# Zero-adjusted reparameterized Birnbaum–Saunders regression model

Vera Tomazella [a], Gustavo H.A. Pereira [a], Juvêncio S. Nobre [b],
Manoel Santos-Neto [a,c,*]

[a] *Departamento de Estatística, Universidade Federal de São Carlos, Brazil*
[b] *Departamento de Estatística e Matemática Aplicada, Universidade Federal do Ceará, Brazil*
[c] *Departamento de Estatística, Universidade Federal de Campina Grande, Brazil*

## ARTICLE INFO

## ABSTRACT

In this paper we present the zero-adjusted reparameterized Birnbaum–Saunders regression model. This new model generalizes at least seven reparameterized Birnbaum–Saunders regression models. Finally, an application to real data shows the potential of the model.

© 2019 Elsevier B.V. All rights reserved.

## 1. Introduction

Although distributions for zero-inflated count data have attracted the most attention, the study of continuous data with zero inflation has, in fact, a longer history. Concerns about statistical analysis of zero-inflated data were first identified in Aitchinson (1955) that proposed a mixture of zero with a lognormal model resulting in the well-known *delta distribution*. In the literature, there are various examples of inflated distributions, for example, Feueverger (1979), Farewell (1986), Meeker (1987), Lambert (1992), Shankar et al. (1997), Iwasaki and Daidoji (2009), Heller et al. (2006), Ospina and Ferrari (2008, 2012), Tong et al. (2013), Pereira et al. (2012), among others. Tu et al. (2014) briefly reviewed the concept of zero inflation and surveyed existing analytical methods for zero-inflated data. When a set of Poisson (PO) distributed count data includes more zeros than the amount expected under a PO probability mass function (PMF), it is declared to be "zero-inflated". Beyond the "zero-inflated" term the idea of joining a degenerate model at zero with a continuous model can be named "zero-adjusted" or "zero-augmented" (Galvis et al., 2014). In this article, we use the term "zero-adjusted", as in Heller et al. (2006) and Leiva et al. (2016).

Leiva et al. (2016) proposed a methodology for inventory logistics that allows demand data that have zeros to be modeled by means of a new discrete–continuous mixture distribution, which is constructed by using a probability mass at zero and

---

a continuous component related to the RBS distribution (ZARBS). In this article, we shall allow the mean and the precision parameter of the BS distribution and the probability of a point mass at 0 to be related to linear or non-linear predictors through link functions. This model generalizes the models proposed by Leiva et al. (2014) and Santos-Neto et al. (2016).

Apart from this introduction, the paper is organized as follows. Section 2 presents a general class of ZARBS regression models. Section 3 is devoted to diagnostic analysis. Section 4 contains two applications using real data and concluding remarks are given in Section 5. Some technical details are collected in an appendix.

## 2. ZARBS model

We say that a random variable $Y$ follows a ZARBS (Leiva et al., 2016) distribution with parameters $\mu$, $\sigma$ and $\nu$, denoted by ZARBS($\mu, \sigma, \nu$), if the distribution of $Y$ admits the following PDF

$$f_Y(y|\mu, \sigma, \nu) = \left\{ [1-\nu] \frac{\exp(\frac{\sigma}{2})\sqrt{\sigma+1}}{4 y^{3/2} \sqrt{\pi\mu}} \left[ y + \frac{\sigma\mu}{\sigma+1} \right] \exp\left( -\frac{\sigma}{4} \left[ \frac{\{\sigma+1\}y}{\sigma\mu} + \frac{\sigma\mu}{\{\sigma+1\}y} \right] \right) \right\}^{1-\mathbb{I}(y=0)} \times \nu^{\mathbb{I}(y=0)},$$

where $y > 0, \sigma > 0, \mu > 0$ and $0 < \nu < 1$ and $\mathbb{I}(\mathcal{S})$ is the indicator function of the set $\mathcal{S}$. We have that the mean and variance of $Y \sim$ ZARBS($\mu, \sigma, \nu$) are, respectively, given by $E[Y] = (1-\nu)\mu$ and $Var[Y] = (1-\nu)\mu^2 \{\nu + CV[T]^2\}$, where $CV[T] = \frac{\sqrt{(2\sigma+5)}}{(\sigma+1)}$. The log-PDF of the ZARBS distribution is given by

$$\ell(\mu, \sigma, \nu) = \log(\nu)\mathbb{I}(y=0) + \log(1-\nu)[1 - \mathbb{I}(y=0)] + [1 - \mathbb{I}(y=0)] \log(f_T(y)). \tag{1}$$

