (\epsilon_{0:t} | \mathcal{D}_{0:t}^i)} \quad (\epsilon_{0:t} | \mathcal{D}_{0:t}^i) = \frac{(\epsilon_{0:t} | \mathcal{D}_{0:t}^i)}{(\epsilon_{0:t} | \mathcal{D}_{0:t}^i)} \quad (14)$$

are the (normalized) importance weights (see [14] and [32]). Furthermore, in SIS,

$$\tilde{w}_i \sim \tilde{w}_{i-1} (\epsilon_{0:t} | \mathcal{D}_{0:t}^i) \quad (\epsilon_{0:t} | \mathcal{D}_{0:t}^i) = \frac{(\epsilon_{0:t-1} | \mathcal{D}_{0:t-1}^i) (\epsilon_t | \mathcal{D}_{0:t-1}^i)}{(\epsilon_{0:t-1} | \mathcal{D}_{0:t-1}^i)} \quad (15)$$

(see (1.6) in [14], which is adapted to a possibly non-Markov setting by replacing $(\cdot \mid -1)$ with $(\cdot \mid 0: -1)$). The two probability terms in the numerator of $(\cdot \mid 0:)$ in (15) constitute the SSM, whereas the denominator relates to the proposal distribution.

We suggest the following proposal distribution and the resulting SIS algorithm for our model. Take

$$(\cdot \mid 0: -1 \quad 0:) \stackrel{\mathcal{D}}{=} \mathcal{N}_{x_t}(\cdot \mid \cdot^2) \quad (16)$$

where \mathcal{N} denotes a normal distribution restricted to the set \cdot , and

$$= \cdot : \quad -1(\cdot \mid -1) \quad -1(\cdot \mid \cdot) \quad (17)$$

The role of \cdot stems from the fact

$$= (\cdot) \quad (18)$$

(i.e., the count value \cdot is obtained if and only if \cdot ; see the expression (A.2) for (\cdot)).

In particular, for \cdot generated from the proposal distribution (16), the term (\cdot) in the incremental weight $(\cdot \mid 0:)$ of (15) is always set to unity. The rest of the incremental weights are calculated as

$$\begin{aligned} (\cdot \mid 0:) &= \frac{(\cdot \mid 0: -1)}{(\cdot \mid 0: -1 \quad 0:)} = \frac{\exp\left(-\frac{(z_t - \hat{z}_t)^2}{2r_t^2}\right) (2 \cdot^2)^{1/2}}{\exp\left(-\frac{(z_t - \hat{z}_t)^2}{2r_t^2}\right) [(2 \cdot^2)^{1/2} \mathbb{P}(\cdot \mid (\cdot^2) \quad \cdot_t)]} \\ &= \mathbb{P}(\mathcal{N}(\cdot^2) \quad \cdot_t) = \left(\frac{-1(\cdot \mid \cdot_t)}{\cdot} \right) \left(\frac{-1(\cdot \mid \cdot_{t-1})}{\cdot} \right) =: (\cdot) \quad (19) \end{aligned}$$

The choice of the proposal distribution is largely motivated by $\mathbb{P}(\cdot = \cdot) = 1 \cdot_k(\cdot)$ and the explicit form in (19) for the incremental weights $(\cdot \mid 0:)$. Optimality considerations are mentioned in Remark B.3.

The following steps summarize our SIS algorithm.

Sequential Importance Sampling (SIS): For $t = 1, \dots, T$, where N represents the number of particles, initialize the weight $w_0 = 1$ and the latent series x_0 by

$$x_0 \stackrel{\mathcal{D}}{=} \mathcal{N}_{x_0}(0, 1) \quad (20)$$

Then, recursively over $t = 1, \dots, T$, perform the following steps:

1: Compute \hat{y}_t with the DL or other algorithm using the previously generated values of x_{t-1} .

2: Update the series x_t and the importance weight w_t via

$$x_t \stackrel{\mathcal{D}}{=} \mathcal{N}_{x_t}(x_{t-1}, \Sigma_t) \quad w_t = w_{t-1} \frac{p(y_t | x_t)}{p(y_t | x_{t-1})} \quad (21)$$

where Σ_t is defined in (19).

Remark 2.3. For $t = 1, \dots, T$, the constructed path $x_{t=0}$ is one of the N independent particles used to approximate the conditional expectation in (13). Equation (21) ensures that for each t , the path $x_{t=0}$ obeys the restriction $\Sigma_t(x_t) = \Sigma_t$ and matches the temporal structure of y_t . These two properties show that $x_{t=0}$ is a realization of the latent Gaussian stationary series producing $y_t = \hat{y}_t$ for all t . Finally, we note where the model parameters enter into the SIS algorithm. The marginal distribution parameters θ enter through the form of Σ_t in (19), whereas the temporal dependence parameters η enter through the one-step-ahead prediction coefficients α_{t-1} , in the calculation of \hat{y}_t in Step 1 of the algorithm, and through the prediction error ϵ_t .

To compute the model likelihood, several known formulas applicable in the (general) SIS setting are needed. The relation

$$\frac{p(y_{1:T} | x_0)}{p(y_{1:T} | x_0)} = \prod_{t=0}^{T-1} p(y_{t+1} | x_t)$$

produces

$$\mathbb{E}[\mathcal{L}(z_{0:t}) | z_0] = \mathbb{E}[\mathcal{L}(z_{0:t}) | z_0]$$

