

Smooth tests of goodness of fit for the distributional assumption of regression models

Peer-reviewed author version

Rayner, J. C. W.; Rippon, Paul; Suesse, Thomas & THAS, Olivier (2022) Smooth tests of goodness of fit for the distributional assumption of regression models. In: AUSTRALIAN & NEW ZEALAND JOURNAL OF STATISTICS, 64 (1) , p. 67 -85.

DOI: 10.1111/anzs.12361

Handle: <http://hdl.handle.net/1942/37334>

Smooth tests of goodness of fit for the distributional assumption of regression models

J. C. W. Rayner^{1,3}, Paul Rippon², Thomas Suesse³ and Olivier Thas^{3,4,5*}

University of Newcastle, University of Wollongong, Hasselt University and Ghent University

Summary

We focus on regression models that consist of (1) a model for the conditional mean of the outcome and (2) a distributional assumption about the distribution of the outcome, both conditional on the regressors. Generalised linear models (GLM) form a well known example. The choice of the outcome distribution is often motivated by prior or background knowledge of the researcher, or it is simply chosen for convenience. We propose smooth goodness of fit tests for testing the distributional assumption in regression models. The tests arise from embedding the regression model in a smooth family of alternatives, and constructing appropriate score tests that correctly account for nuisance parameter estimation. The tests are customised, focussed and comprehensive. We present several examples to illustrate the wide applicability of our method. A small simulation study demonstrates that our tests have power to detect important deviations from the hypothesised model.

Key words: GLM; model diagnostics; Poisson regression; score test; ZIP regression

1. Introduction

Consider regression models for a univariate outcome variable Y and p regressors x_i ($i = 1, \dots, p$), stacked into the vector $\mathbf{x}^\top = (x_1, \dots, x_p)$. For notational comfort, the first regressor $x_1 = 1$ if an intercept is required in the model, and, similarly, regressors x_i may also represent interaction effects or other transformations of regressors. We will focus on models that are specified as

$$E(Y | \mathbf{x}) = g^{-1}(\mathbf{x}^\top \boldsymbol{\beta}) = \mu(\mathbf{x}^\top \boldsymbol{\beta}) \quad (1)$$

$$Y | \mathbf{x} \sim f(\cdot; \mu(\mathbf{x}^\top \boldsymbol{\beta}), \boldsymbol{\gamma}). \quad (2)$$

*Author to whom correspondence should be addressed.

¹Centre for Computer-Assisted Research Mathematics and its Applications, University of Newcastle, Australia

²School of Mathematical and Physical Sciences, University of Newcastle, Australia

³National Institute of Applied Statistics Research Australia (NIASRA), University of Wollongong, Australia

⁴I-BioStat, Data Science Institute, Hasselt University, Belgium

⁵Department of Applied Mathematics, Computer Science and Statistics, Ghent University, Belgium

Email: olivier.thas@uhasselt.be

14 Equation (1) specifies the conditional mean of Y as a function of the linear predictor $\mathbf{x}^\top \boldsymbol{\beta}$
15 through the link function g . We will often use the notation $\mu(\mathbf{x}^\top \boldsymbol{\beta})$ to denote the conditional
16 mean. Equation (2) states that the density function of the conditional distribution of Y given
17 \mathbf{x} is given by f , which depends on the conditional mean μ and on an additional t -dimensional
18 nuisance parameter $\boldsymbol{\gamma}^\top = (\gamma_1, \dots, \gamma_t)$. The $\boldsymbol{\gamma}$ -parameters may be related to the conditional
19 variance of the outcome. For example, for linear regression $\text{var}(Y|\mathbf{x}) = \gamma$. Note that the
20 model also allows the conditional variance $\text{var}(Y|\mathbf{x})$ to depend on the mean $\mu(\mathbf{x}^\top \boldsymbol{\beta})$.

21 Generalised linear models (GLM), which were first introduced by Nelder & Wedderburn
22 (1972), form a special class. They arise if f belongs to a one-dimensional (dispersion)
23 exponential family.

24 Equations (1) and (2) completely specify the conditional outcome distribution, and hence
25 the maximum likelihood (ML) framework can be used for inference on the target parameter
26 $\boldsymbol{\beta}$. The distributional component is often motivated by prior or background knowledge on
27 the probabilistic mechanism that generated the data. For example, count outcomes are known
28 to be often well described by a Poisson distribution (f is Poisson), and binary outcomes
29 often behave like a Bernoulli distribution (f is binomial/Bernoulli). These two examples
30 result in a GLM because Poisson and binomial distributions belong to the exponential
31 family. In modern genomics applications, RNASeq and microbiome 16S RNA sequencing
32 experiments are believed to give count outcomes that can be described by negative binomial
33 (NB) distributions; see e.g., Love, Huber & Anders (2014) and McMurdie & Holmes
34 (2014), which are overdispersed Poisson distributions that contain an overdispersion nuisance
35 parameter. The overdispersion is explained by the biological variability on top of the technical
36 variability that is described by the Poisson distribution. For single cell RNASeq experiments,
37 several papers suggest that the count outcomes should be modelled with a zero-inflated
38 negative binomial (ZINB) distribution; see e.g. Risso et al. (2017). Also in other biological
39 applications, counts often show more zeroes than expected under a Poisson distribution and
40 a zero-inflated Poisson (ZIP) distribution has been suggested to be more appropriate than
41 a Poisson (Thas & Rayner 2005). The zero-inflation is often explained by a second data-
42 generating mechanism that causes the zero counts. The ZIP, NB and ZINB distributions do
43 not belong to the exponential family, but they are still regression models of the form (1) and
44 (2), and hence they fall within the scope of this paper.

45 Valid asymptotic statistical inference on $\boldsymbol{\beta}$ requires a correct specification of the
46 conditional mean model, and hence several papers have proposed diagnostic methods for
47 detecting violations to the mean model (1); see e.g. Stute & Zhu (2002); Khmaladze et al.
48 (2004); Hart (2013). However, a correct specification of the variance, or the mean-variance
49 relationship is also required. The Quasi-Likelihood approach (Wedderburn 1974) builds
50 upon a semiparametric model with only mean and variance(-mean) specifications. When

51 the latter is misspecified, the variance of the mean model parameter estimators can be
52 estimated with a robust sandwich estimator, but this shows poor small sample behaviour
53 (Kauermann & Carroll 2001). For GLMs, Huang & Rathouz (2017) demonstrated that the
54 mean model parameters show orthogonality to the outcome distribution, opening the door to
55 first nonparametrically estimate the outcome distribution, and subsequently use this estimate
56 in an empirical likelihood for the estimation of the mean model parameters. Also see Huang
57 (2014). Upon using the orthogonality, the authors showed that asymptotically no efficiency is
58 lost as compared to ML in the correctly specified GLM. In small samples, however, the ML
59 estimator still shows better performance.