Let $Y_1, \ldots, Y_n$ be a random sample from $Y \sim$ ZARBS($\mu, \sigma, \nu$). Then, the corresponding likelihood function for $\boldsymbol{\theta} = [\mu, \sigma, \nu]^\top$ is given by

$$L(\mu, \sigma, \nu) = \prod_{i=1}^n f_Y(y_i|\mu, \sigma, \nu) = \nu^{n_0}[1-\nu]^{n-n_0} \prod_{i=1}^n f_T(y_i|\mu, \sigma)^{1-\mathbb{I}(y_i=0)}, \tag{2}$$

where $n_0 = \sum_{i=1}^n \mathbb{I}(y_i = 0)$ is the number of observations equal to zero. Hence, the corresponding log-likelihood function obtained from (2) can be expressed as $\ell(\boldsymbol{\theta}) = \ell(\nu) + \ell(\mu, \sigma)$, where

$$\ell(\nu) = n_0 \log(\nu) + [n - n_0] \log(1 - \nu), \tag{3}$$

$$\ell(\mu, \sigma) = [n - n_0]c(\mu, \sigma) - \frac{3}{2} \sum_{y_i > 0} \log(y_i) - \frac{[\sigma+1]}{4\mu} \sum_{y_i > 0} y_i - \sum_{y_i > 0} \frac{\mu\sigma^2}{4[\sigma+1]y_i} + \sum_{y_i > 0} \log\left( y_i + \frac{\mu\sigma}{[\sigma+1]} \right), \tag{4}$$

with $c(\mu, \sigma) = -[1/2]\log(16\pi) + [\sigma/2] - [1/2]\log(\mu) + [1/2]\log(\sigma+1)$. Let $Y_1, \ldots, Y_n$ be a random sample such that each $Y_i \sim$ ZARBS($\mu_i, \sigma_i, \nu_i$). Suppose the mean, precision and mixture parameters of $Y_i$ satisfy the following functional relations: $g_1(\mu_i) = \eta_i = f_1(\mathbf{x}_i; \boldsymbol{\beta})$, $g_2(\sigma_i) = \tau_i = f_2(\mathbf{z}_i; \boldsymbol{\alpha})$, and $g_3(\nu_i) = \xi_i = f_3(\mathbf{w}_i; \boldsymbol{\gamma})$, for $i = 1, \ldots, n$, where $\boldsymbol{\beta} = [\beta_1, \ldots, \beta_p]^\top$, $\boldsymbol{\alpha} = [\alpha_1, \ldots, \alpha_q]^\top$ and $\boldsymbol{\gamma} = [\gamma_1, \ldots, \gamma_r]^\top$ are vectors of unknown parameters to be estimated, for $p + q + r < n, \boldsymbol{\eta} = [\eta_1, \ldots, \eta_n]^\top$, $\boldsymbol{\tau} = [\tau_1, \ldots, \tau_n]^\top$ and $\boldsymbol{\xi} = [\xi_1, \ldots, \xi_n]^\top$ are predictor vectors, and $f_j(\cdot; \cdot), j = 1, 2, 3$ are linear or nonlinear twice continuously differentiable functions in the second argument, such that the derivative matrices $\widetilde{\mathbf{X}} = \partial\boldsymbol{\eta}/\partial\boldsymbol{\beta}, \widetilde{\mathbf{Z}} = \partial\boldsymbol{\tau}/\partial\boldsymbol{\alpha}$ and $\widetilde{\mathbf{W}} = \partial\boldsymbol{\xi}/\partial\boldsymbol{\gamma}$ are full ranks for all $\boldsymbol{\beta}, \boldsymbol{\alpha}$ and $\boldsymbol{\gamma}$. Moreover, $\tilde{\mathbf{x}}_i = [\tilde{x}_{i1}, \ldots, \tilde{x}_{ip_1}]^\top, \tilde{\mathbf{z}}_i = [\tilde{z}_{i1}, \ldots, \tilde{z}_{ip_2}]^\top$ and $\tilde{\mathbf{w}}_i = [\tilde{w}_{i1}, \ldots, \tilde{w}_{ip_3}]^\top$ are vectors that contain the values of $p_1, p_2$ and $p_3$ explanatory variables, respectively. In this model, the link functions $g_j: \mathbb{R}^+ \rightarrow \mathbb{R}, j = 1, 2$ are strictly monotone, positive, and at least twice differentiable and $g_3: (0, 1) \rightarrow \mathbb{R}$ is strictly monotonic and twice differentiable.

