# PERIODIC GAMMA AUTOREGRESSIVE MODEL: AN APPLICATION TO THE BRAZILIAN HYDROELECTRIC SYSTEM

Diogo Braga[1,2] and Wilson Calmon[2,*]

**Abstract.** Hydrological time series forecasting play a crucial role in the Brazilian Power System since most of the power generated comes from hydroelectric power plants. A minor improvement in the predictive ability of water inflows time series might lead both to: (i) lower costs to final consumers; and (ii) higher power reliability. This paper explores Periodic Gamma Autoregressive Models (PGAR) to model brazilian water inflows time series. This type of time series has some features which seem to be more adaptable to Gamma models, like nonnegative random values and asymmetric pattern. The main purpose of this study is to compare periodic Normal and Lognormal models to PGAR. The results suggest that: (i) both Gamma and Lognormal models perform better than Normal model; and (ii) the Gamma model is a good alternative to the Lognormal model.

## 1. Introduction

The Brazilian power generation system is known worldwide as highly dependent on hydroelectric generation. In recent years, more than 75% of total electricity consumed came from hydropower plants. Despite the reduction in overall power generation costs, the system reliance on hydroelectric generation demands high accuracy on forecasting water inflows to hydropower reservoirs. An effective estimative of water flows should lead to a better management of water resources, avoiding the use of more expensive sources of energy (such as fossil-fuel power stations) and assuring a satisfactory power supply.

The role of water reservoirs in the brazilian power generation park is carefully described in [4]. They are important to system reliability since they allow the interchange of water inflows during wet and dry seasons. In other words, reservoirs allow a better control of water flows, increasing the reliability of energy supply. The nacional interconnected system (SIN) is basically divided into four subsystems: North, South, Southeast/Midwest and Northeast. They are connected by transmission lines in a way that makes possible to move energy through regions, reducing significantly the risk of power outages.

There are two groups of hydropower plants currently in operation in Brazil. The first and most important are those plants which admit large reservoirs, since it makes unpredictable climate changes more manageable and

[1] Departement of system Engineering and Computer Science, Federal University of Rio de Janeiro, Brazil. diogobmb@id.uff.br

[2] Department of Economics, Federal University Fluminense, Niterói, RT, Brazil. calmonwilson@gmail.com

[*] Wilson works at Department of Statistics for Federal University Fluminense.

operational power planning doable. The other one are the run-of-the-river power plants, which in some cases do not have any water storage. They do have some advantages when compared to conventional hydropower plants (low social and environmental impact, for example), but their contribution to the system coordination is very marginal.

The natural affluent energy (ENA) is defined as the total amount of energy produced by hydroelectric power plants. The operator of the national electricity system (ONS) is the brazilian agency that manages the water use at the SIN. Particularly, ONS is responsible for coordinating and controlling every aspect of power generation and transmission, ensuring the balance between supply and demand whilst minimizing the total cost.

As pointed out by [8,13], the operation of reservoirs has traditionally been multiannual, developed to manage a long period of drought, up to five years. The basic idea that guides water management in a large hydrothermal system is the tradeoff between an arbitrary use of water in the short term and a huge increase in operational costs in the long term.

Essentially, the seasonal pattern of hydroelectric inflows should be taken into account in order to capture their variability from one year to another or even at the same year. Another aspect that brings more complexity to brazilian operational system is the diversity of regions and their respective weather and hydrological conditions. Not necessarily the dry season will be the same at every power station. Actually, the opposite might be considered true.

The operational power planning relies heavily on stochastic optimization techniques and time series analysis. In essence, the latter feeds the former, in such a way that a minor forecast mistake could lead to wrong construction of scenarios, causing a large waste of money and a serious damage on system reliability. Thus a high quality modelling and forecasting of water inflows is critical to the whole system.

Geophysical data, by their nature, are often treated as seasonally stationary. For any given season[3], river inflows, precipitation or the temperature are very similar from one year to another, but they could vary widely between seasons (see [15]). Take, for instance, the temperature for July at southern hemisphere. It seems very reasonable that the temperature across the years should variate around the same mean for this month. In addition, the temperature in July should be not only related to preceding months but also to previous year. Periodic models aim to capture such cyclical variations. As emphasized by [15,16,19,22,25], hydrological time series show autocorrelation structures that vary over the seasons, which make periodic models a reasonable choice to analyse them. Vieira *et al.* [33], for example, applies periodic multivariate models to Brazilian hydrological time series. Essentially, periodic models generalizes their non-periodic versions, allowing us to capture directly those season patterns aforementioned. Furthermore, the existence of periodic dependency and their different forms can be statistically tested.