60 Given the arguments and examples provided in the previous paragraphs, we conclude
61 that assessing the distributional assumption, contained in (2), is also of scientific importance.
62 Relatively few methods have been proposed for testing this distributional assumption; see,
63 for example, Dean & Lawless (1989) and Peña & Slate (2006). In practice, formal hypothesis
64 testing can be complemented with graphical inspection of the model fit. Residual plots may
65 be used for assessing the mean model, but in general QQ-plots of residuals cannot be used
66 for gauging the distributional assumption, unless the distribution belongs to a location-shift
67 class (e.g. the normal distribution). QQ-plots and other visualisations based on residuals were
68 strongly promoted by John Tukey in many of his works. We cannot agree more with him
69 that this is indeed of primordial importance for all data analyses. He also developed formal
70 statistical tests for assessing the normality assumption in linear models, but these methods
71 are restricted to additive two-way analysis of variance (Anscombe & Tukey 1963).

72 In this paper we propose smooth tests of goodness of fit for testing the distributional
73 assumption contained in (2). Smooth tests are well established for testing goodness of fit in
74 the one-sample problem and they can properly account for nuisance parameter estimation.
75 We refer to Rayner, Thas & Best (2009) for a comprehensive overview of the general theory
76 and for several examples, including the normal, Poisson, NB and ZIP distributions. Here we
77 extend the smooth testing method to the regression context as described above. We formulate
78 the theory for regression models of the form (1) and (2), and we show how simplifications
79 arise for special cases, including the class of GLMs (Section 2). In Section 3 we give the
80 Poisson, normal, and ZIP distributions as examples. Small simulation studies for Poisson and
81 ZIP models are presented in Section 4, and conclusions are formulated in Section 5.

82 **2. Smooth tests**

83 The general construction of the test is given in Section 2.1, and special cases resulting
84 in simplifications are provided in Section 2.2.

85 2.1. The smooth test for general regression models

86 The development of the theory is very similar to the development detailed in Rayner,
87 Thas & Best (2009, chapters 6 and 8). We therefore limit the exposition here to the overall
88 procedure, with a focus on the details for the regression setting, and to the most important
89 results. Proofs are deferred to Appendix A.1 in Supporting and Supplementary Material.

90 *Derivation of the smooth test*

The construction starts with nesting the density function of the regression model (the conditional distribution of the outcome Y , given the regressor \mathbf{x}) in a family of distributions indexed by the parameter $\boldsymbol{\theta}^\top = (\theta_1, \dots, \theta_k)$ and then deriving the score test for testing $\boldsymbol{\theta} = \mathbf{0}$ against $\boldsymbol{\theta} \neq \mathbf{0}$. The number k refers to the order of the alternative. The order k embedding density of $Y \mid \mathbf{x}$ is given by

$$f_k(y; \mu(\mathbf{x}^\top \boldsymbol{\beta}), \gamma, \boldsymbol{\theta}) = C(\mu(\mathbf{x}^\top \boldsymbol{\beta}), \gamma, \boldsymbol{\theta}) \exp \left(\sum_{i=1}^k \theta_i h_i(y; \mu(\mathbf{x}^\top \boldsymbol{\beta}), \gamma) \right) f(y; \mu(\mathbf{x}^\top \boldsymbol{\beta}), \gamma), \quad (3)$$

91 where C is a normalisation constant and $\{h_i\}$ is a set of functions that are orthonormal
92 to the regression model with density function f . Since both C and $\{h_i\}$ differ from their
93 counterparts of the one-sample smooth tests in the sense that they depend on the regressor \mathbf{x} ,
94 we provide some more details, but first the log-likelihood function is given for a sample of n
95 independently sampled observations $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$:

$$\begin{aligned} l(\boldsymbol{\beta}, \gamma, \boldsymbol{\theta}) &= \sum_{j=1}^n \log f_k(y_j; \mu(\mathbf{x}_j^\top \boldsymbol{\beta}), \gamma, \boldsymbol{\theta}) \\ &= \sum_{j=1}^n \log C(\mu(\mathbf{x}_j^\top \boldsymbol{\beta}), \gamma, \boldsymbol{\theta}) + \sum_{i=1}^k \theta_i \sum_{j=1}^n h_i(y_j; \mu(\mathbf{x}_j^\top \boldsymbol{\beta}), \gamma) \\ &\quad + \sum_{j=1}^n \log f(y_j; \mu(\mathbf{x}_j^\top \boldsymbol{\beta}), \gamma). \end{aligned} \quad (4)$$

96 Note that the last term equals the log-likelihood of the regression model.

For all $\mathbf{x} = \mathbf{x}_j$, $j = 1, \dots, n$, the normalisation constants $C(\mu(\mathbf{x}^\top \boldsymbol{\beta}), \gamma, \boldsymbol{\theta})$ must guarantee that the area under density (3) equals one. Thus for a given set of parameter values, and for sample of n observations, n normalisation constants are required. For some outcome distributions the normalisation constant may not exist, but the type of score tests that we will derive still do exist (Mardia & Kent 1991; Kallenberg, Ledwina & Rafajlowicz 1997). For the orthonormal functions $\{h_i\}$, again for a given set of parameter values, the set of functions

must be calculated for each of the n regressors \mathbf{x}_j . In particular, the orthonormality condition reads as

$$\int_{-\infty}^{+\infty} h_u(y; \mu(\mathbf{x}^\top \boldsymbol{\beta}), \boldsymbol{\gamma}) h_v(y; \mu(\mathbf{x}^\top \boldsymbol{\beta}), \boldsymbol{\gamma}) f(y; \mu(\mathbf{x}^\top \boldsymbol{\beta}), \boldsymbol{\gamma}) dy = \delta_{uv}, \quad (5)$$

97 with $\delta_{uv} = 1$ if $u = v$ and $\delta_{uv} = 0$ otherwise.

98 In what follows the notation will often be simplified by omitting the dependence of
99 C , f , f_k and $\{h_i\}$ on all parameters and regressors. For example, (5) may be written
100 as $E_0(h_u(Y; \mu, \boldsymbol{\gamma}) h_v(Y; \mu, \boldsymbol{\gamma}) | \mathbf{x}) = \delta_{uv}$, in which $E_0(\cdot | \mathbf{x})$ denotes the conditional
101 expectation under the null hypothesis $\boldsymbol{\theta} = \mathbf{0}$. We will also use μ_j as shorthand notation for
102 $\mu(\mathbf{x}_j^\top \boldsymbol{\beta})$.

103 The smooth test requires the score test statistic for testing $\boldsymbol{\theta} = \mathbf{0}$, and hence the score
104 statistic for $\boldsymbol{\theta}$ is required, as well as the information matrix based on the log-likelihood (4)
105 with all expectations evaluated under the null hypothesis and conditional on the regressors. In
106 this setting, the parameters $\boldsymbol{\gamma}$ and $\boldsymbol{\beta}$ are both considered as nuisance parameters and therefore
107 we stack them into a single vector, $\boldsymbol{\eta}^\top = (\boldsymbol{\gamma}^\top, \boldsymbol{\beta}^\top)$. The nuisance parameter only needs to
108 be estimated under the null hypothesis: its maximum likelihood estimator, which is denoted
109 by $\hat{\boldsymbol{\eta}}$, arises from the hypothesised regression model. The following theorem gives the score
110 statistics and the required information matrices. We will use \mathcal{X} to denote the set of of n
111 regressor vectors \mathbf{x}_j in the sample.