The log-likelihood function for this model is given by $\ell(\mu_i, \sigma_i, \nu_i) = \ell(\nu_i) + \ell(\mu_i, \sigma_i)$, with the expressions of $\ell(\nu_i)$ and $\ell(\mu_i, \sigma_i)$ given as (3). Here, $\mu_i = g_1^{-1}(\eta_i), \sigma_i = g_2^{-1}(\tau_i)$ and $\nu_i = g_3^{-1}(\xi_i)$. We note that $\ell(\mu_i, \sigma_i)$ is the log-likelihood function for $[\boldsymbol{\beta}, \boldsymbol{\alpha}]^\top$ in a nonlinear RBS regression model with varying precision. Furthermore, $\ell(\nu_i)$ represents the log-likelihood function of a regression model for binary responses. The score function, which is given by $\mathbf{U}_{\boldsymbol{\theta}} = [\mathbf{U}_{\boldsymbol{\beta}}^\top, \mathbf{U}_{\boldsymbol{\alpha}}^\top, \mathbf{U}_{\boldsymbol{\gamma}}^\top]^\top$, where $\boldsymbol{\theta} = [\boldsymbol{\beta}^\top, \boldsymbol{\alpha}^\top, \boldsymbol{\gamma}^\top]^\top$ and $\mathbf{U}_{\boldsymbol{\beta}} = \widetilde{\mathbf{X}}^\top \mathbf{A}_\kappa [\mathbf{y}^* - \boldsymbol{\mu}^*]$, $\mathbf{U}_{\boldsymbol{\alpha}} = \widetilde{\mathbf{Z}}^\top \mathbf{B}_\kappa [\mathbf{y}^\bullet - \boldsymbol{\sigma}^\bullet]$, and $\mathbf{U}_{\boldsymbol{\gamma}} = \widetilde{\mathbf{W}}^\top \mathbf{C} [\mathbf{y}^\circ - \boldsymbol{\nu}^\circ]$, and these results are shown in (A2) of Appendix A. Thus, the Fisher information matrix and the its elements are presented in (A.5) of Appendix A.

Theoretical results of this paper have been implemented in the **R** (R-Team, 2017) software. In special, the **RBS** package provides functions for fitting RBS regression models using the **gamlss** package. The current version can be downloaded from GitHub via

```
1  devtools::install_github("santosneto/RBS")
```

The proposed specification for the ZARBS regression model follows the traditionally used framework in continuous zero-adjusted models. Another structure for continuous zero-adjusted models was proposed by Labrecque-Synnott and Angers (2009). The structure proposed by them has advantages over the traditional one when values of the response variable equal to zero or close to zero are, from a practical point of view, equivalent. In the case of the applications of this paper and in most of the problems in which the ZARBS regression model can be used, the response variable does not have these characteristics. For this reason, we proposed the ZARBS regression model using the traditional structure, which is easier to interpret.

## 3. Diagnostic analysis

### 3.1. Residuals

To assess goodness of fit and departures from the assumptions of the ZARBS regression model, we propose the randomized quantile residual (Dunn and Smyth, 1996). It is a randomized version of the Cox and Snell (1968) residual and is given by

$$r_i^q = \Phi^{-1}(f_i)\mathbb{I}(y_i > 0) + \Phi^{-1}(u_i)\mathbb{I}(y_i = 0),$$

where $\Phi^{-1}(\cdot)$ is the inverse function of the standard normal cumulative distribution function (CDF), $f_i = F(y_i, \hat{\mu}_i, \hat{\sigma}_i, \hat{\nu}_i)$ is the estimated value using the CDF of the ZARBS distribution and $u_i$ is the observed value of a random variable $U_i$ with uniform distribution in $(0, \hat{\nu}_i)$, where $\hat{\nu} = \widehat{\Pr}(Y_i = 0)$. If the model is correctly specified, then $r_i^q$ is asymptotically standard normally distributed.