According to Fernandez and Salas [15], although Gaussian linear models have been useful tools for modelling hydrological time series, they may present important limitations. Among other things, Gaussian models cannot properly deal with asymmetry, which is very common in geophysical data. Also, despite their inherent ability to recognize the "shape" of water inflows time series, they might generate negative entries in forecasting experiments or scenario generation, which it is not acceptable in this case. Furthermore, a typical property of hydrological time series is that they are not time-reversible, fact that contrasts with Gaussian processes.

This paper explores the adoption of periodic gamma autoregressive model (PGAR) in the sense of conditional distributions for modeling ENA time series and compares its performance to traditional Normal and Lognormal conditional distributions. Accordingly, this investigation may provide subsidies to improve the forecasting accuracy of hydrological time series and so the construction of scenarios for stochastic optimization, minimizing the total cost for final consumers. At first sight, Gamma distribution appears to be a better choice for this kind of models since it can appropriately deal with hydrological time series properties, at least theoretically.

---

[3]For the purpose of this paper, the term season will not be associated to weather patterns, as in the usual way. In contrast, season will refer to periods of time, such as months and weeks.

Data set were obtained at the ONS website[4]. This study is restricted to the analysis of ENA time series for the southeast subsystem, which is the biggest and more relevant subsystem at the SIN. The data set contain 79 observations for each month, from 1931 to 2009.

In Section 2 PAR Models are introduced for both Gaussian and Gamma distributions and experiments designs are drawn for forecasting analysis. Results are discussed in Section 3, while Section 4 concludes.

## 2. PERIODIC AUTOREGRESSIVE MODEL (PAR)

Periodic models are frequently used to study environmental and water resources data (see [16, 19]). In fact, in addition to time lag between observations, the season of the year is central to analyze the autocorrelation structure of hydrological time series. Periodic models typically exhibit weak stationarity within a given season[5], but it is not stationary over the year (the water inflow pattern differs widely between seasons, for example).

PAR models may be described as AR models for each season of the year. Suppose, for simplicity, that data are released monthly, where $m = 1, \ldots, 12$ over a period of R years, where $r = 1, \ldots, R$. Accordingly, $T = 12R$ is the number of available observations, where $t = 1, \ldots, T$. Therefore, one can represent a PAR(1) model as follows:

$$z_t - \mu_t = \phi_1^{(t)}\Big(z_{t-1} - \mu_{t-1}\Big) + a_t \tag{2.1}$$

where $\mu_t$ is the mean of a time series $z_t$, $\mu_t = \mu_{t+12}$, $\phi_1^{(t)} = \phi_1^{(t+12)}$ is the autoregressive coefficient[6] and $a_t$ is a white noise.

PAR(p) model may be written

$$z_t - \mu_t = \sum_{i=1}^{p} \phi_i^{(t)}\Big(z_{t-i} - \mu_{t-i}\Big) + a_t \tag{2.2}$$

where p is the order of PAR model and $\phi_p^t$ must be nonzero for some $t$. It is also important to observe that the innovation variance is m-periodic.

Taking lag operator $B$ one may rewrite (2.2) in a more concise way

$$\Phi^{(t)}(B)\Big(z_t - \mu_t\Big) = a_t \tag{2.3}$$

where $\Phi^{(t)}(B) = 1 - \phi_1^{(t)}B - \phi_2^{(t)}B^2 - \ldots - \phi_p^{(t)}B^p$.

For the purpose of this work two approaches are going to be used:

I) Conditional Distribution of $\{z_t\}$ is Normal

$$z_t | z_{t-1}, z_{t-2}, \ldots \sim N\left(\sum_{i=1}^{p} \phi_i^{(t)} z_{t-i}, \sigma_t^2\right). \tag{2.4}$$

---

[4] www.ons.org.br

[5]A stochastic process $\{y_t\}$ is said to be weakly stationary if

$$E(y_t) = a;$$
$$\mathrm{Cov}(y_t, y_{t+h}) = \gamma(h)$$

where $\gamma(.)$ is the covariance function of $\{y_t\}$, a is constant independent of t and h is integer. Periodic weak stationarity is defined in [16].

[6]One has, for instance, to the $25^{(th)}$ observation $\phi_1^{(1)} = \phi_1^{(13)} = \phi_1^{(25)} = \phi_1^{(25+12)}$. That is, coefficients for january are equal for every $t$, $t = 1, \ldots, T$. Note also that $\mu_t$ is 12-periodic.

One may rewrite (2.4) as follows

$$z_t|z_{t-1}, z_{t-2}, \ldots \sim \left[\sum_{i=1}^{p} \phi_i^{(t)} z_{t-i} + \sigma_t^2 N(0,1)\right]. \tag{2.5}$$

If one considers $\sigma_t^2 N(0,1)$ as the error term of the Normal conditioned model PAR(p), (2.5) should represent an additive model, very similar to (2.2).