112 **Theorem 1.** Score statistics and information matrices. *The score statistic for θ_i in
113 model (4) is given by $U_i = U_i(\boldsymbol{\eta}) = (\partial/\partial\theta_i)l(\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\theta}) = \sum_{j=1}^n h_i(y_j; \mu_j, \boldsymbol{\gamma})$. Let $g'(\mu) =$
114 $(d/d\mu)g(\mu)$. We use the notation $E_k(\cdot)$ and $E_0(\cdot)$ to refer to the expectation w.r.t density
115 functions $f_k(\cdot; \mu, \boldsymbol{\gamma}, \boldsymbol{\theta})$ and $f(\cdot; \mu, \boldsymbol{\gamma})$, respectively. The elements of the information matrix
116 are given by:*

$$\begin{aligned} (\mathbf{I}_{\theta\theta})_{uv} &= -E_k \left(\frac{\partial^2}{\partial\theta_u \partial\theta_v} l | \mathcal{X} \right) \Big|_{\boldsymbol{\theta}=\mathbf{0}} = n\delta_{uv} \\ (\mathbf{I}_{\theta\boldsymbol{\gamma}})_{uv} &= -E_k \left(\frac{\partial^2}{\partial\theta_u \partial\gamma_v} l | \mathcal{X} \right) \Big|_{\boldsymbol{\theta}=\mathbf{0}} = \sum_{j=1}^n E_0 \left(h_u(Y_j; \mu_j) \frac{\partial}{\partial\gamma_v} \log f(Y_j; \mu_j) | \mathbf{x}_j \right) \\ (\mathbf{I}_{\theta\boldsymbol{\beta}})_{uv} &= -E_k \left(\frac{\partial^2}{\partial\theta_u \partial\beta_v} l | \mathcal{X} \right) \Big|_{\boldsymbol{\theta}=\mathbf{0}} = \sum_{j=1}^n E_0 \left(h_u(Y_j; \mu_j) \frac{\partial}{\partial\mu_j} \log f(Y_j; \mu_j) | \mathbf{x}_j \right) \frac{x_{jv}}{g'(\mu_j)}. \end{aligned}$$

117 The theorem does not give expressions for $\mathbf{I}_{\boldsymbol{\beta}\boldsymbol{\beta}}$, $\mathbf{I}_{\boldsymbol{\gamma}\boldsymbol{\gamma}}$ and $\mathbf{I}_{\boldsymbol{\gamma}\boldsymbol{\beta}}$, because they are not
118 affected by our embedding density (3), i.e. they are the information matrices under the
119 hypothesised regression model.

120 For the calculation of the smooth test statistic, the expectations in the expressions of
121 the information matrices of Theorem 1 are evaluated; this depends on the exact form of

122 the regression model, and the nuisance parameters need to be replaced by their maximum
 123 likelihood estimates $\hat{\boldsymbol{\eta}}$, after which these matrices are denoted by $\hat{\mathbf{I}}$.

124 Several examples of regression models will follow in Section 3. The next lemma gives
 125 the smooth test statistic and its limiting null distribution. The proof is a straightforward
 126 application of maximum likelihood theory and is omitted here (see e.g. Boos & Stefanski
 127 (2013) for a good exposition to maximum likelihood theory).

128 **Lemma 1.** Smooth test statistic and its asymptotic null distribution. *Let \mathbf{V} denote the*
 129 *vector $(1/\sqrt{n})(U_1(\boldsymbol{\eta}), \dots, U_k(\boldsymbol{\eta}))^\top$, and let $\hat{\mathbf{V}}$ denote the same vector but with the nuisance*
 130 *parameter $\boldsymbol{\eta}$ replaced with its maximum likelihood estimator $\hat{\boldsymbol{\eta}}$. The smooth test statistic for*
 131 *testing $\boldsymbol{\theta} = \mathbf{0}$ against $\boldsymbol{\theta} \neq \mathbf{0}$ in model (4) is given by $\hat{S}_k = n\hat{\mathbf{V}}^\top \left(\hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\theta}} - \hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\eta}}\hat{\mathbf{I}}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1}\hat{\mathbf{I}}_{\boldsymbol{\eta}\boldsymbol{\theta}}^\top \right)^- \hat{\mathbf{V}}$*
 132 *in which $(\cdot)^-$ denotes a generalised inverse. Let r be the rank of $\hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\theta}} - \hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\eta}}\hat{\mathbf{I}}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1}\hat{\mathbf{I}}_{\boldsymbol{\eta}\boldsymbol{\theta}}^\top$. Given*
 133 *that $r \geq 1$, under the null hypothesis, as $n \rightarrow \infty$, $\hat{S}_k \xrightarrow{d} \chi_r^2$.*

134 Note that the second term in the matrix $\hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\theta}} - \hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\eta}}\hat{\mathbf{I}}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1}\hat{\mathbf{I}}_{\boldsymbol{\eta}\boldsymbol{\theta}}^\top$ corrects the estimated
 135 information matrix of the parameter $\boldsymbol{\theta}$ for the nuisance parameter estimation. For many
 136 regression models some of the elements of $\hat{\mathbf{V}}$ will be exactly zero as a consequence of the
 137 estimation of the model parameters; in Section 3 examples will be given. It is then convenient
 138 to first remove these components from $\hat{\mathbf{V}}$, or, equivalently, remove the corresponding
 139 terms $\theta_i h_i(y; \mu(\mathbf{x}^\top \boldsymbol{\beta}), \boldsymbol{\gamma})$ from model (3). In this case the estimated covariance matrix
 140 $\hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\theta}} - \hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\eta}}\hat{\mathbf{I}}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1}\hat{\mathbf{I}}_{\boldsymbol{\eta}\boldsymbol{\theta}}^\top$ will be of full rank r , which is k minus the number of components removed
 141 from $\hat{\mathbf{V}}$ (or from the model).

142 A single element from $\hat{\mathbf{V}}$, say \hat{V}_i , corresponds the parameter θ_i of the density function
 143 (3) and to score statistic U_i (Theorem 1), and hence it can serve as the basis of a test statistic
 144 for testing $\theta_i = 0$ against $\theta_i \neq 0$. With $\hat{\sigma}_i^2$, the i th diagonal element of $\hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\theta}} - \hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\eta}}\hat{\mathbf{I}}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1}\hat{\mathbf{I}}_{\boldsymbol{\eta}\boldsymbol{\theta}}^\top$, the
 145 hypotheses can be tested upon using the asymptotic null distribution, i.e. $\hat{V}_i/\hat{\sigma}_i \xrightarrow{d} N(0, 1)$,
 146 as $n \rightarrow \infty$. The statistic $\hat{V}_i/\hat{\sigma}_i$ is referred to as the i th component of \hat{S}_k . Examples will follow
 147 later in this paper.

148 We conclude this section with a note on the convergence of the test statistics
 149 to their asymptotic null distributions. Rippon (2012) performed simulation studies
 150 for assessing the empirical type I error rates as a function of the sample size.
 151 Results for the special case of the Poisson regression can be found at [https://ogma.newcastle.edu.au/vital/access/services/Download/uon:](https://ogma.newcastle.edu.au/vital/access/services/Download/uon:12622/ATTACHMENT02?view=true\#page=93)
 152 [12622/ATTACHMENT02?view=true\#page=93](https://ogma.newcastle.edu.au/vital/access/services/Download/uon:12622/ATTACHMENT02?view=true\#page=93). These results demonstrate slow
 154 convergence of the asymptotic approximations. Sample sizes of at least 100 are needed for
 155 good type I error rate control. For normal linear regression models, Peña & Slate (2006)
 156 concluded from their simulation study, that the convergence of test statistics similar to our

157 \hat{V}_3 and \hat{V}_4 is also very slow. The parametric bootstrap, on the other hand, works well (see
 158 simulation studies in Section 4).