In particular (with $z_0 = 1$),

$$\mathbb{E}[\mathcal{L}(z_{0:t})] = \mathcal{L}(z_0) \quad (22)$$

To conduct prediction, we use Equation (1.2) in [14] to get

$$\mathbb{E}[\mathcal{L}(z_{0:t+1})] = \mathbb{E}\left[\mathbb{E}[\mathcal{L}(z_{0:t+1}) | z_{0:t}]\right] =: \mathbb{E}[\mathcal{L}_{t+1}(z_{0:t})] \quad (23)$$

where

$$\mathcal{L}_{t+1}(z) = \mathbb{E}[\mathcal{L}(z_{0:t+1}) | z_{0:t} = z] = \int_{\mathbb{R}} \mathcal{L}(z_{0:t+1}) \frac{1}{\sqrt{2\sigma_{t+1}^2}} \exp\left(-\frac{(z_{t+1}-z)^2}{2\sigma_{t+1}^2}\right) dz_{t+1} \quad (24)$$

since $z_{t+1} | z_{0:t} \stackrel{\mathcal{D}}{=} \mathcal{N}(z_{t+1} | z, \sigma_{t+1}^2)$. In view of (23) and (13), the following prediction approximation arises:

$$\mathbb{E}[\mathcal{L}(z_{0:t+1})] \approx \sum_{i=1}^M \omega_i \mathcal{L}_{t+1}(z_{0:t}^{(i)}) =: \mathbb{E}[\mathcal{L}(z_{0:t+1})] \approx \sum_{i=1}^M \omega_i \mathcal{L}(z_{0:t+1}^{(i)}) \quad (25)$$

Appendix B further connects our model and algorithm to the popular GHK sampler, hidden Markov models (HMMs), and PF and SMC techniques.

The SIS algorithm has a fundamental weakness called *weight degeneracy*: as the algorithm propagates through an increasing number of iterations, a large number of the normalized weights become negligible. As a result, only a few particles contribute in the likelihood approximation. Following the developments in the SMC (see [14], [31] and [8]) and HMM literatures (Sections 10.4.1 and 10.4.2 in [13]), we modify the SIS algorithm by adding a resampling step (all future simulations and computations use resampling).

Sequential Importance Sampling with Resampling (SISR): Proceed as in the SIS algorithm, but modify Step 2 and add a resampling Step 3 as follows:

2: Modify Step 2 of the SIS by setting

$$\tilde{x}_t \stackrel{\mathcal{D}}{=} \mathcal{N}(x_t, \Sigma_t) \quad \tilde{x}_t = x_{t-1} + \tilde{\epsilon}_t \quad \tilde{\epsilon}_t = \sum_{i=1}^M \tilde{\epsilon}_t^i \quad (26)$$

3: For each particle $i = 1, \dots, M$, draw, conditionally and independently given $(x_{t-1}, \tilde{\epsilon}_t^i)$, a multinomial trial (n_i, p_i) for each i and with the success probabilities $p_i = \tilde{w}_t^i$ and $\sum_{i=1}^M p_i = 1$.

While the resampling step removes particles with low weights, mitigating degeneracy issues, it introduces additional estimator variance. We follow standard practice and resample only when the variance of the weights exceeds a certain threshold, quantified by the so-called *effective sample size* defined as $\text{ESS}(w_t) = (\sum_{i=1}^M w_t^i)^{-1}$, and the resampling step is executed when $\text{ESS}(w_t) < 2$ as in [15]. See also Section 2.5.3 in [31] for a justification of the ESS based on the Delta method.

3 Inference

The model in (1) contains the parameters θ in the marginal count distribution and η in the dependence structure of x_t . This section addresses inference questions, including parameter estimation and goodness-of-fit assessment. Three methods are presented for parameter estimation: Gaussian pseudo-likelihood, implied Yule-Walker moment methods, and full likelihood. Gaussian pseudo-likelihood estimators, a time series staple, pretend that the series is Gaussian and maximize its Gaussian-based likelihood. These estimators only involve the mean and covariance structure of the series, are easy to compute, and will provide a comparative basis for likelihood estimators. They can also be used as initial

guesses in gradient step-and-search likelihood optimizations. Implied Yule-Walker techniques are moment based estimators applicable to the commonly encountered case where \mathbf{X} is a causal autoregression. Likelihood estimators, the statistical gold standard and the generally preferred estimation technique, are based on the PF and SMC methods of the last section. Finally, we will not delve into a detailed statistical inference for the aforementioned methods: while consistency and asymptotic normality are expected in some of the examined cases (e.g. likelihood estimation with an autoregressive \mathbf{X}), a rigorous theoretical treatment is beyond the scope of this paper.

3.1 Gaussian pseudo-likelihood estimation

As in Section 2.3, we work with observations \mathbf{x}_t for the times $t = 0, \dots, T$ and set $\mathbf{X} = (\mathbf{x}_0, \dots, \mathbf{x}_T)'$. Denote the likelihood of the model in (1) by

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}) = \mathbb{P}(\mathbf{x}_0 = \mathbf{x}_0, \mathbf{x}_1 = \mathbf{x}_1, \dots, \mathbf{x}_T = \mathbf{x}_T) = \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}; \mathbf{x}_0, \dots, \mathbf{x}_T) \quad (27)$$

While this likelihood is a multivariate normal probability, it is difficult to calculate or approximate when T is large. For most count model classes, true likelihood estimation is difficult to conduct as joint distributions are generally intractable [11]. While Section 3.3 below devises a well performing PF/SMC likelihood approximation (see also [39]), we first consider a simple Gaussian pseudo-likelihood (GL) approach. In a pseudo GL approach, parameters are estimated via

$$\hat{(\boldsymbol{\theta}, \boldsymbol{\eta})} = \underset{\boldsymbol{\theta}, \boldsymbol{\eta}}{\operatorname{argmax}} \frac{-\frac{1}{2}(\mathbf{X} - \boldsymbol{\mu}_{\boldsymbol{\theta}})' \boldsymbol{\Sigma}^{-1}(\boldsymbol{\theta}, \boldsymbol{\eta})(\mathbf{X} - \boldsymbol{\mu}_{\boldsymbol{\theta}})}{(2\pi)^{(T+1) \cdot d} |\boldsymbol{\Sigma}(\boldsymbol{\theta}, \boldsymbol{\eta})|^{1/2}} \quad (28)$$