### 3.2. Local influence

The concept of local influence was introduced by Cook (1986). This approach based on normal curvature is an important diagnostic tool for assessing local influence of minor perturbations to a statistical model. The normal curvature for $\boldsymbol{\theta}$ in the direction vector $\mathbf{d}$, with $\|\mathbf{d}\| = 1$, is expressed as $C_{\mathbf{d}} = 2|\mathbf{d}^\top \boldsymbol{\Delta}^\top \mathbf{h}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Delta} \mathbf{d}|$, where $\boldsymbol{\Delta} = \partial^2 \ell(\boldsymbol{\theta}|\boldsymbol{\omega})/\partial\boldsymbol{\theta}\partial\boldsymbol{\omega}$ is the matrix of perturbations. A local influence diagnostic tool is generally based on index plots. The index plot of the direction of maximum curvature, $\mathbf{d}_{\max}$, corresponding to the maximum eigenvalue of $\mathbf{F}_{\boldsymbol{\theta}} = \boldsymbol{\Delta}^\top \mathbf{h}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Delta}$, $C_{\mathbf{d}_{\max}}$, evaluated at $\boldsymbol{\theta} = \widehat{\boldsymbol{\theta}}$ and $\boldsymbol{\omega} = \boldsymbol{\omega}_0$, can contain valuable diagnostic information. Also, the index plot of $C_i = 2|f_{ii}|$, where $f_{ii}$ is the $i$th diagonal element of $\mathbf{F}_{\boldsymbol{\theta}}$, can be used as a diagnostic technique to evaluate the existence of influential observations. Those cases when $\widehat{C}_i > 2 \cdot \overline{C}$, where $\overline{C} = \sum_{i=1}^n \widehat{C}_i/n$, are considered as potentially influential.

#### Case-weights perturbation

We define a perturbation $\boldsymbol{\omega} = [\omega_1, \ldots, \omega_n]^\top$, $\omega_i \geq 0$, to perturb the contribution of each case to the log-likelihood. In the ZARBS regression model, we have that the form of $\Delta_{\boldsymbol{\beta}}$, $\Delta_{\boldsymbol{\alpha}}$ and $\Delta_{\boldsymbol{\gamma}}$ is, respectively, given by $\Delta_{\boldsymbol{\beta}} = \widetilde{\mathbf{X}}^\top \mathbf{A}_\kappa \left[(y_i^* - \mu_i^*)\delta_{ii}^n\right]$, $\Delta_{\boldsymbol{\alpha}} = \widetilde{\mathbf{Z}}^\top \mathbf{B}_\kappa \left[(y_i^\bullet - \sigma_i^\bullet)\delta_{ii}^n\right]$ and $\Delta_{\boldsymbol{\gamma}} = \widetilde{\mathbf{W}}^\top \mathbf{C} \left[(y_i^\circ - \nu_i^\circ)\delta_{ii}^n\right]$, evaluated at $\widehat{\boldsymbol{\theta}}$ and $\boldsymbol{\omega}_0 = \mathbf{1}_n$, that is, an $n \times 1$ column vector with all elements equal to 1 and $\delta_{ii}^n$ is a diagonal matrix $n \times n$ (see details in Appendix A).

## 4. Applications

### 4.1. ZARBS regression model

In this application the data set refers to the study carried out by the Institute of Biomedical Sciences, University of São Paulo, that studied the fumonisin production by Fusarium verticillioides in corn grains (Rocha et al., 2017). Here, we study the relation between the variables Fusarium verticillioides ($x_1$), Time ($x_2$), Precipitation ($x_3$) and Water activity ($x_4$) and the fumonisin B2 (FB2) production ($y$). Based on the initial analysis, we will assume the response $Y_i \sim \text{ZARBS}(\mu_i, \sigma_i, \nu_i)$ for **FUMCorn** data. Firstly, we consider the following regression model

$$\log(\mu_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4}, \quad \sigma_i = \alpha_0,$$
$$\log\left(\frac{\nu_i}{1 - \nu_i}\right) = \gamma_0 + \gamma_1 w_{i1} + \gamma_2 w_{i2} + \gamma_3 w_{i3} + \gamma_4 w_{i4}, \quad i = 1, 2, \ldots, 364, \tag{5}$$