II) Conditional Distribution of $\{z_t\}$ is Gamma

$$z_t|z_{t-1}, z_{t-2}, \ldots \sim \text{Gamma}\left(\sum_{i=1}^{p} \phi_i^{(t)} z_{t-i}, \kappa^{(t)}\right) \tag{2.6}$$

where $\sum_{i=1}^{p} \phi_i^{(t)} z_{t-i}$ is the conditional mean and $\kappa^{(t)} = \kappa^{(t+m)}$ is the periodic coefficient of variation. Then one may write

$$z_t|z_{t-1}, z_{t-2}, \ldots \sim \left[\sum_{i=1}^{p} \phi_i^{(t)} z_{t-i}\right] \text{Gamma}\left(1, \kappa^{(t)}\right). \tag{2.7}$$

In this case, if one considers $\text{Gamma}(1, \kappa^{(t)})$ as innovation, then one has a model with multiplicative error term, which implies important changes in relation to the model described in (2.2)[7].

The model construction process follows the so-called Box−Jenkins approach applied to periodic models. Two techniques can be conveniently used to identify PAR models. The first one combines the sample periodic autocorrelation function (PeACF) and sample periodic partial autocorrelation function (PePACF). The other one is the exhaustive enumeration, which examines several regressions for each season and, using a specific model selection criteria, such as AIC or BIC[8], identifies the most appropriate model. This procedure specifies a maximum value $p^*$ for each season in order to allow the determination of the best model. In this paper we work with what we labeled as "internal zeros" for exhaustive enumeration. The first thing to do is defining which is the maximum lag allowed for each AR model in each season. In this paper we make $p^* = 23$, which allows AR models for each season to have their lag structure depending on the previous 23 months. The idea of internal zeros combined with exhaustive enumeration works in the following way: if, for instance, an AR with $p = 14$ is chosen for july, then this model may have significant coefficients for lags $1, 2, 3, 12, 13$ and $14$. The remaining time lags, comprising the fourth to the eleventh, are set to zero. Accordingly, for $p = 15$, non-zero coefficients apply for time lags equal to $1, 2, 3, 4, 12, 13, 14$ and $15$. Lastly, if $p = 23$, every coefficient for the AR model might be different from zero. This kind of technique is used in order to avoid unrealistic regressions. In essence, it would not be justifiable from the physical point of view (see [19]) if AR regression for july had valid autoregressive parameters for may and no parameter different from zero for june. Thus, in terms of exhaustive enumeration, the analysis for the AR model examines several combinations for each season at time lags equal to $1, 2, \ldots, 23$.

The models we are dealing with here are special cases of generalized linear models (GLM). Link functions were restricted to identity. As noted by [21], link functions for Gamma distribution may be both inverse and identity as well as for Normal and Lognormal distributions.

[9, 19, 22] suggest the combination of three methods for residuals analysis:

- Residual Autocorrelation Function (RACF).
- Ljung−Box test.
- ACF test for square residuals.

---

[7]Residual mean fluctuate around 1 in a multiplicative model.

[8]BIC favors more parsimonious models, since it tends to penalize more heavily the inclusion of new variables (model selection and identification is well covered in [12]).

More attention was given to the first two methods, as it will be noted at the section dedicated to analyze the results.

## 2.1. Forecasting procedures

The role of forecasting analysis in any statistical modelling exercise is to provide an overview of how accurate the models chosen for the duty were and compare their performance in terms of certain criteria, eventually selecting some of them as references. In general it is not possible to assert that a specific model is better than the other ones, since there are many ways to evaluate it, not only from a statistical viewpoint but also in real world applications. However, forecasting analysis may provide important informations related to the accuracy of a new model when compared to traditional ones. The forecasting study conducted in this paper follows three approaches:

(i) The first procedure conveniently omits the last 3 years of data (2007–2009) in order to evaluate forecasting quality. The models are fitted to a shorter time series (1931–2006) and forecasting for different models are then calculated.

(ii) The second technique relies on simulation experiments. Every model is estimated using the complete set of data. For each of the three models 200 scenarios for 78 years were generated using the first two years of ENA time series and the set of estimated parameters as initial conditions.

(iii) The last procedure basically uses the basis of the second. The idea is to take the mean, standard deviation and skewness of each scenario and draw their distributions for each model from 1933 to 2009. Then one is able to analyse how well each model has performed regarding the closeness of these statistics to the one's generated from the original time series.

Exploratory data analysis are integrated to both procedures in order to interpret the results for periodic Gamma and Gaussian models. For practical purposes, forecasting studies are crucial to policymakers since they provide decision support for water resource management in hydroelectric systems.