where $\boldsymbol{\mu}_{\boldsymbol{\theta}} = (\boldsymbol{\theta}, \boldsymbol{\theta})'$ is a $(T+1) \cdot d$ -dimensional constant mean vector. These estimators maximize the series likelihood assuming the data are Gaussian, each component having

mean θ , and all components having covariance matrix $(\theta \ \eta) = (\ \)_{=0}$. Time series analysts have been maximizing Gaussian pseudo likelihoods for decades, regardless of the series marginal distribution, with often satisfactory performance. The next section and Appendix C present a case where this approach works reasonably well, and one where it does not. For large n , the pseudo GL approach is equivalent to least squares estimation, where the sum of squares $\sum_{=0}(\ \ \mathbb{E}[\ 0 \ -1])^2$ is minimized (see Chapter 8 in [6]). The covariance structure of $\ \$ was efficiently computed in Section 2; the mean θ is usually explicitly obtained from the marginal distribution $\ \$ posited. Numerical optimization of (28) yields a Hessian matrix that can be inverted to obtain standard errors for the model parameters. These standard errors can be asymptotically corrected for distributional misspecification via the sandwich methods of [19].

3.2 Implied Yule-Walker estimation for latent AR models

Suppose that $\ \$ follows the causal AR($\ \$) model $\ = \ \ _1 \ -1 + \ \ + \ \ - + \ \$, where $\ \$ consists of IID $\mathcal{N}(0 \ \ ^2)$ variables. Here, $\ \ ^2$ depends on the autoregressive coefficients $\ \ _1$ in a way that induces $\mathbb{E}[\ ^2] = 1$. The Yule-Walker equations are

$$\phi = \Gamma^{-1}\gamma \tag{29}$$

where $\Gamma = (\ \ \)_{=1}$, $\gamma = (\ (1) \ \ \)'$, and $\phi = (\ \ _1 \ \ \)'$. From (6), note that

$$\ \ \ = \ \ ^{-1}(\ \ \) \tag{30}$$

the inverse being justified via the strictly increasing nature of $\ \ \$ in $\ \ \$.

Equations (29) and (30) suggest the following estimation procedure. First, estimate the CDF parameter θ directly from the counts; standard methods (e.g. method of moments) are

typically available for this task. The estimated parameter $\boldsymbol{\theta}$ defines an estimated link $\ell(\cdot)$ through its estimated power series coefficients. From a numerical power series reversion procedure, one can now efficiently construct the inverse estimator $\ell^{-1}(\cdot)$.

Next, in view of (30) and (29), set

$$\ell(\cdot) = \ell^{-1}(\ell(\cdot)) \quad \boldsymbol{\phi} = \boldsymbol{\Gamma}^{-1} \boldsymbol{\gamma} \quad (31)$$

where $\ell(\cdot)$ is the lag- s sample autocorrelation of x_s , and $\boldsymbol{\Gamma}$ and $\boldsymbol{\gamma}$ are defined analogously to the above using $\ell(\cdot)$ in place of $\ell(\cdot)$.

3.3 Particle filtering and sequential Monte Carlo likelihoods

Using (23) and its notation, the true likelihood in (27) is

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}) = \ell(0) \prod_{s=1}^{\infty} \ell(0: -1) = \ell(0) \prod_{s=1}^{\infty} \mathbb{E}[1_{\{x_s\}}(\cdot) | 0: -1] = \ell(0) \prod_{s=1}^{\infty} \mathbb{E}[\ell(\cdot) | 0: -1] \quad (32)$$

where (23) was used with $1_{\{x_s\}}(\cdot) = \ell(\cdot)$ and $\ell(\cdot)$ is defined and numerically computed akin to (19). The particle approximation of the likelihood is then

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}) = \ell(0) \prod_{s=1}^{\infty} \mathbb{E}[\ell(\cdot) | 0: -1]; \quad (33)$$

this uses the notation in (13) and supposes that the particles are generated by one of the methods in Section 2.3. The approximate PF maximum likelihood estimates satisfy

$$(\boldsymbol{\theta}, \boldsymbol{\eta}) = \underset{\boldsymbol{\theta}, \boldsymbol{\eta}}{\operatorname{argmax}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}) \quad (34)$$

Remark 3.1. With the SIS algorithm, (33) reduces to

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}) = \ell(0) \frac{1}{\sum_{s=1}^{\infty} 1} \sum_{s=1}^{\infty} 1 \quad (35)$$

which is consistent with (22). The work [35] also essentially implements (35). In contrast to [35], our approach includes a resampling step in the likelihood approximations, considers other estimation approaches (pseudo GL and implied Yule-Walker), and provides model diagnostic tools more specific to count series (the PIT histograms in Section 3.4 below).

To optimize the estimate $\mathcal{L}(\boldsymbol{\theta} | \boldsymbol{\eta})$, we employ a large number of particles (growing linearly with n) and common random number (CRN) techniques, a standard practice that serves to smooth $\mathcal{L}(\boldsymbol{\theta} | \boldsymbol{\eta})$ somewhat by expressing its random quantities through parameter-dependent transformations of uniform random variables that remain constant for likelihood evaluations across distinct parameters. While the CRN procedure works well in SIS, it fails to ward against discontinuous $\mathcal{L}(\boldsymbol{\theta} | \boldsymbol{\eta})$ in our preferred SISR algorithm. An elegant solution to this issue for univariate state processes is proposed in [34]: first reorder the (real-valued) particles and then replace the discontinuous resampling CDF with a piecewise linear approximation. More recent and well performing (but less straightforward) approaches such as the sequential quasi Monte Carlo and the SMC² algorithm are reviewed in detail in Chapters 13, 14, and 18 of [8] (see also the Chapter 19 references on controlled sequential Monte Carlo methods). We do not pursue these issues further here.