where $\boldsymbol{\beta} = [\beta_0, \beta_1, \beta_2, \beta_3, \beta_4]^\top$ and $\boldsymbol{\gamma} = [\gamma_0, \gamma_1, \gamma_2, \gamma_3, \gamma_4]^\top$ being the regression coefficients and $x_{ij} \equiv w_{ij}(j = 1, 2, 3, 4)$ are the values of the regressor $\mathbf{X}$. We fit the ZARBS model by using the **RBS** package presented in Section 2. The MLE's of the model parameters are (with approximate estimated standard errors): $\hat{\beta}_0 = -8.509051(15.928556)$, $\hat{\beta}_1 = 0.016231(0.006918)$, $\hat{\beta}_2 = 0.057854(0.015256)$, $\hat{\beta}_3 = -0.007053(0.047547)$, $\hat{\beta}_4 = 3.900954(16.048109)$, $\hat{\alpha}_0 = 1.1329(0.2349)$, $\hat{\gamma}_0 = -7.938496(14.660591)$, $\hat{\gamma}_1 = -0.025994(0.007319)$, $\hat{\gamma}_2 = -0.046594(0.015765)$, $\hat{\gamma}_3 = 0.174875(0.047704)$, $\hat{\gamma}_4 = 11.492422(14.747538)$. We note that predictors $x_3$ and $x_4$ seem to be not statistically significant at 5% in the mean model and the explanatory variable $x_4$ is not statistically significant in proportion model. The normal probability plot with envelope for the quantile residual, is shown in Fig. 1. This figure does not show unusual features, so that the assumption that the response variable follows a ZARBS distribution does not seem to be unsuitable.

Fig. 1 displays index plots of $C_i$, from where observation #314 was detected as potentially influential. We will investigate their impact on the model inference when it is removed. The inferences do not change at the significance level of 5%. That is, the explanatory variables $x_3$ and $x_4$ should be removed from the mean model and explanatory variable $x_4$ should be removed from proportion model. Thus, the final ZARBS regression model is given by $\log(\hat{\mu}_i) = \underset{(0.6919)}{-4.7261} + \underset{(0.0053)}{0.0150}x_{1i} +$

$\underset{(0.0115)}{0.0573}x_{2i}$ and $\log\left(\frac{\hat{\nu}_i}{1-\hat{\nu}_i}\right) = \underset{(0.6879)}{3.4719} - \underset{(0.0074)}{0.0262}w_{1i} - \underset{(0.0132)}{0.0532}w_{2i} + \underset{(0.0464)}{0.1678}w_{3i}$, with $\hat{\sigma} = 1.1266(0.2333)$. More details about this application can be found in https://www.santosnetoce.com.br/spl. A second application is also available.
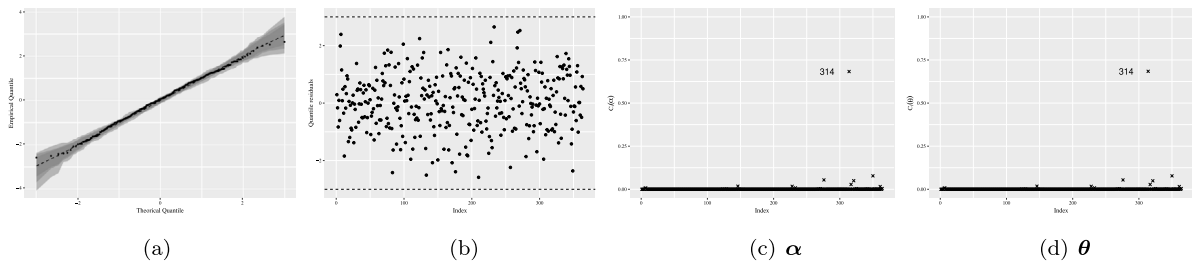
**Fig. 1.** QQ-plot with simulated envelope of quantile residuals (a), index versus quantile residuals (b) and $C_i$'s plots.