## 2.2. Separate (non nested) tests and statistical measures for Gamma and Lognormal models

It is noteworthy that under regular circumstances the Gaussian template might lack accuracy in generating scenarios when compared to the others models since, for example, it might produce negative values that are not compatible with ENA series. For this reason we present a more detailed comparison between Lognormal and Gamma models. This task is addressed, on one hand, by performing a bootstrap exercise under the theory of separate or non nested tests and, on the other hand, by assessing their predictive powers through statistical measures, such as mean absolute percentage error and the modified version of classical root mean square error.

The first procedure relies on a direct application of modified log-likelihood ratio (LR) tests to the context of separated or non nested models. Typically, the usual LR test is used to compare two hypothesis $H_0$ (null) and $H_1$ (alternative), where in $H_0$ the parameter should belong to a subset of parametric space $\Theta_0$. Not rare, the test is conduced by using the following likelihood ratio statistics[9]

$$LR = 2\left(l(\hat{\theta}) - l(\tilde{\theta})\right) \qquad (2.8)$$

which is asymptotically $\chi_r^2$ (see details in [14]). The submodel obtained by restricting the parameter to the subset $\Theta_0$ is the Restricted Model. We could think about this test as a procedure to compare both restricted and unrestricted models.

Without loss of generality, let $f(y, \alpha)$ and $g(y, \beta)$ denote the Lognormal and Gamma densities for a random vector $Y$[10], respectively, where $\alpha$ is the parameter associated to the Lognormal model while $\beta$ is designed to

---

[9]$l(.)$ denotes the log likelihood, $\hat{\theta}$ is the unrestricted maximum likelihood estimate and $\tilde{\theta}$ is the unrestricted maximum likelihood estimate of $\theta$.

[10]Possibly conditional to the regressor $X$ or, as in the time series context, to lags of $Y$.

the gamma model. Someone could adopt mistakenly the same testing strategy to both Lognormal and Gamma models. However, as pointed out by [27, 28], when one tries to compare models associated to Lognormal and Gamma distributions is necessary to redefine the usual likelihood ratio, since both Lognormal and Gamma families of distributions are separate, in the sense that an arbitrary member of one family cannot be obtained as a limit of members of the other. Cox [10, 11] suggests to use modified (separate) likelihood ratio test statistics, as by

$$SLR_f = \log f(y, \hat{\alpha}) - \log g(y, \hat{\beta}) - E_{\hat{\alpha}}\{\log f(y, \alpha) - \log g(y, \beta_\alpha)\}$$
$$SLR_g = \log g(y, \hat{\beta}) - \log f(y, \hat{\alpha}) - E_{\hat{\beta}}\{\log g(y, \beta) - \log f(y, \alpha_\beta)\}. \tag{2.9}$$

The statistic $SLR_f$ adopts Lognormal distribution as the null hypothesis, while $SLR_g$ considers Gamma as the null hypothesis. It is important to say that both $\hat{\alpha}$ and $\hat{\beta}$ are the maximum likelihood estimates (MLE) of $\alpha$ and $\beta$, respectively. $\beta_\alpha$ [$\alpha_\beta$] denotes the *plim* of $\hat{\beta}$ [$\hat{\alpha}$] under the null of Lognormal [Gamma]. Moreover, $E_{\hat{\alpha}}$ [$E_{\hat{\beta}}$] indicates that the expected value is evaluated under the null of Lognormal [Gamma].

In [10, 11] is showed that, under some regularity conditions, $SLR_f$ and $SLR_g$ are asymptotically normally distributed with mean zero. The results for variances are presented in [26]. Therefore, in order to test Lognormal against Gamma (or Gamma against Lognormal), one possible way is to compare the ratio between the referred statistic and its standard deviations with the appropriate quantile of standard normal distribution. However, this task is far from being easy, since we must calculate, for example, very complex quantities, such as $E_{\hat{\alpha}}\{\log g(y, \beta_\alpha)\}$. Although much hard work had been done in this direction (see [26, 27]), we opt by an alternative solution, employing a boostrap procedure.

In our case, we have defined

$$BLR_f = \log f(y, \hat{\alpha}) - \log g(y, \hat{\beta}) \tag{2.10}$$

to test $H_0$: Lognormal *vs.* $H_1$: Gamma and

$$BLR_g = \log g(y, \hat{\beta}) - \log f(y, \hat{\alpha}) \tag{2.11}$$

to test $H_0$: Gamma *vs.* $H_1$: Lognormal.