In our numerical implementations, gradient-free algorithms from the R package `optimx` [36] are used, which follows standard practices in optimizing noisy objective functions. These routines allow for boundary constraints and performed well in modest computing times for our sample sizes. On the other hand, we found less success with the more popular gradient-based quasi-Newton algorithm `L-BFGS-B` (gradients were computed via finite differences) as convergence instabilities and high-variance estimates were encountered. However, promising recent developments for optimizing noisy objectives in [4] and [38] were not explored. A comprehensive investigation of these approaches and of the rich gradient-

based SMC inference literature for our framework as in [27] is deferred to future work.

3.4 Model diagnostics

The goodness-of-fit of count models is commonly assessed through probability integral transform (PIT) histograms and related tools [10, 29]. These are based on the predictive distributions of \mathbf{y}_t , defined at time t by

$$P_t(\mathbf{y}_t) = \mathbb{P}(\mathbf{y}_0 = \mathbf{y}_0 \mid \mathbf{y}_{-1} = \mathbf{y}_{-1}) = \mathbb{P}(\mathbf{y}_0 \in [0, 1]) \quad (36)$$

This quantity can be estimated through the PF/SMC methods in Section 2.3 as

$$P_t(\mathbf{y}_t) = \sum_{s=0}^t \mathbb{E}[1_{\{\mathbf{y}_s\}}(\mathbf{y}_t \mid \mathbf{y}_{0:-1})] = \sum_{s=0}^t \mathbb{E}[1_{\{\ell_s\}}(\mathbf{y}_t \mid \mathbf{y}_{0:-1})] \quad (37)$$

which uses (24) and (25) and supposes that the particles are generated by the SIS, SISR, or other algorithms. Similar to $1_{\{x_s\}}(\mathbf{y}_t) = 1_{\{x\}}(\mathbf{y}_t)$, note that $1_{\{\ell_s\}}(\mathbf{y}_t) = \tilde{1}_{\{\ell\}}(\mathbf{y}_t)$, where

$$\tilde{1}_{\{\ell\}}(\mathbf{y}_t) = \left(\frac{-1(\mathbf{y}_t)}{+1} \right) \quad \left(\frac{-1(\mathbf{y}_{-1})}{+1} \right) \quad (38)$$

and $\tilde{1}_{\{x\}}(\mathbf{y}_t) = 1_{\{x\}}(\mathbf{y}_t)$.

The (non-randomized) sample mean PIT is defined as

$$\bar{P}_t(\mathbf{y}_t) = \frac{1}{t+1} \sum_{s=0}^t P_s(\mathbf{y}_t) \quad [0, 1] \quad (39)$$

where

$$P_s(\mathbf{y}_t) = \begin{cases} 0 & \text{if } \mathbf{y}_s \in [0, 1) \\ \frac{-1(\mathbf{y}_s) - 1}{-1(\mathbf{y}_s) - 1} & \text{if } \mathbf{y}_s \in [1, 1) \\ 1 & \text{if } \mathbf{y}_s \in [1, 1) \end{cases} \quad (40)$$

which is estimated by replacing $\hat{\mu}_t$ by $\hat{\mu}_{t-1}$ in practice. The PIT histogram with K bins is defined as a histogram with the height $\frac{1}{K} \mathbb{1}(\hat{F}_t(x) \in [\frac{k-1}{K}, \frac{k}{K}])$ for bin $k = 1, \dots, K$.

Another possibility considers model residuals based on

$$\hat{\mu}_t = \mathbb{E}[\mu_t | \mathcal{Y}_{t-1}] = \frac{\exp(-\frac{1}{2}(\hat{\mu}_{t-1} - \mu_t)^2 / \sigma^2)}{\int \exp(-\frac{1}{2}(\hat{\mu}_{t-1} - \mu)^2 / \sigma^2) d\mu} \quad (41)$$

which is the estimated mean of the latent Gaussian process at time t given \mathcal{Y}_{t-1} only (not the entire past), where (41) follows by direct calculations for the model (1) (using the estimated parameters $\hat{\theta}$ of the marginal distribution in the \mathcal{Y}_{t-1} s). For a fitted underlying time series model with parameter η , the residuals are then defined as the standard time series residuals $\hat{\epsilon}_t$ of this model fitted to the series \hat{y}_t , after centering by the sample mean.

3.5 Nonstationarity and covariates

As discussed in Section 2.2, covariates can be accommodated by allowing a time-varying parameter θ in the marginal distribution. With covariates, θ at time t is denoted by $\theta(t)$. The GL and PF/SMC procedures are modified for $\theta(t)$ as follows.

For the GL procedure, the covariance $\text{Cov}(\mu_{(1)}(x_1) - \mu_{(2)}(x_2))$ is needed, where $\mu_{(i)}$ is subscripted to signify dependence on $\theta(t)$. But as in (5),

$$\text{Cov}(\mu_{(1)}(x_1) - \mu_{(2)}(x_2)) = \sum_{k=1}^{\infty} \frac{1}{k!} \theta_{(1)}^{(k)}(x_1) \theta_{(2)}^{(k)}(x_2) \quad (42)$$

where again, the subscript $\theta(t)$ is added to the μ s to indicate dependence on t . Numerically, evaluating (42) is akin to the task in (5); in particular, both calculations are based on the Hermite coefficients $\theta^{(k)}$.

For the PF/SMC approach, the modification is somewhat simpler: one just needs to replace θ by $\theta(t)$ at time t when generating the underlying particles. For example, for the

SIS algorithm, $\theta(\cdot)$ enters only through the μ s in (19), (20), and (21). This is because the covariates enter only through θ , the parameter controlling marginal distributions.

4 A simulation study

To evaluate our estimation methods, a simulation study considering several marginal distributions and dependence structures was conducted. Here, the classic Poisson count distribution \mathcal{P} is examined (mixed Poisson and negative binomial simulations are presented in Appendix C), with μ taken from the ARMA(p, q) class. All simulation cases are replicated 200 times for three distinct series lengths: $n = 100, 200$, and 400. For notation, estimates of a parameter θ from Gaussian pseudo-likelihood (GL), implied Yule-Walker (IYW), and PF/SMC methods are denoted by $\hat{\theta}_{GL}$, $\hat{\theta}_{IYW}$, and $\hat{\theta}_{PF/SMC}$, respectively.