159 2.2. Special cases

160 *Generalised linear models (GLM)*

Generalised linear models form a special class of regression models (1) and (2) by restricting the conditional distribution of the outcome variable to the exponential family. Many GLMs belong to a one-parameter exponential family for which the single parameter is related to the conditional mean and hence to the β -parameter through $\mu(\mathbf{x}; \beta) = E(Y | \mathbf{x}) = g^{-1}(\mathbf{x}^\top \beta)$. Examples include logistic regression (binomial distribution) and Poisson regression (Poisson distribution), among others. For this class of models, there is no nuisance parameter γ , and hence some of the information matrices simplify. In particular,

$$(\mathbf{I}_{\theta\beta})_{uv} = -E_0 \left(\frac{\partial^2}{\partial \theta_u \partial \beta_v} l \mid \mathcal{X} \right) = \sum_{j=1}^n \frac{x_{jv}}{\text{var}_0(Y_j \mid \mathbf{x}_j) g'(\mu_j)} E_0(h_u(Y_j; \mu_j)(Y_j - \mu_j)).$$

Furthermore, if the canonical link function is used, $(\mathbf{I}_{\theta\beta})_{uv}$ further simplifies to

$$(\mathbf{I}_{\theta\beta})_{uv} = -E_0 \left(\frac{\partial^2}{\partial \theta_u \partial \beta_v} l \mid \mathcal{X} \right) = \sum_{j=1}^n x_{jv} E_0(h_u(Y_j; \mu_j)(Y_j - \mu_j)).$$

161 The use of the exponential family also allows the use of Iterative Reweighted Least Squares
 162 (IRLS) as a general algorithm for β -parameter estimation.

163 GLMs may include a dispersion parameter. These models assume that the outcome
 164 distribution belongs to the exponential dispersion family. Normal and gamma regression
 165 models belong to this class. Although in most GLM literature the dispersion parameter is
 166 not estimated by maximum likelihood, but rather by the method of moments, we will further
 167 assume that the dispersion parameter is the γ -nuisance parameter. Note that for the normal
 168 distribution the maximum likelihood estimator and the method of moment estimator are
 169 equivalent.

170 *Score functions are linear combinations of the basis functions*

For many distributions the score functions of the nuisance parameters and the mean μ_j can be expressed as a linear combination of the orthonormal basis functions $h_i(y; \mu, \gamma)$,

$i = 1, \dots, k$. In particular,

$$\frac{\partial}{\partial \gamma_v} \log f(y; \mu_j, \gamma) = \sum_{i=1}^k a_{ijv} h_i(y; \mu_j, \gamma) \quad \text{and} \quad \frac{\partial}{\partial \mu_j} \log f(y; \mu_j, \gamma) = \sum_{i=1}^k b_{ij} h_i(y; \mu_j, \gamma), \quad (6)$$

171 for some sets of constants $\{a_{ijv}\}$ and $\{b_{ij}\}$. Often, many of these constants are zero. The
 172 normal and exponential distributions are two examples.

Upon using the orthonormality, the elements of the information matrix now become

$$\begin{aligned} (\mathbf{I}_{\theta\gamma})_{uv} &= \sum_{j=1}^n a_{ujv}, & (\mathbf{I}_{\theta\beta})_{uv} &= \sum_{j=1}^n \frac{x_{jv} b_{uj}}{g'(\mu_j)}, \\ (\mathbf{I}_{\gamma\beta})_{uv} &= \sum_{j=1}^n \frac{x_{jv} (\sum_{i=1}^k a_{iju} b_{ij})}{g'(\mu_j)}, & (\mathbf{I}_{\gamma\gamma})_{uv} &= \sum_{j=1}^n \sum_{i=1}^k a_{iju} a_{ijv}, \\ (\mathbf{I}_{\beta\beta})_{uv} &= \sum_{j=1}^n \frac{x_{ju} x_{jv} \sum_{i=1}^k b_{ij}^2}{(g'(\mu_j))^2}. \end{aligned}$$

173 Note that for GLMs from a one-parameter exponential family, the score function of μ_j
 174 is a first order polynomial in y , proportional to $y - \mu_j$.

175 For several common distributions the orthonormal basis consists of orthonormal
 176 polynomials in $y - \mu_j$. Rayner, Thas & Best (2009, Appendix C) gives explicit forms for
 177 many examples.

178 3. Examples

179 In the sections following, smooth tests for several specific regression models will be
 180 discussed in more detail. All tests are constructed as earlier described. Poisson regression,
 181 which is treated in Section 3.1 is an example of a GLM with no nuisance parameters and with
 182 a score function linearly related to the orthonormal polynomials (Sections 2.2 and 2.2 apply).
 183 Logistic regression belongs to the same type of regression models; some details are given in
 184 Appendix A.2 in the Supporting and Supplementary Material. The normal linear regression
 185 model is discussed in Section 3.2. Again Sections 2.2 and 2.2 apply, but now a nuisance
 186 parameter is present (the error term variance). In Section 3.3 the smooth test for zero-inflated
 187 Poisson (ZIP) regression models are developed. ZIP regression models do not belong to the
 188 class of GLMs, but still our general theory applies.

189 For all these regression models, a numerical example is provided. The p -values are
 190 computed by means of a parametric bootstrap procedure with 2,000 bootstrap runs. In

191 simulation studies presented in Section 4, we will demonstrate that this bootstrap procedure
 192 succeeds in controlling the type I error rate.

193 3.1. Poisson regression

194 *Test statistic*

Poisson regression fits into the GLM framework with a Poisson distribution for the outcome variable and a canonical log-link, i.e. $g(\mu(\mathbf{x}; \boldsymbol{\beta})) = \log(\mu(\mathbf{x}; \boldsymbol{\beta})) = \mathbf{x}^\top \boldsymbol{\beta}$. No nuisance parameters other than the $\boldsymbol{\beta}$ -parameters are involved. The Poisson–Charlier polynomials are known to form an orthonormal basis w.r.t. the Poisson distribution. The polynomials of order one and two are given by $h_1(y; \mu) = (y - \mu)/\sqrt{\mu}$ and $h_2(y; \mu) = [(y - \mu)^2 - y]/(\mu\sqrt{2})$. Higher order polynomials can be found in Rayner, Thas & Best (2009, Appendix C). The score function for the mean is given by $(\partial/\partial\mu) \log f(y; \mu) = (y - \mu)/\mu = h_1(y; \mu)/\sqrt{\mu}$. From this expression, we see that the simplifications of Section 2.2 apply with $b_{1j} = 1/\sqrt{\mu_j}$ and $b_{ij} = 0$ for $i = 2, \dots, k$. Hence, we find

$$(\mathbf{I}_{\theta\beta})_{1v} = \sum_{j=1}^n x_{jv} \sqrt{\mu_j} \quad (\mathbf{I}_{\beta\beta})_{uv} = \sum_{j=1}^n x_{ju} x_{jv} \mu_j \quad \text{and} \quad (\mathbf{I}_{\theta\beta})_{uv} = 0, \quad \text{for } u = 2, \dots, k.$$