The bootstrap allows us to use the original likelihood test statistic without appealing to asymptotic normal approximation. The main purpose is constructing or recovering the distribution of the test statistics under the null hypothesis, as described in [30]. We adopted the following scheme[11] to generate the proper results:

(1) Generate $R$ samples of size $T$ by sampling from fitted null model $g(y, \hat{\beta})$.
(2) For each $r$th simulated sample, $\hat{\alpha}_*^r(\hat{\beta})$ and $\hat{\beta}_*^r$ represent the parameter estimates obtained by maximising respectively the log likelihoods

$$f(y_*^r(\hat{\beta}), \hat{\alpha}) \text{ and } g(y_*^r(\hat{\beta}), \beta)$$

where $y_*^r(\hat{\beta})$ denotes the *rth* bootstrap sample conditional upon $\beta = \hat{\beta}$.
We then compute the simulated log likelihood ratio statistic

$$BLR_g^{*r} = \log g(y_*^r(\hat{\beta}), \hat{\beta}_*^r) - \log f(y_*^r(\hat{\beta}), \hat{\alpha}_*^r(\hat{\beta}))$$

(3) By constructing the empirical cdf of $\{BLR_g^{*r} : 1 \le r \le R\}$, we can compare the observed statistic $BLR_g$ with critical values obtained from the R independent (conditional) values of $BLR_g^{*r}$. The p-value based upon the bootstrap procedure is given by

$$P_R = \frac{\sum_{r=1}^{R} \mathbb{1}(BLR_g^{*r} \ge BLR_g)}{R},$$

where $\mathbb{1}$ is an indicator function.

---

[11]Here we present the general procedure where the null hypothesis is the Gamma model, but the reverse, in which Lognormal is the null hypothesis, is analogous.

The bootstrap procedure does not require the evaluation of limits in probability as presented in the former approach which is even more complex when one deals with time series and avoids the approximation problems associated with asymptotic theory, as discussed in [14].

The second procedure compares the performance of both Gamma and Lognormal models taking into account two deviation measures, the mean absolute percentage error (MAPE) and the modified version of classical root mean square error (RMSE), which we will denote by RMSE*[12]. The novelty is that the traditional statistic was divided by the sample mean of observed values. This allows us to compare results associated to different months and promotes a statistical measure with similar interpretation to MAPE.

## 3. Results

Box–Jenkins approach for model construction and forecasting procedures guide the analysis of results. Figure 1 shows the pattern of ENA measured in average megawatt (MWa) for southeast subsystem. As previously said, ENA time series comprises 79 years, from 1931 to 2009.

### 3.1. Periodic Gaussian model

Gaussian models are split into two different approaches for dealing with data. The first one simply takes the original data, while the second applies natural logarithmic transformation to the data set. Hereafter they are labeled as Gaussian Model and Lognormal Model, respectively.

Model identification of the Gaussian model was carried out through the combination of periodic PACF and exhaustive enumeration with BIC. Looking at Figure 2a, it is interesting to note that, for every month, PePACF cuts off at lag 1. This suggests that every season has at least an autoregressive model of order 1. For april, there is an indication of an $AR(2)$, whereas for july and october it appears to be a $AR(3)$. Almost all results are confirmed by BIC method, which selects a PAR model of order 14 – the fifth season has significant coefficients for lags 1, 2, 3, 12, 13 and 14 (this should be expected since PePACF cuts off at lag 14 for may). In addition, july and october confirm an $AR(3)$ and april selects an $AR(2)$. The remaining seasons have all an $AR(1)$. To check the validity of the estimated PAR models one may use the periodic ACF of residuals and Ljung−Box test. Both methods are depicted in Figures 2b and 2c, respectively. Taking into account the definition of periodic ACF of residuals, one may have serious doubt against the null hypothesis if the 5% confidence interval bars are cut off. In Figure 2c, there are 3 red lines representing the significance levels for 1%, 5% and 10%. Small p-values indicates that residuals should be autocorrelated. In general, only january and april provide some evidence against the null hypothesis, suggesting that residuals should be autocorrelated in these seasons. For the other seasons, in contrast, residuals seem to be non-correlated.

Looking into Lognormal model, periodic PACF indicates basically almost every results already seen at PePACF for the Gaussian model. BIC method selects the same order of the previous model (PAR(14)), with

---

[12]According to [32], if $\{y_t\}$ denotes the observed time series and $\hat{y}_t$ denotes the forecast, then the MAPE and the RMSE are defined respectively by:

$$MAPE \equiv \frac{1}{T} \sum_{t=1}^{T} \left| \frac{y_t - \hat{y}_t}{y_t} \right|$$

$$RMSE \equiv \sqrt{\frac{1}{T} \sum_{t=1}^{T} (y_t - \hat{y}_t)^2}.$$

Here we choose to use the modified version of the RMSE, as follows

$$RMSE^* \equiv \frac{RMSE}{\frac{1}{T} \sum_{t=1}^{T} y_t}.$$

FIGURE 1. ENA for the Southeast subsystem.

only one difference. October does not show significant coefficients for an $AR(3)$ but for an $AR(1)$. Besides, Figures 3b and 3c indicate that residuals are independent in every season in this model.