We now consider the classical case where μ has a Poisson marginal distribution for each t with mean $\mu = 0$. To obtain μ_t , the AR(1) process $\mu_t = \rho \mu_{t-1} + (1 - \rho^2)^{1/2} \epsilon_t$, was simulated and transformed via (1) with $\mu_t = \mathcal{P}; \mathbb{E}[\mu_t^2] = 1$ was induced by taking $\text{Var}(\epsilon_t) = 1$. Twelve parameter schemes resulting from all combinations of $\rho \in \{-0.2, -0.5, -0.75, 0.2, 0.5, 0.75\}$ and $\mu \in \{0.25, 0.75\}$ were considered.

Figure 1 displays box plots of the parameter estimates when $n = 200$. In estimating μ , all methods perform reasonably well. When the lag-one correlation in μ (and hence also that in μ_t) is negative (right panel), $\hat{\mu}_{GL}$, $\hat{\mu}_{IYW}$, and $\hat{\mu}_{PF/SMC}$ have smaller variability than the positively correlated case (left panel — note the different y-axis scales on the panels). This is expected: the mean of μ is μ , and the variability of the sample mean, one good estimator of the mean for a stationary series, is smaller for negatively correlated series than for positively correlated ones. Note that $\hat{\mu}_{GL}$ is biased toward zero for both negatively and

positively correlated series, whereas $\hat{\phi}_{IYW}$ and $\hat{\phi}_{PF}$ only show bias when ϕ is positive for the sample sizes $T = 100$ and $T = 200$. Overall, the PF/SMC estimates were the least biased. All estimates of ϕ have roughly similar variances. Simulations with $\lambda = 5$ and $\lambda = 10$ produced analogous results with smaller values of λ yielding less variable estimates. This is again expected as the variance of the Poisson distribution is also λ . Graphics of these box plots are omitted for brevity's sake.

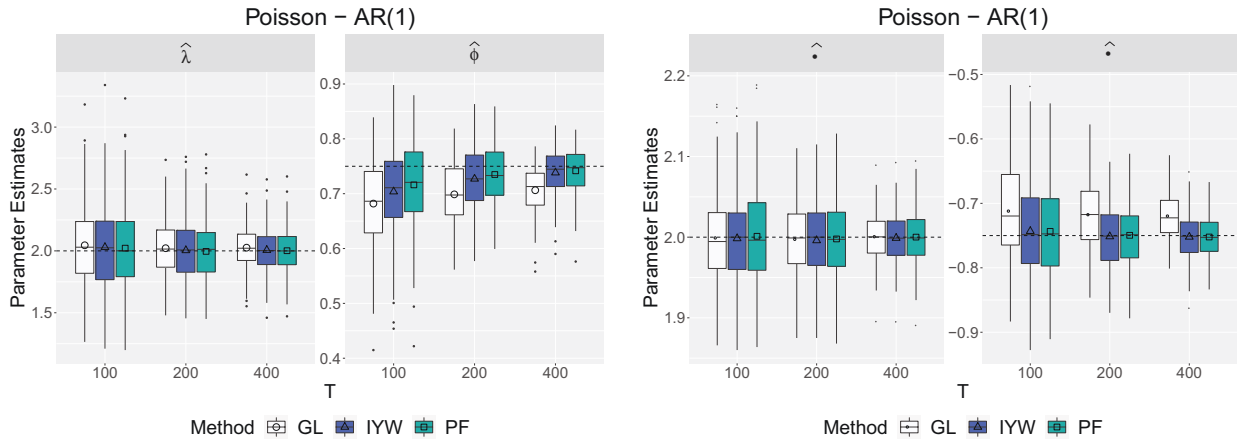


Figure 1: Gaussian likelihood, implied Yule-Walker, and PF/SMC parameter estimates for 200 synthetic Poisson-AR(1) series with lengths $T = 100, 200,$ and 400 . The true parameter values (indicated by horizontal dashed lines) are $\lambda = 2$ and $\phi = 0.75$ (left panel), and $\lambda = 2$ and $\phi = -0.75$ (right panel).

5 An application

This section applies our methods to a weekly count series of product sales at Dominicks Finer Foods, a now defunct U.S. grocery chain that operated in Chicago, IL and adjacent areas

from 1918 - 2013. Soft drink sales of an unnamed brand from a single store will be analyzed over a two-year span commencing on September 10, 1989. The series is plotted in Figure 2 (leftmost plot) and is part of a large and well-studied retail dataset, publicly available at <https://www.chicagobooth.edu/research/kilts/datasets/dominicks> (Source: The James M. Kilts Center for Marketing, University of Chicago).¹ Our goal here is not an in-depth retail

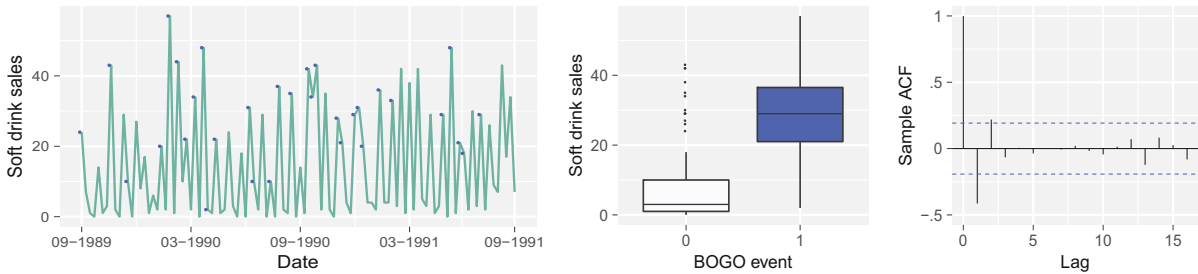


Figure 2: *Left: Weekly sales of a soft drink product sold at a single store of the grocery store Dominick’s Finer Foods from 09-10-1989 to 09-03-1991. The dots indicate the weekly sales were at least one “Buy one and get one free” (BOGO) sales promotion event took place. Middle: Boxplots of sales grouped by the BOGO covariate (0: weekly sales with no BOGO event, 1: weekly sales with at least one BOGO day during the week). Right: Sample ACF of the series with 95% pointwise bands for zero correlation.*

analysis, but to illustrate our methods with a real world example of an overdispersed time series of small counts that has negative autocorrelation and dependence on a covariate.