195 With these expressions, we find that $\hat{\mathbf{I}}_{\theta\theta} - \hat{\mathbf{I}}_{\theta\eta} \hat{\mathbf{I}}_{\eta\eta}^{-1} \hat{\mathbf{I}}_{\theta\eta}^\top$ is a diagonal matrix with first
 196 element equal to ω^2 (see further for its definition) and all other diagonal elements equal
 197 to one. The diagonal structure results in a decomposition of the order k smooth test
 198 statistic: $\hat{S}_k = \hat{V}_1^2/\omega^2 + \hat{V}_2^2 + \dots + \hat{V}_k^2$, where $\hat{V}_i = \sum_{j=1}^n h_i(y_j; \mathbf{x}_j^\top \hat{\boldsymbol{\beta}})/\sqrt{n}$ and where
 199 ω^2 can be conveniently expressed using matrix notation. Let \mathbf{X} denote the usual $n \times p$
 200 design matrix and \mathbf{D} a diagonal matrix with elements $\sqrt{\mu_1}, \dots, \sqrt{\mu_n}$. With \mathbf{I}_n the $n \times n$
 201 identity matrix and $\mathbf{1}_n$ a column vector with all elements set to one, we write $\omega^2 =$
 202 $\mathbf{1}_n^\top \left(\mathbf{I}_n - \mathbf{D} \mathbf{X} (\mathbf{X}^\top \mathbf{D}^2 \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{D} \right) \mathbf{1}_n / n$.

We discuss the first few components in some detail. The numerator of the first component is given by the square of

$$\hat{V}_1 = \frac{1}{\sqrt{n}} \sum_{j=1}^n h_1(y_j; \mathbf{x}_j^\top \hat{\boldsymbol{\beta}}) = \frac{1}{\sqrt{n}} \sum_{j=1}^n \frac{y_j - \hat{\mu}_j}{\sqrt{\hat{\mu}_j}} = \frac{1}{\sqrt{n}} \sum_{j=1}^n \frac{y_j - \exp(\mathbf{x}_j^\top \hat{\boldsymbol{\beta}})}{\sqrt{\exp(\mathbf{x}_j^\top \hat{\boldsymbol{\beta}})}},$$

203 in which $y_j - \hat{\mu}_j$ is the residual of the j th observation. The denominator $\sqrt{\hat{\mu}_j}$ is the standard
 204 error of the residual if the true $\boldsymbol{\beta}$ would have been used instead of its MLE $\hat{\boldsymbol{\beta}}$. To correct
 205 for the estimation of $\boldsymbol{\beta}$, the factor $1/\omega^2$ appears in the first component. Hence, \hat{V}_1^2/ω^2 can
 206 be interpreted as a goodness of fit statistic for the specification of the mean model, which
 207 includes both the specification of the linear predictor $\eta = \mathbf{x}^\top \boldsymbol{\beta}$ and the link function $g(\cdot)$.

The second component can be written as the square of

$$\hat{V}_2 = \frac{1}{\sqrt{n}} \sum_{j=1}^n h_2(y_j; \mathbf{x}_j^\top \hat{\boldsymbol{\beta}}, \hat{\sigma}^2) = \frac{1}{\sqrt{n}} \sum_{j=1}^n \frac{(y_j - \hat{\mu}_j)^2 - y_j}{\sqrt{2}\hat{\mu}_j} = \frac{1}{\sqrt{n}} \sum_{j=1}^n \frac{(y_j - \exp(\mathbf{x}_j^\top \hat{\boldsymbol{\beta}}))^2 - y_j}{\sqrt{2} \exp(\mathbf{x}_j^\top \hat{\boldsymbol{\beta}})}.$$

208 In the numerator we interpret $(y_j - \hat{\mu}_j)^2 - y_j$ as the residual of the j th observation w.r.t.
 209 the model-based specification of the variance of the outcome. In particular, the conditional
 210 expectation of $(Y_j - \hat{\mu}_j)^2$, given \mathbf{x}_j , is, by definition, the conditional variance of the
 211 outcome Y_j , and the conditional expectation of Y_j given \mathbf{x}_j trivially is the conditional
 212 mean of the outcome. Hence the conditional expectation of the residual is a contrast
 213 between the conditional variance and mean, which, if the Poisson regression model holds
 214 true, should be equal. Hence, the second component can be interpreted as a statistic
 215 measuring the goodness of fit of the second moment of the Poisson regression model.
 216 The higher order components $\hat{V}_3^2, \dots, \hat{V}_k^2$ are interpreted in a similar fashion. Finally, we
 217 note that \hat{V}_2 is closely related to a test statistic proposed by Dean & Lawless (1989) for
 218 detecting overdispersion in Poisson regression. Their test statistic is given by $\sum_{j=1}^n [(y_i -$
 219 $\exp(\mathbf{x}_j^\top \hat{\boldsymbol{\beta}}))^2 - y_i] / \sqrt{2n \sum_{j=1}^n [\exp(\mathbf{x}_j^\top \hat{\boldsymbol{\beta}})]^2}$.

220 **Numerical example**

221 Spinelli, Lockhart & Stephens (2002) presented an example in which the expected
 222 frequency of cases of bladder cancer in male aluminium workers is analysed with a Poisson
 223 regression model with age and exposure to coal tar pitch volatiles as regressors. The age is
 224 included as a factor variable with 11 levels, referring to age groups, each spanning five years.
 225 The exposure is included as a continuous regressor, but it is actually an ordinal variable
 226 taking four values. The model also includes an offset, which is set to the logarithm of the
 227 total person years at risk. The dataset includes 4213 workers. One of the conclusions from
 228 the data analysis is that the exposure has a significant effect ($p = 0.00184$) on the number of
 229 bladder cancer cases per person year, correcting for age. The effect is estimated as a factor
 230 2.089 increase in the expected number of bladder cancer cases per person year, when the
 231 exposure level increases with one level.

232 The smooth test of order $k = 4$ has been applied to this example. The p -values were
 233 computed using the parametric bootstrap with 2,000 bootstrap samples. Table 1 shows the
 234 results. The overall order $k = 4$ test gives $p = 0.632$, and hence at the 5% level of significance
 235 there is no evidence against the Poisson assumption. Neither do any of the component tests
 236 suggest any deviation from the Poisson assumption.

Table 1. Results of the order $k = 4$ smooth test applied to the bladder cancer example. The two-sided p -values are computed from 2,000 parametric bootstrap runs (denoted by p_B) and from the asymptotic distribution (p_A). Results are shown for two models: “Age Factor” refers to the model with age included as a factor variable, and “Age Ordinal” refers to the model with age as an ordinal regressor.