## 3.2. Periodic Gamma model

The general procedure to look into Gamma periodic model was basically the same as for Gaussian and Lognormal models. Only one adjustment was needed since conditional mean of PGAR model follows an additive model whereas its error term is multiplicative. BIC method selects a PGAR model of order 14 (again may is the only season to present significant coefficients for time lags $1, 2, 3, 12, 13$ and $14$), with the same order for the other months as Lognormal model. Nonetheless, almost every $p$-values associated to Ljung−Box test are above the 10% significance line, exception made to december, which achieves a p-value of 9.6%. These aspects are displayed in Figure 4b. Accordingly, there is only a weak evidence against the null hypothesis for the last season in this model, whereas for the others residuals seem to be independent. Periodic ACF of residuals also confirms this result.

## 3.3. Forecasting analysis

Forecasting experiments are largely used to compare the performance of two or more models and select the best ones according to certain criteria. Graphical tools are recognized for providing powerful insights related to statistical modelling and forecasting. In the case of this study, the methods described in Section 1 are accompanied of box plots in order to measure the predictive power of PAR models.

Regarding the first method, Figure 5 shows both the pattern of seasonal fluctuations of ENA time series forecasted from each model and original data set. It is interesting to note their remarkable resemblance in terms of predictive ability, at least when applied to the southeast subsystem. This is show in Figure 5d. The difference between them are almost indistinguishable. Lastly, all three models seem to capture poorly the last semester of 2009.

(A) Periodic PACF



(B) Periodic ACF of Residuals

(C) Ljung–Box Test

FIGURE 2. Periodic Gaussian model.

According to the second approach, there is a more pronounced difference between models. One possible drawback of the Gaussian model raised in this paper is the fact that they might forecast negative entries even though the data set have all their values greater than zero. This is precisely the case of the southeast subsystem, shown in Figure 6. On the left side of each figure there is the boxplot of the simulated model from the first year whereas on the right one is able to visualize the boxplot of the original data set. Despite outliers have been omitted for all simulations, they draw attention to a interesting fact: only Gaussian model reproduced outliers both above and below boxplot, which is reasonable to occur in these models, but contrasts with hydrological time series. In addition, the Gaussian model is the only one to generate outliers with negative values. Even whiskers for september, for example, almost touch the zero dotted line.

Regarding Lognormal and Gamma models, the boxplots are show in Figures 6b and 6c, respectively. Their predictive abilities and the way they reproduce time series patterns and trends are very similar. One might note that Gamma model captures a little better the seasons from 6 to 10, which in general represents the driest period in most regions in Brazil.

Another way to go deeply into forecasting analysis is evaluating every simulated models in terms of their capacity to reproduce the main parameters of unconditional distribution of data. In this case the analysis relies

(A) Periodic PACF



(B) Periodic ACF of Residuals



(C) Ljung–Box Test

FIGURE 3. Periodic Lognormal model.

on mean, standard deviation (SD) and skewness. These aspects become more evident in Figure 7. The three graphics on top of Figure 7 are related to Gaussian model whereas the others are the analogue of Lognormal and Gamma models, respectively. Besides, the red line in each one of nine graphics represents those three statistics obtained from the original data set. The blue "x" mark refers to the sample mean of simulated scenarios for mean, SD and skewness. It seems clear that all models are able to recognize accurately mean and standard deviation, with some minor exceptions, mainly from Lognormal model. But Gaussian model fails to identify skewness. Almost every red line from the original data set is above the box, some of them even outside the whiskers. It is not surprise that skewness from the Gaussian model fluctuates around zero, but again this does not apply to hydrological data. In fact, this is possibly what makes it generate negative values when applied to simulation experiments.

On the other hand, PGAR seems very competitive in relation to Lognormal model. It is not only able to recognize the mean and standard deviation in most seasons, but also it appears to be the best model to capture the skewness from ENA time series. Moreover, there seems to be a trade-off between Gaussian and Lognormal models regarding standard deviation and skewness. While the former achieves good results for standard deviation

(A) Periodic ACF of Residuals

(B) Ljung–Box Test

FIGURE 4. Periodic Gamma model.



(A) Periodic Gaussian Model

(B) Periodic Lognormal Model

(C) Periodic Gamma Model

(D) All Three Models

FIGURE 5. Predictive Power 2007–2009.

(A) Periodic Gaussian Model



(B) Periodic Lognormal Model

(C) Periodic Gamma Model

FIGURE 6. Simulated Path – Southeast subsystem 1933–2009.

but poor accuracy for skewness, the later does the opposite. This may be considered an advantage for Gamma model, since its performance is balanced between them two.

Gaussian model, despite of reproducing fairly well the seasonal pattern of hydrological time series, has some issues when used to simulate ENA time series. Besides, generally speaking, Gamma model has performed very well when compared to Lognormal model. These aspects will be further analysed on Section 3.4.