The covariate we use is a zero-one “buy one get one free” (BOGO for short) sales promotion event S_t , $S_t = 1$ implying that the BOGO promotion was offered at least one day during week t . The dots in the left plot of Figure 2 signify that the week had at

¹In the dataset manual, the series in Figure 2 (leftmost plot) is the sales of the product with universal product code (UPC) 4640055081 from store 81.

least one BOGO day. The middle plot shows the soft drinks sales distribution grouped by t , visually suggesting that a BOGO event increases soft drink sales. The rightmost plot shows the sample ACF of the series and reveals negative dependence at lag one. The lag one sample autocorrelation of the residuals after a linear regression of the series on the BOGO covariate is also negative, but comparatively smaller in magnitude.

To model overdispersion, negative binomial and generalized Poisson marginal distributions will be considered. Although similar, these two distributions can yield different conclusions [25]. Following standard generalized linear modeling practice, both distributions are parametrized via the series mean (although our setup allows covariates to enter through other parameters as well). More specifically, for the negative binomial marginal, the standard pair (μ, ϕ) used in Appendix C is now mapped to the parameter pair (μ, ϕ) , where $\mu = \mu(1 - \phi)$ is the mean of the process and $\phi = 1 - \phi$ is the overdispersion parameter. Similarly, the generalized Poisson distribution of Appendix C is parametrized through the pair (μ, ϕ) as in [17], relation (2.4). In this parametrization, μ is the mean of the series, whereas the sign of ϕ controls the type of dispersion, with positive values indicating overdispersion. To incorporate the BOGO covariate x_t into the model, the mean of the series is allowed to depend on time t through the typical GLM log-link $\mu_t = \exp(\beta_0 + \beta_1 x_t)$, while the parameters β_0 and β_1 are kept fixed in time t .

An exploratory examination of the sample ACFs and PACFs of the series along with diagnostic plots of residuals obtained by fitting all ARMA(p, q) models with $p, q \leq 5$ suggest an AR(3) model as a suitable choice for y_t . Table 1 in Appendix D shows the AICc and BIC for both marginal distributions obtained via PF/SMC and GL methods (we omit IYW results for simplicity). The AR(3) model was selected by AICc and BIC in both fits. Interestingly, both the sample ACF and PACF of the series show one large non-zero value

at lag one, but relatively smaller values at other lags (except perhaps the lag two value, which barely exceeds the 95% 1.96 $\sqrt{\text{var}}$ dashed confidence threshold for zero correlation).

We also considered a white noise latent series (labeled as WN Table 1 in Appendix D), which renders our model a standard GLM. The PF/SMC WN estimates from both distributions (omitted here for brevity) closely agree with parameter estimates obtained from exact generalized linear models (using, for example, functions from the R package MASS). As expected, the WN model yielded the highest AICc and BIC values among all considered dependence structures, thus confirming the need for a model with temporal dependence.

Table 1 shows parameter estimates and standard errors from fitting a negative binomial-AR(3). (Table 2 in Appendix D is for a generalized Poisson-AR(3) model.) All marginal distributions and estimation methods yielded $\hat{\rho}_1 < 0$. Although a formal asymptotic theory is beyond the scope of our presentation here, asymptotic normality is expected. Assuming this, the PF/SMC standard errors (the ones believed most trustworthy) suggest that all parameters are significantly non-zero at level 95%. The findings suggest the negative binomial distribution is preferred over the generalized Poisson, that the correlation in the series at lag one is negative, and that a BOGO event indeed increases sales.

	1	2	3	0	1

Table 1:

Turning to residual diagnostics, the plots in Figure 3 for the negative binomial-AR(3)

fit suggest that the model has captured both the marginal distribution and the dependence structure. The residuals here were computed using (41).

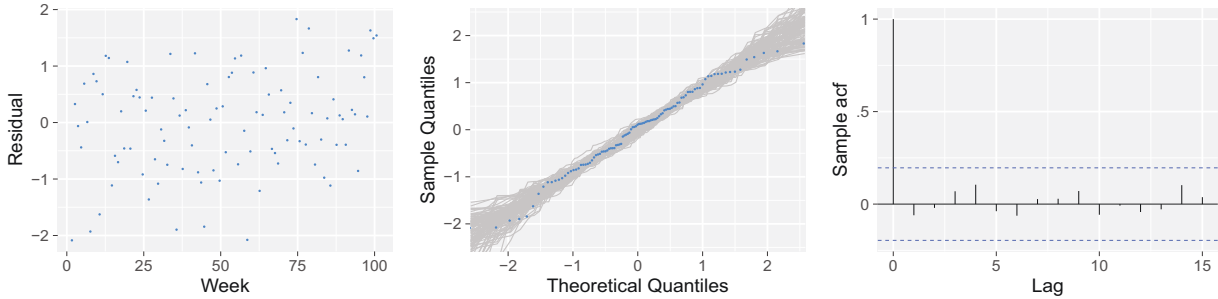


Figure 3: *The leftmost plot displays the estimated residuals against time. The middle graph is a QQ plot for normality of the estimated residuals. The shaded region in the QQ plot shows 100 realizations from a normal distribution with size, mean and standard deviation matching the residual sample counterparts. The right plot displays the sample autocorrelations of the estimated residuals.*