Statistic	Age Factor			Age Ordinal		
	value	p_A	p_B	value	p_A	p_B
\hat{S}_4	0.631	0.960	0.653	1.306	0.860	0.798
\hat{V}_1/ω	0.123	0.902	0.951	1.025	0.305	0.020
\hat{V}_2	-0.733	0.464	0.842	-0.502	0.616	0.390
\hat{V}_3	0.280	0.779	0.958	0.013	0.999	0.694
\hat{V}_4	0.003	0.998	0.998	0.049	0.961	0.962

237 Table 1 also shows the results of the smooth test applied to a Poisson regression model
 238 in which age is not included as a factor variable, but as a continuous regressor which takes
 239 the values 1 up to 11, referring to the 11 age classes (age is here thus an ordinal variable). The
 240 first component test gives a two-sided p -value of 0.020, and hence it suggests that the mean
 241 model is not correctly formulated. Although the results are not presented here, we generally
 242 also advise looking at the conventional residual plots for assessing the correctness of the mean
 243 model.

244 Finally, Table 1 also shows the p -values calculated from the asymptotic null
 245 distributions. Although the conclusions at the 5% level of significance are almost the
 246 same, these results illustrate the discrepancy between the bootstrap and the asymptotic
 247 approximation for the sample size of this example ($n = 44$).

248 3.2. Linear regression with normal error terms

249 **Test statistic**

250 For the normal linear regression model, the link function is the identity function, and the
 251 residual variance, say σ^2 , is the only nuisance parameter in the normal distribution; thus $\gamma =$
 252 σ^2 . Hence, the conditional mean is written as $\mu(\mathbf{x}^\top \boldsymbol{\beta}) = \mathbf{x}^\top \boldsymbol{\beta}$. The system of orthonormal
 253 polynomials is given by the Hermite polynomials (Rayner, Thas & Best 2009, Appendix C),
 254 of which the first few are given by $h_0(y; \mathbf{x}^\top \boldsymbol{\beta}, \sigma^2) = 1$, $h_1(y; \mathbf{x}^\top \boldsymbol{\beta}, \sigma^2) = (y - \mathbf{x}^\top \boldsymbol{\beta})/\sigma$
 255 and $h_2(y; \mathbf{x}^\top \boldsymbol{\beta}, \sigma^2) = [(y - \mathbf{x}^\top \boldsymbol{\beta})^2 - \sigma^2]/(\sigma^2 \sqrt{2})$.

For the normal regression model, the score functions for $\boldsymbol{\beta}$ and σ^2 can be written as
 ($u = 1, \dots, p$)

$$\frac{\partial}{\partial \beta_u} \log f(y; \mu_j, \sigma^2) = x_{ju}(y_j - \mathbf{x}_j^\top \boldsymbol{\beta}) \quad \frac{\partial}{\partial \sigma^2} \log f(y; \mu_j, \sigma^2) = [(y_j - \mathbf{x}_j^\top \boldsymbol{\beta})^2 - \sigma^2].$$

Note that the right hand sides of the equations involve the Hermite polynomials h_1 and h_2 and that the score functions take the form of (6) and hence the simplifications of Section 2.2 apply. Also note that the maximum likelihood estimator (MLE) of σ^2 is the solution of $\sum_{j=1}^n (\partial/\partial\sigma^2) \log f(y_j; \mu_j, \sigma^2) = \sigma^2 \sum_{j=1}^n h_2(y_j; \mu_j; \sigma^2) = 0$. This implies that the second component statistic $\hat{V}_2 = S_2(\hat{\beta}, \hat{\sigma}^2)/\sqrt{n} = 0$. The interpretation is that the smooth test cannot detect a wrongly specified variance, because the variance is estimated by matching the variance parameter σ^2 to the empirical variance (up to the asymptotically negligible factor $n/(n-1)$). We therefore continue with the construction of the smooth test, with the second orthonormal Hermite polynomial removed from model (3). Upon applying the methods described in this paper, including the simplifications of Section 2.2, we find again, as for Poisson regression, that $\hat{\mathbf{I}}_{\theta\theta} - \hat{\mathbf{I}}_{\theta\eta} \hat{\mathbf{I}}_{\eta\eta}^{-1} \hat{\mathbf{I}}_{\eta\theta}^\top$ is the identity matrix with the first element replaced by $\omega^2 = \frac{1}{n} \mathbf{1}_n (\mathbf{I}_n - \mathbf{H}) \mathbf{1}_n^\top$, in which $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ is the hat matrix of the linear regression model. Thanks to the diagonal structure of the matrix, the order k smooth test statistic becomes $\hat{S}_k = \hat{V}_1^2/\omega^2 + \hat{V}_3^2 + \dots + \hat{V}_k^2$, with $\hat{V}_i = \sum_{j=1}^n h_i(y_j; \mathbf{x}_j^\top \hat{\beta}, \hat{\sigma}^2)/\sqrt{n}$ ($i = 1, 3, 4, \dots, k$). For $i = 1$, we write

$$\hat{V}_1 = \frac{1}{\sqrt{n}} \sum_{j=1}^n h_1(y_j; \mathbf{x}_j^\top \hat{\beta}, \hat{\sigma}^2) = \frac{1}{\sqrt{n}} \sum_{j=1}^n \frac{y_j - \mathbf{x}_j^\top \hat{\beta}}{\hat{\sigma}} = \frac{1}{\hat{\sigma}\sqrt{n}} \sum_{j=1}^n (y_j - \mathbf{x}_j^\top \hat{\beta}) = 0,$$

256 in which the equality to zero is a consequence of the residuals summing to zero when ML
 257 or least-squares was used for parameter estimation in a linear regression model that includes
 258 an intercept. Hence, for such regression models the first order Hermite polynomial is also
 259 removed from model (3). The final test statistic is then given by $\hat{S}_k = \hat{V}_3^2 + \dots + \hat{V}_k^2$.

260 The test statistic thus shows a natural decomposition into $k-2$ components. Because
 261 of the polynomial nature of the orthonormal functions, the components can be roughly
 262 interpreted in terms moments; see Henze & Klar (1996); Henze (1997); Klar (2000); Thas
 263 (2010) for detailed discussions on this issue. For example, a large \hat{V}_3^2 is an indication that
 264 the skewness of the true outcome distribution does not agree with the skewness of the
 265 hypothesised normal outcome distribution. Since the latter is zero (symmetric distribution), a
 266 large \hat{V}_3^2 suggests that the true outcome distribution is skewed. Similarly, a large \hat{V}_4^2 suggests
 267 that the kurtosises of the true distribution and the hypothesised normal distribution do not
 268 agree. The two squared components (\hat{V}_3^2 and \hat{V}_4^2) are equivalent to the components \hat{S}_3^2 and
 269 \hat{S}_4^2 of Peña & Slate (2006).

270 **Numerical example**

271 Davison (2003, example 8.25) reports data on an experiment in which 48 animals were
 272 randomly allocated to 12 groups of four animals. Each group was given one of three poisons

Table 2. Results of the order $k = 4$ smooth test applied to the poison data. Smooth test results for the normal distributional assumption are shown. The two-sided p -values are computed from the asymptotic distribution (p_A) and from 2,000 parametric bootstrap runs (denoted by p_B). Results for the ANOVA models with the survival time and with the reciprocal survival time are shown.