## 5. Concluding remarks

In this paper, we have introduced the ZARBS regression model. We have modeled the mean, the precision parameter and the probabilities of occurrences of zeros, through linear and/or nonlinear predictors, using appropriate link functions. In addition, we have proposed randomized quantile residuals for our model. An iterative estimation procedure and its computational implementation have also been discussed. Moreover, we have considered a perturbation scheme which allows the identification of observations that exert unusual influence on the estimation process. In general, the results of the applications have shown the potentiality of proposed methodology.

## Acknowledgments

## Appendix A. Supplementary data

Supplementary material related to this article can be found online at https://doi.org/10.1016/j.spl.2019.01.019.

## References

Aitchinson, J., 1955. On the distribution of a positive random variable having a discrete probability mass at the origin*. J. Amer. Statist. Assoc. 50, 901–908.
Cook, R.D., 1986. Assessment of local influence. J. R. Stat. Soc. B 48, 133–169.
Cox, D., Snell, E., 1968. A general definition of residuals. J. R. Stat. Soc. B 2, 248–275.
Dunn, P., Smyth, G., 1996. Randomized quantile residuals. J. Comput. Graph. Statist. 5, 236–244.
Farewell, V.T., 1986. Mixture models in survival analysis: Are they worth the risk?. Can. J. Stat. / La Rev. Can. de Stat. 14, 257–262.
Feuerverger, A., 1979. Some methods of analysis for weather experiments. Biometrika 66, 655–658.
Galvis, D.M., Bandyopadhyay, D., Lachos, V.H., 2014. Augmented mixed beta regression models for periodontal proportion data. Stat. Med. 33, 3759–3771.
Heller, G., Stasinopoulos, M., Rigby, B., 2006. The zero-adjusted inverse gaussian distribution as a model for insurance claims, in: Proceedings of the 21th International Workshop on Statistical Modelling, volume 226233.
Iwasaki, M., Daidoji, K., 2009. Zero-inflated probability models and their applications to the analysis of test scores. Kodo Keiryogaku (Jap. J. Behav.) 36, 25–34.
Labrecque-Synnott, F., Angers, J.-F., 2009. An Extension of zero-modified models to the continuous case. In: Technical Report. CRM-3289.
Lambert, D., 1992. Zero-inflated Poisson regression with an application to defects in manufacturing. Technometrics 34, 1–14.
Leiva, V., Santos-Neto, M., Cysneiros, F.J.A., Barros, M., 2014. Birnbaum-Saunders statistical modelling: A new approach. Stat. Model. 14, 21–48.
Leiva, V., Santos-Neto, M., Cysneiros, F.J.A., Barros, M., 2016. A methodology for stochastic inventory models based on a zero-adjusted Birnbaum-Saunders distribution. Appl. Stoch. Models Bus. Ind. 32, 74–89.
Meeker, W.Q., 1987. Limited failure population life tests: Application to integrated circuit reliability. Technometrics 29, 51–65.
Ospina, R., Ferrari, S.L.P., 2008. Inflated beta distributions. Statist. Papers 51, 111.
Ospina, R., Ferrari, S.L.P., 2012. A general class of zero-or-one inflated beta regression models. Comput. Statist. Data Anal. 56, 1609–1623.
Pereira, G.H.A., Botter, D.A., Sandoval, M.C., 2012. The truncated inflated beta distribution. Comm. Statist. Theory Methods 41, 907–919.
R-Team, 2017. R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria.
Rocha, L.O., Reis, G.M., Fontes, L.C., Piacentini, K.C., Barroso, V.M., Reis, T.A., Pereira, A.A., Corrêa, B., 2017. Association between fum expression and fumonisin contamination in maize from silking to harvest. CR. Prot. 94, 77–82.
Santos-Neto, M., Cysneiros, F., Leiva, V., Barros, M., 2016. Reparameterized birnbaum-saunders regression models with varying precision. Electron. J. Stat. 10, 2825–2855.
Shankar, V., Milton, J., Mannering, F., 1997. Modeling accident frequencies as zero-altered probability processes: An empirical inquiry. Accid. Anal. Prev. 29, 829–837.
Tong, E.N.C., Mues, C., Thomas, L., 2013. A zero-adjusted gamma model for mortgage loan loss given default. Int. J. Forecast. 29, 548–562.
Tu, W., Liu, H., data, Zero-inflated., 2014. Zero-inflated data in: Wiley StatsRef: Statistics Reference Online. John Wiley & Sons, Ltd.