We next assess the predictive ability of the two fits via the non-randomized histograms shown in Figure 4 and discussed in detail in Section 3.4. We selected ten bins at the points $h/10, h = 1, \dots, 10$ as is typical in the literature. The negative binomial PIT plot suggests a satisfactory predictive ability with most bar heights being close to 0.1 (1 over the number of bins). In comparison, the generalized Poisson fit deviates more from the uniform distribution, with somewhat more pronounced peaks and valleys. We remind the reader here that PIT plots are known to be sensitive for smaller series lengths. Quantifying this uncertainty (for each bin) through a statistical test is beyond the scope of this paper. Nevertheless, we gauged the variability of the uniform distribution’s bin heights through a small experiment. Specifically, 500 synthetic realizations of sample size $T = 104$ were generated and the percentiles of all bin heights were collected. The 5th and 95th percentiles

ranged in the intervals (0.048, 0.058) and (0.145, 0.154) respectively, suggesting that the peaks and valleys of the negative binomial PIT plot (which are within these percentiles) are mild; that is, uniformity is plausible and the marginal distribution fits seems adequate.

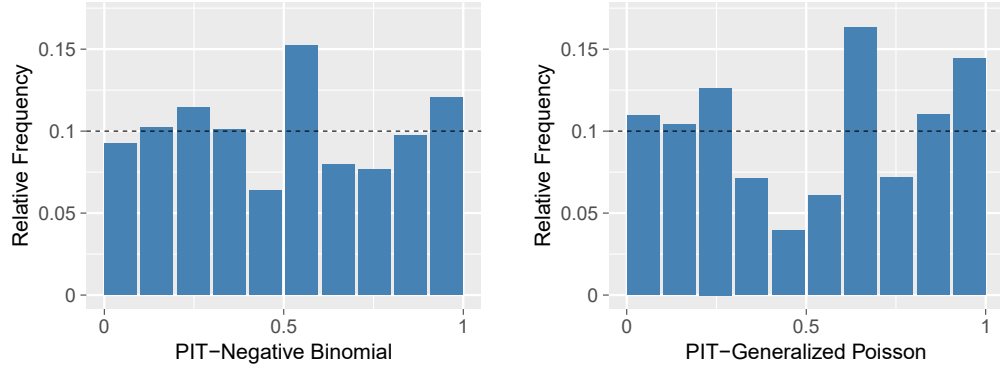


Figure 4:

6 Conclusions and comments

This paper developed the theory and methods for a stationary count time series model made from a latent Gaussian process. By using Hermite expansions, a very general model class was devised. In particular, the autocorrelations in the series can be positive or negative, and in a pairwise sense, span the range of all achievable correlations. The series can have any marginal distribution desired, thereby improving classical DARMA and INARMA count time series methods. On inferential levels, autocovariances of the model were extracted from Hermite expansions, allowing for Gaussian pseudo-likelihood and implied Yule-Walker inference procedures. A PF/SMC likelihood approach was also developed and produced estimators that were demonstrated to outperform the Gaussian

pseudo-likelihood and implied Yule-Walker estimators in most cases. These results complement the importance sampling methods for copula likelihoods in [39]. The methods were used in a simulation study and were applied in a regression analysis of a count series of weekly grocery sales that exhibited overdispersion, a negative lag one correlation, and dependence on a buy one get one free covariate. Model fits and predictive abilities of the methods were illustrated with generalized Poisson and negative binomial marginal distributions.

While the paper provides a reasonably complete treatment for count time series models, additional research is needed. Some statistical issues, like asymptotic normality of parameter estimators, were not addressed here. PF/SMC algorithms that optimize model likelihoods, which can be unwieldy, also merit further exploration. The paper only considers univariate methods. Multivariate count time series models akin to those in [40] could be developed by replacing the univariate \mathbf{Z}_t with a multivariate Gaussian process \mathbf{Z} , whose components have a standard normal marginal distribution, but are cross-correlated for each t . The details for such a construction would proceed akin to the methods developed here. Also, while the count case is considered here, the same methods will produce stationary time series having any general prescribed continuous distribution. Finally, the same methods should prove useful in constructing spatial and spatio-temporal processes having any prescribed marginal distribution. While [12, 21] recently addressed this issue in the spatial setting, additional work is needed, including exploring spatial Markov properties and likelihood evaluation techniques. To the best of our knowledge, no comprehensive analogous work has been conducted for space-time count modeling to date.

References

- [1] Asmussen, S. (2014). Modeling and performance of bonus-malus systems: stationarity versus age-correction. *Risks* 2, 49–73.
- [2] Belyaev, M., E. Burnaev, and Y. Kapushev (2015). Gaussian process regression for structured data sets. In A. Gammerman, V. Vovk, and H. Papadopoulos (Eds.), *Statistical Learning and Data Sciences: Third International Symposium, SLDS 2015*. Switzerland: Springer International Publishing.
- [3] Benjamin, M. A., R. A. Rigby, and D. M. Stasinopoulos (2003). Generalized autoregressive moving average models. *Journal of the American Statistical Association* 98, 214–223.
- [4] Berahas, A. S., R. H. Byrd, and J. Nocedal (2019). Derivative-free optimization of noisy functions via quasi-Newton methods. *SIAM Journal on Optimization* 29, 965–993.
- [5] Blight, P. A. (1989). Time series formed from the superposition of discrete renewal processes. *Journal of Applied Probability* 26, 189–195.
- [6] Brockwell, P. J. and R. A. Davis (1991). *Time Series: Theory and Methods* (Second ed.). New York City: Springer-Verlag.
- [7] Chen, H. (2001). Initialization of NORTA: Generation of random vectors with specified marginals and correlations. *Inform Journal on Computing* 13, 312–331.
- [8] Chopin, N. and O. Papaspiliopoulos (2020). *An Introduction to Sequential Monte Carlo Methods*. New York City: Springer.