Statistic	survival time			reciprocal survival time		
	value	p_A	p_B	value	p_A	p_B
\hat{S}_4	19.4546	0.001	0.008	1.1467	0.887	0.874
\hat{V}_3	3.2628	0.001	<0.001	0.9664	0.334	0.370
\hat{V}_4	2.9679	0.003	0.012	-0.4612	0.645	0.760

273 and one of four treatments, resulting in a balanced design. The outcome is the survival time
 274 in 10-hour units. An analysis based on an additive two-factor analysis of variance (ANOVA)
 275 model (i.e. regression model with dummies coding for the two factors) revealed that both
 276 poison and treatment have significant effects on the mean outcome. This dataset will be
 277 referred to as the poison data.

278 Table 2 shows the results from the smooth test of order $k = 4$, and its component tests. At
 279 the 5% level of significance we can conclude that the normality assumption is not satisfied.
 280 Both the third and fourth order component tests give highly significant results, suggesting
 281 that perhaps within each or some poison/treatment groups the outcome shows a skewed
 282 distribution with too heavy or too light tails. Given the rather small samples sizes in each
 283 group (4 animals), our analysis gives a warning that the ANOVA p -values may not be trusted.
 284 Davison (2003) suggested applying a Box–Cox transformation to the outcome to resolve the
 285 issue.

286 We have also applied the smooth tests to the same model, but with the reciprocal survival
 287 time as outcome. The results are also presented in Table 2. Now no significant goodness of
 288 fit is observed.

289 Finally, Table 2 also shows the p -values calculated from the asymptotic null
 290 distributions. Once more the results illustrate some disagreement between the bootstrap
 291 and the asymptotic approximation for the sample size of this example ($n = 48$), but the
 292 differences are not as large as in the previous example.

293 3.3. Zero inflated Poisson regression

294 *Test statistic*

295 The Zero Inflated Poisson (ZIP) distribution is a mixture distribution of a Poisson
 296 distribution and a point probability at zero. It thus allows for an excess of zeroes as compared
 297 to what is expected under a Poisson distribution. The probability of an excess zero is

298 quantified through an additional parameter, which is considered here as a nuisance parameter.
 299 The ZIP distribution does not belong to the exponential family, and hence ZIP regression
 300 models are not within the class of GLMs. Neither does the simplification of Section 2.2
 301 apply.

302 The mean of the Poisson component is related to the linear predictor $\mathbf{x}^\top \boldsymbol{\beta}$ through
 303 the log-link, i.e. $g(\mu(\mathbf{x}; \boldsymbol{\beta})) = \log(\mu(\mathbf{x}; \boldsymbol{\beta})) = \mathbf{x}^\top \boldsymbol{\beta}$. The probability of an excess zero is
 304 denoted by the parameter γ , which acts as a nuisance parameter. The score functions for γ
 305 and μ are given by

$$\begin{aligned} \frac{\partial}{\partial \mu} \log f(y; \mu, \gamma) &= \frac{y - (1 - \gamma)\mu}{\mu} - \gamma \left(1 - \frac{\delta_0(y)}{\gamma + (1 - \gamma) \exp(-\mu)} \right) \\ \frac{\partial}{\partial \gamma} \log f(y; \mu, \gamma) &= \frac{1}{1 - \gamma} \left(\frac{\delta_0(y)}{\gamma + (1 - \gamma) \exp(-\mu)} - 1 \right), \end{aligned}$$

306 where $\delta_0(y) = 1$ if $y = 0$ and $\delta_0(y) = 0$ otherwise. With these score functions, the MLEs of
 307 $\boldsymbol{\beta}$ and γ can be obtained (e.g. iterative estimation scheme).

308 The polynomial of order one is given by $h_1(y; \mu, \gamma) = [y - (1 -$
 309 $\gamma)\mu] / \sqrt{(1 - \gamma)\mu + \gamma(1 - \gamma)\mu^2}$. Higher order orthonormal polynomials can be computed
 310 using the Emerson recursion relation (see e.g. Rayner, Thas & De Boeck (2008)).
 311 Polynomials up to order 4 are explicitly given in Appendix C of Rayner, Thas & Best (2009).
 312 With these polynomials the score functions cannot be written in the form of (6), and hence
 313 the simplification of Section 2.2 does not apply here. Therefore we need all the elements
 314 of the information matrix as given in Lemma 1. These elements require the following
 315 expectations

$$\begin{aligned} E_0 \left(h_u(Y_j; \mu_j) \frac{\partial}{\partial \gamma} \log f(Y_j; \mu_j, \gamma) \mid \mathbf{x}_j \right) &= \frac{h_u(0; \mu_j, \gamma)}{1 - \gamma} \\ E_0 \left(h_u(Y_j; \mu_j) \frac{\partial}{\partial \mu_j} \log f(Y_j; \mu_j, \gamma) \mid \mathbf{x}_j \right) &= \delta_1(u) \sqrt{(1 - \gamma)/\mu_j + \gamma(1 - \gamma)} + h_u(0; \mu_j, \gamma) \gamma \\ E_0 \left(\frac{\partial}{\partial \gamma} \log f(Y_j; \mu_j, \gamma) \frac{\partial}{\partial \mu_j} \log f(Y_j; \mu_j) \mid \mathbf{x}_j \right) &= \frac{1}{1 - \gamma} \left(\frac{\gamma}{\gamma + (1 - \gamma) \exp(-\mu_j)} - 1 \right). \end{aligned}$$

316 With these expressions, and with the parameters replaced with their MLEs, the matrix
 317 $\hat{\mathbf{I}}_{\theta\theta} - \hat{\mathbf{I}}_{\theta\eta} \hat{\mathbf{I}}_{\eta\eta}^{-1} \hat{\mathbf{I}}_{\theta\eta}^\top$ can be calculated. In contrast to the two previous examples, this matrix
 318 is not diagonal, and thus the smooth test statistic \hat{S}_k cannot be written as the sum of its
 319 components.

320 **Numerical example**

321 Kostic et al. (2015) investigated the gut microbiome of 33 infants who were genetically
 322 predisposed to develop type I diabetes (T1D). The infants were followed during 3 to 4

323 years in a longitudinal study. At regular visits stool samples were taken for microbiome
324 analysis through 16S rRNA sequencing, resulting in abundance data for 2,239 OTUs (OTU
325 = operational taxonomic unit, which is a proxy for a microorganism species identification).
326 Here we consider only the data of the last visit of each of the 33 infants; we also know the
327 age of the child at this last visit (in days) and whether the child was diagnosed with T1D or
328 not. One of the original research questions in this study was to test for differential abundance
329 of microbial species between T1D cases and the healthy infants, while correcting for age.
330 This should be tested for each OTU separately. The microbiome data come as counts from
331 the sequencing technology. For each biological sample, the sum of the counts of all OTUs
332 is known as the library size, which varies substantially between the samples and which is
333 considered to be an irrelevant technical artefact. A typical data analysis starts with assuming
334 a count distribution, and modelling the log-transformed mean parameter of this distribution
335 as $\log(\text{lib. size}) + \beta_0 + \beta_1 \text{T1D} + \beta_2 \text{AGE}$, with $\log(\text{lib. size})$ an offset, T1D the 0/1 disease
336 indicator, and AGE the age of the child. Several count distributions have been proposed for
337 OTU count data: ZIP, negative binomial and zero inflated negative binomial (see e.g. Xu
338 et al. (2015)). Here we test the ZIP distributional assumption in the regression model. We
339 only present the results for two OTUs. The data are shown in Table 1 in Supporting and
340 Supplementary Material.