## 3.4. Gamma and Lognormal comparison

From the previous results and analysis it seems fairly conclusive that both Lognormal and Gamma models are more suitable to model hydrological time series than Gaussian model. In this subsection we go further in

FIGURE 7. Gaussian (*top*), Lognormal (*middle*) and Gamma (*bottom*).

this analysis presenting two additional approaches to compare them two. The first method is a statistical testing procedure that might indicate which model, under certain circumstances, prevails. The second method is based on individual statistics measures that are commonly used to evaluate the performances of statistical models.

For the first method we considered $R = 2000$ bootstrap samples in each analysis. Under the Gamma's null hypothesis, we obtained a $p$-value of $46, 45\%$, which indicates that we do not reject the Gamma model against

TABLE 1. MAPE and RMSE* for Gamma and Lognormal models.

|  | MAPE | | RMSE* | |
|  | Gamma | Lognormal | Gamma | Lognormal |
|---|---|---|---|---|
| Jan | 0.223 | 0.223 | 0.269 | 0.27 |
| Feb | 0.183 | 0.182 | 0.211 | 0.219 |
| Mar | 0.154 | 0.164 | 0.186 | 0.169 |
| Apr | 0.133 | 0.132 | 0.168 | 0.17 |
| May | 0.089 | 0.09 | 0.144 | 0.145 |
| Jun | 0.091 | 0.093 | 0.147 | 0.146 |
| Jul | 0.068 | 0.069 | 0.1 | 0.099 |
| Aug | 0.088 | 0.09 | 0.119 | 0.118 |
| Sep | 0.172 | 0.172 | 0.224 | 0.226 |
| Oct | 0.222 | 0.221 | 0.231 | 0.234 |
| Nov | 0.147 | 0.152 | 0.195 | 0.191 |
| Dec | 0.091 | 0.098 | 0.155 | 0.155 |
| Total | 0.138 | 0.141 | 0.218 | 0.218 |

the alternative Lognormal. Under the Lognormal's null hypothesis, on the other hand, we obtained a p-value of $53, 25\%$, therefore we do not reject the Lognormal model against the alternative Gamma. Although this result seems inconclusive, it is in line with our previous results. Both Gamma and Lognormal models are very competitive when dealing with hydrological data and have similar performances.

In terms of the second method, both models were used to generate one-period-ahead forecast for the period starting from January 1991 until December 2009. We have gotten one-step-ahead forecast for each month, totalling 228 forecasts. Table 1 exhibits the MAPE and RMSE* computed for each month and their respective aggregate values. Not surprisingly, both models show very similar performances. The Gamma model is slightly better at MAPE and the values of total RMSE* in each model are equal until the third decimal place and very close in monthly comparisons.

Additionally, we underline that, in terms of the results presented in Table 1, there was no significant difference in taking exponential transformation to Lognormal forecasts nor applying logarithmic transformation to ENA series. Also, forecasts based on more than one-step-ahead do not promote substantial distinction between models, despite increasing error rates (which is expected).

## 4. CONCLUSION

Periodic models have been the basis of hydrological time series modelling since the early 1960s (see [19]). In most applications, Gaussian conditional distribution was used to fit PAR models. The main purpose of this paper was to provide an application of PAR models following Gamma conditional distribution and compare its forecasting performance to Gaussian models. In this sense, Gamma model appeared to be very useful for ENA time series modelling, having predictive power very similar to Lognormal model. It would be interesting to reproduce this kind of study in real world stochastic optimization modelling and verify how well it would perform.

## REFERENCES

[1] H. Akaike, A New Look at the Statistical Model Identification. *IEEE Trans. Automat. Control* **19** (1974) 716–723.

[2] S. Almon, The Distributed Lag Between Capital Appropriations and Expenditures. *Econometrica* **33** (1965) 178–196.

[3] M. Barros, F.M. Mello and R.C. Souza, Aquisição de Energia no Mercado Cativo Brasileiro: Simulações dos Efeitos da Regulação sobre o Risco das Distribuidoras. *Pesquisa Operacional* **29** (2009) 303–322.

[4] B. Bezerra, L.A. Barroso, M. Brito, F. Porrua, B. Flach and M.V. Pereira, Measuring the Hydroelectric Regularization Capacity of the Brazilian Hydrothermal System. *Power and Energy Society General Meeting*, *Minneapolis* (2010) 1–7.

[5] G.E.P. Box and D.R. Cox, An Analysis of Transformation. *J. R. Stat. Soc. Ser. B* **26** (1964) 211–252.

[6] G.E.P. Box and D.A. Pierce, Distribution of Residual Autocorrelations in Autoregressive-Integrated Moving Average Time Series Models. *J. Am. Stat. Soc.* **65** (1970) 1509–1526.

[7] P.J. Bickel and K.A. Doksum, *Mathematical Statistics: Basic Ideas and Selected Topics*, vol I. Prectice-Hall (2006).