- [9] Cui, Y. and R. B. Lund (2009). A new look at time series of counts. *Biometrika* 96, 781–792.
- [10] Czado, C., T. Gneiting, and L. Held (2009). Predictive model assessment for count data. *Biometrics* 65, 1254–1261.
- [11] Davis, R. A., S. H. Holan, R. B. Lund, and N. Ravishanker (Eds.) (2016). *Handbook of Discrete-Valued Time Series*. Boca Raton, Florida, USA: CRC Press.
- [12] De Oliveira, V. (2016). Hierarchical Poisson models for spatial count data. *Journal of Multivariate Analysis* 122, 393–408.
- [13] Douc, R., E. Moulines, and D. S. Stoer (2014). *Nonlinear Time Series: Theory, Methods, and Applications with R Examples*. Boca Raton, Florida, USA: CRC Press.
- [14] Doucet, A., N. De Freitas, and N. Gordon (2001). An introduction to sequential Monte Carlo methods. In *Sequential Monte Carlo Methods in Practice*, pp. 3–14. New York City: Springer.
- [15] Doucet, A. and A. M. Johansen (2009). A tutorial on particle filtering and smoothing: Fifteen years later. *Handbook of Nonlinear Filtering* 12, 656–704.
- [16] Dumsmuir, W. T. M. (2016). Generalized linear autoregressive moving-average models. In *Handbook of Discrete-valued Time Series*, pp. 51–76. Boca Raton, Florida, USA: CRC Press.
- [17] Famoye, F. (1993). Restricted generalized Poisson regression model. *Communications in Statistics-Theory and Methods* 22, 1335–1354.

- [18] Fokianos, K. (2012). Count time series models. In *Handbook of Statistics*, Volume 30, pp. 315–347. Amsterdam: Elsevier.
- [19] Freedman, D. (2006). On the so-called Huber sandwich estimator and robust standard errors. *The American Statistician* 60, 299–302.
- [20] Grigoriu, M. (2007). Multivariate distributions with specified marginals: applications to wind engineering. *Journal of Engineering Mechanics* 133, 174–184.
- [21] Han, Z. and V. De Oliveira (2016). On the correlation structure of Gaussian copula models for geostatistical count data. *Australian & New Zealand Journal of Statistics* 58, 47–69.
- [22] Han, Z. and V. De Oliveira (2020). Maximum likelihood estimation of Gaussian copula models for geostatistical count data. *Communications in Statistics - Simulation and Computation* 49, 1957–1981.
- [23] Jacobs, P. A. and P. A. W. Lewis (1978a). Discrete time series generated by mixtures I: Correlational and runs properties. *Journal of the Royal Statistical Society* 40, 94–105.
- [24] Joe, H. (1996). Time series models with univariate margins in the convolution-closed infinitely divisible class. *Journal of Applied Probability* 33, 664–677.
- [25] Joe, H. and R. Zhu (2005). Generalized Poisson distribution: the property of mixture of Poisson and comparison with negative binomial distribution. *Biometrical Journal: Journal of Mathematical Methods in Biosciences* 47, 219–229.
- [26] Kachour, M. and J. F. Yao (2009). First order rounded integer valued autoregressive (RINAR(1)) processes. *Journal of Time Series Analysis* 30, 417–448.

- [27] Kantas, N., A. Doucet, S. S. Singh, J. Maciejowski, and N. Chopin (2015). On particle methods for parameter estimation in state-space models. *Statistical Science* 30, 328–351.
- [28] Kedem, B. (1980). Estimation of the parameters in stationary autoregressive processes after hard limiting. *Journal of the American Statistical Association* 75, 146–153.
- [29] Kolassa, S. (2016). Evaluating predictive count data distributions in retail sales forecasting. *International Journal of Forecasting* 32, 788–803.
- [30] Lennon, H. (2016). *Gaussian copula modelling for integer-valued time series*. Ph. D. thesis, The University of Manchester.
- [31] Liu, J. S. (2008). *Monte Carlo Strategies in Scientific Computing*. New York City: Springer Science & Business Media.
- [32] Liu, J. S. and R. Chen (1998). Sequential Monte Carlo methods for dynamic systems. *Journal of the American Statistical Association* 93, 1032–1044.
- [33] Livsey, J., R. B. Lund, S. Kechagias, and V. Pipiras (2018). Multivariate integer-valued time series with flexible autocovariances and their application to major hurricane counts. *Annals of Applied Statistics* 12, 408–431.
- [34] Malik, S. and M. K. Pitt (2011). Particle filters for continuous likelihood evaluation and maximisation. *Journal of Econometrics* 165, 190–209.
- [35] Masarotto, G. and C. Varin (2012). Gaussian copula marginal regression. *Electronic Journal of Statistics* 6, 1517–1549.
- [36] Nash, J. C. and R. Varadhan (2011). Unifying optimization algorithms to aid software system users: optimx for R. *Journal of Statistical Software* 43, 1–14.

- [37] Pipiras, V. and M. S. Taqqu (2017). *Long-Range Dependence and Self-Similarity*, Volume 45. Boca Raton, Florida, USA: Cambridge University Press.
- [38] Shi, H.-J. M., M. Q. Xuan, F. Oztoprak, and J. Nocedal (2021). On the numerical performance of derivative-free optimization methods based on finite-difference approximations. *arXiv preprint arXiv:2102.09762*.
- [39] Smith, M. S. and M. A. Khaled (2012). Estimation of copula models with discrete margins via Bayesian data augmentation. *Journal of the American Statistical Association* 107, 290–303.
- [40] Song, P., M. Li, and P. Zhang (2013). Vector generalized linear models: a Gaussian copula approach. In P. Jaworski, F. Durante, and W. Hardle (Eds.), *Copulae in Mathematical and Quantitative Finance*. Heidelberg, Germany: Springer.
- [41] Tong, Y. L. (1990). *The Multivariate Normal Distribution*. New York City: Springer-Verlag.
- [42] Whitt, W. (1976). Bivariate distributions with given marginals. *The Annals of Statistics* 4, 1280–1289.
- [43] Zheng, T., H. Xiao, and R. Chen (2015). Generalized ARMA models with martingale difference errors. *Journal of Econometrics* 189, 492–506.