[8] M. Brito, M. Vianna, B. Bezerra, M.V. Pereira and L.A. Barroso, *Uma Metodologia para Analisar o Impacto das Usinas a Fio D'água na Capacidade de Regularização do Sistema Hidrotérmico Brasileiro*, Anais do XX SNPTEE (2009).

[9] P.J. Brocwell and R.A. Davis, Introduction to Time Series Forecasting. *Texts in Statistics*. Springer (2002).

[10] D.R. Cox, Tests of Separate Families of Hypotheses. Vol. 1 of *Proc. 4th Berkeley Symposium* (1961) 105–123.

[11] D.R. Cox, Further Results on Tests of Separate Families of Hypotheses. *J. R. Stat. Soc., Ser. B* **24** (1962) 406–424.

[12] G. Claeskens and N.L. Hjort, Model Selection and Model Averaging, *Cambridge Series in Statistical and Probabilistic Mathematics.* Cambridge University Press (2008).

[13] F.S. Costa, M.E.P. Maceira and J.M. Damazio, Modelos de Previsão Hidrológica Aplicados ao Planejamento de Operações do Sistema Elétrico Brasileiro, *Revista Brasileira de Recursos Hídricos* **12** (2007) 21–30.

[14] R. Davidson and J.G. MacKinnon, *Econometric Theory and Methods.* Oxford University Press (2004).

[15] B. Fernandez and J.D. Salas, Periodic Gamma Autoregressive Process for Operational Hydrology. *Water Resources Research* **22** (1986) 1385–1396.

[16] P.H. Franses and R. Paap, Periodic Time Series Models. *Advanced Texts in Econometrics*. Oxford (2004).

[17] D.P. Gaver and P.A.W. Lewis, First-order Autoregressive Gamma Sequences and Point Process. *Adv. Appl. Probab.* **12** (1980) 727–745.

[18] L.G. Godfrey and D.S. Poskitt, Testing the Restrictions of the Almon Lag Technique. *J. Am. Stat. Assoc.* **70** (1975) 105–108.

[19] K.W. Hipel and A.I. McLeod, Time Series Modelling of Water Resources and Environmental Systems. *Developments in Water Science.* Elsevier (1994).

[20] A.J. Lawrence, The Innovation Distribution of a Gamma Distributed Autoregressive Process. *Scand. J. Stat.* **9** (1982) 234–236.

[21] P. McCullagh and J. Nelder, Generalized Linear Models. *Monographs on Statistics and Applied Probability.* Chapman & Hall (1989).

[22] A.I. McLeod, Diagnostic Checking of Periodic Autoregressive Models with Application. *J. Time Series Anal.* **15** (1994) 221–233.

[23] J.A. Morgan and J.F Tatar, Calculation of the Residual Sum of Squares for All Possible Regressions. *Technometrics* **14** (1972) 317–325.

[24] M.E. Moss and M.C. Bryson, Autocorrelation Structure of Monthly Streamflows. *Water Resour. Res.* **10** (1974) 737–744.

[25] D.J. Noakes, A.I. McLeod and K.W. Hipel, Forecasting Monthly Riverflow Time Series. *Int. J. Forecast.* **1** (1986) 179–190.

[26] B. de B. Pereira, Empirical Comparisons of Some Tests of Separate Families of Hypothesis. *Metrika* **25** (1978) 219–234.

[27] B. de B. Pereira, Testes para Discriminar entre as Distribuições Lognormal, Gama e Weibull. *J. Inter.–An. Stat. Inst.* **33** (1979) 41–46.

[28] B. de B. Pereira, Choice of a Survival Model for Patients with Brain Tumour. *Metrika* **28** (1981) 53–61.

[29] M.V.F. Pereira and L.M.V.G. Pinto, Stochastic Optimization of a Multireservoir Hydroelectric System: A Decomposition Approach. *Water Resour. Res.* **21** (1985) 779–792.

[30] M.H. Pesaran and M. Weeks, Non-nested Hypothesis Testing: An Overview, in *Companion to Theoretical Econometrics*, edited by B.H. Baltagi. Basil Blackwell, Oxford (2001).

[31] J.D. Salas, J.W. Delleur, V. Yevjevich and W.L. Lane, *Applied Modelling of Hydrological Series.* Water Resources Publications (1980).

[32] R.S. Tsay, *Analysis of Financial Time Series.* Wiley, 3rd edition (2010).

[33] A.M. Vieira, B. de B. Pereira and P.R.H. Sales, Estimação Conjunta dos Parâmetros de um Modelo Multivariado Contemporâneo Periódico Auto-Regressivo − PAR(P). *Rev. Brasileira de Recursos Hídricos* **3** (1998) 5–17.

[34] G. Weiss, Time-reversibility of Linear Stochastic Processes. *J. Appl. Probab.* **12** (1975) 831–836.