341 As for the Poisson regression example, we use a smooth test of order $k = 4$, and p -values
342 were computed based on 2,000 parametric bootstrap runs. Results are shown in Table 3. The
343 table also shows the results for testing the Poisson distribution with the smooth test of Section
344 3.1.

345 For the OTU 195929 data the Poisson model is problematic because the V_3 component is
346 significant at the 5% level. However for the ZIP model all components and S_4 have p -values
347 greater than 0.05 and we can conclude this is an acceptable model. For the OTU 1954177
348 data both the Poisson and ZIP models have two components significant at the 0.05 level and
349 neither model is acceptable.

350 It is worth noting that while all components have the same asymptotic null distributions
351 in all models, in small samples their null distributions no longer coincide. A similar comment
352 applies to S_4 . Thus for OTU 1954177, the S_4 p -values for the ZIP and the Poisson models of
353 0.006 and 0.088, respectively, do not necessarily indicate that the Poisson is a more acceptable
354 model than the ZIP. The log likelihood at the MLEs for the Poisson is -118.37 with three
355 degrees of freedom while that for the ZIP is -99.32 with four degrees of freedom. Since the
356 class of ZIP models includes the class of Poisson models the favoured ZIP model will always
357 be at least as good a model as the favoured Poisson model.

358 As for the two previous examples, the p -values based on the asymptotic null distributions
359 are also reported in Table 3. We observe a strong deviation between the asymptotic and

Table 3. Results of the order $k = 4$ smooth test applied to two OTUs of the infant gut microbiome example. Smooth tests for the ZIP and the Poisson distributional assumption are shown. The two-sided p -values are computed from 2,000 parametric bootstrap runs (denoted by p_B) and from the asymptotic distribution (p_A). For the ZIP, the parameter ω equals 1.

Statistic	OTU 194177					
	Poisson			ZIP		
	value	p_A	p_B	value	p_A	p_B
\hat{S}_4	2782.015	<0.001	0.088	101.941	<0.001	0.006
\hat{V}_1/ω	3.924	<0.001	<0.001	2.256	0.024	0.470
\hat{V}_2	17.894	<0.001	<0.001	9.185	<0.001	<0.001
\hat{V}_3	16.748	<0.001	0.240	6.285	<0.001	0.001
\hat{V}_4	46.538	<0.001	0.212	1.333	0.183	0.118
Statistic	OTU 195929					
	Poisson			ZIP		
	value	p_A	p_B	value	p_A	p_B
\hat{S}_4	63.776	<0.001	0.489	13.301	0.010	0.544
\hat{V}_1/ω	-0.434	0.664	0.257	1.996	0.046	0.541
\hat{V}_2	7.583	<0.001	0.753	3.246	0.001	0.384
\hat{V}_3	-2.063	0.039	0.014	0.426	0.670	0.544
\hat{V}_4	1.348	0.178	0.311	-1.066	0.286	0.173

360 bootstrap p -values, which often results in opposite conclusions at the 5% level of significance.

361 Recall that the asymptotic approximation for the data analysis presented in Table 2 was better.

362 The results of all three data examples make us conclude that the discrepancies will
 363 vary with both the model and the sample size. In the next section we will empirically
 364 demonstrate that the bootstrap succeeds in controlling the type I error rate and is hence to
 365 be recommended.

366

4. Simulation study

367 The type I error rate and power of the order $k = 4$ smooth test and its component tests
 368 are evaluated in a simulation study. By no means do we intend to present a comprehensive
 369 simulation study that includes smooth tests for many different distributional assumptions.
 370 Instead we only show the results for Poisson and ZIP regression models for illustrative
 371 purposes.

372 All results are based on 2,000 Monte Carlo runs, and p -values are computed from 200
 373 parametric bootstrap runs. All tests are performed at the 5% level of significance.

374 4.1. Poisson regression

375 In each Monte Carlo simulation run, we simulated $n = 15$ observations from a negative
 376 binomial (NB) regression model with $\log \mu(x) = 2.6 + 2x$ with x taking values 0, 0.5
 377 and 1, each for $n/3$ of the simulated observations. In each simulation run, a Poisson
 378 model with mean model $\log \mu(x; \beta) = \beta_0 + \beta_1 x$ is fitted to the data (i.e. no mean-model
 379 misspecification). The variance of a negative binomial distribution with mean μ is given by
 380 $\mu + \tau\mu^2$, with τ the overdispersion parameter. With $\tau = 0$, the NB collapses to a Poisson
 381 distribution. For a range of values for τ , the results are shown in the top panel of Figure
 382 1. The graph shows that all bootstrap tests control the type I error rate. As expected, the
 383 overdispersion is best detected with the second order component test (V_2), but the power of
 384 the order 4 smooth test (S_4) is not much less.

385 In a second set of simulations, $n = 25$ observations are simulated with a Poisson
 386 distribution with $\log \mu(x) = 1 + 3x_1 + \zeta x_2$ with (x_1, x_2) taken values in the 5×5 grid
 387 pattern generated with $x_1 \in \{-1, -0.5, 0.5, 1\}$ and $x_2 \in \{-1.2, -0.7, -0.2, 0.3, 0.8\}$, with
 388 $n/25$ observations in each point, and with $\zeta \in [0, 0.5]$. The generated data, however, are
 389 analysed with a misspecified Poisson model with only one regressor (x_1). The results are
 390 shown in the bottom panel of Figure 1. The bootstrap tests control the type I error rate.
 391 The order 4 smooth test has good power. However, despite the first moment of the model
 392 being misspecified, it is the second order component test (V_2) that gives the largest power.
 393 The missing regressor in the model causes the data to appear overdispersed. Hilbe (2011)
 394 called this kind of situation ‘apparent overdispersion’ to distinguish it from cases where the
 395 distributional assumption really is violated. This suggests that one should apply goodness of
 396 fit tests always in combination with other diagnostic tools for assessing the correctness of the
 397 mean model (e.g. Pearson or deviance residual plots).

398 4.2. ZIP regression

399 Here we test the null hypothesis of a ZIP outcome distribution with $\log \mu(x) = 1 + \beta x$.
 400 Two simulation scenarios are considered.

401 In the first case we sampled from a zero inflated negative binomial distribution with a
 402 probability of 0.2 for zero-inflation, and $\log \mu(x) = 1 + x$ in which x is standard normally
 403 distributed. The conditional variance of the outcome is given by $\mu + \tau\mu^2$, with τ the
 404 overdispersion parameter of the NB. The parameter τ takes the values $\tau = 0.0, 0.2, \dots, 2.0$.

405 In the second case we simulated data from a ZIP-regression model with $\log \mu(x) = 1 +$
 406 $x_1 + \zeta x_2$ in which x_1 and x_2 are standard normally distributed and $\zeta = 0.0, 0.1, 0.2, \dots, 1.2$.
 407 For $\zeta \neq 0$ the hypothesised ZIP-regression model has a misspecified mean model.

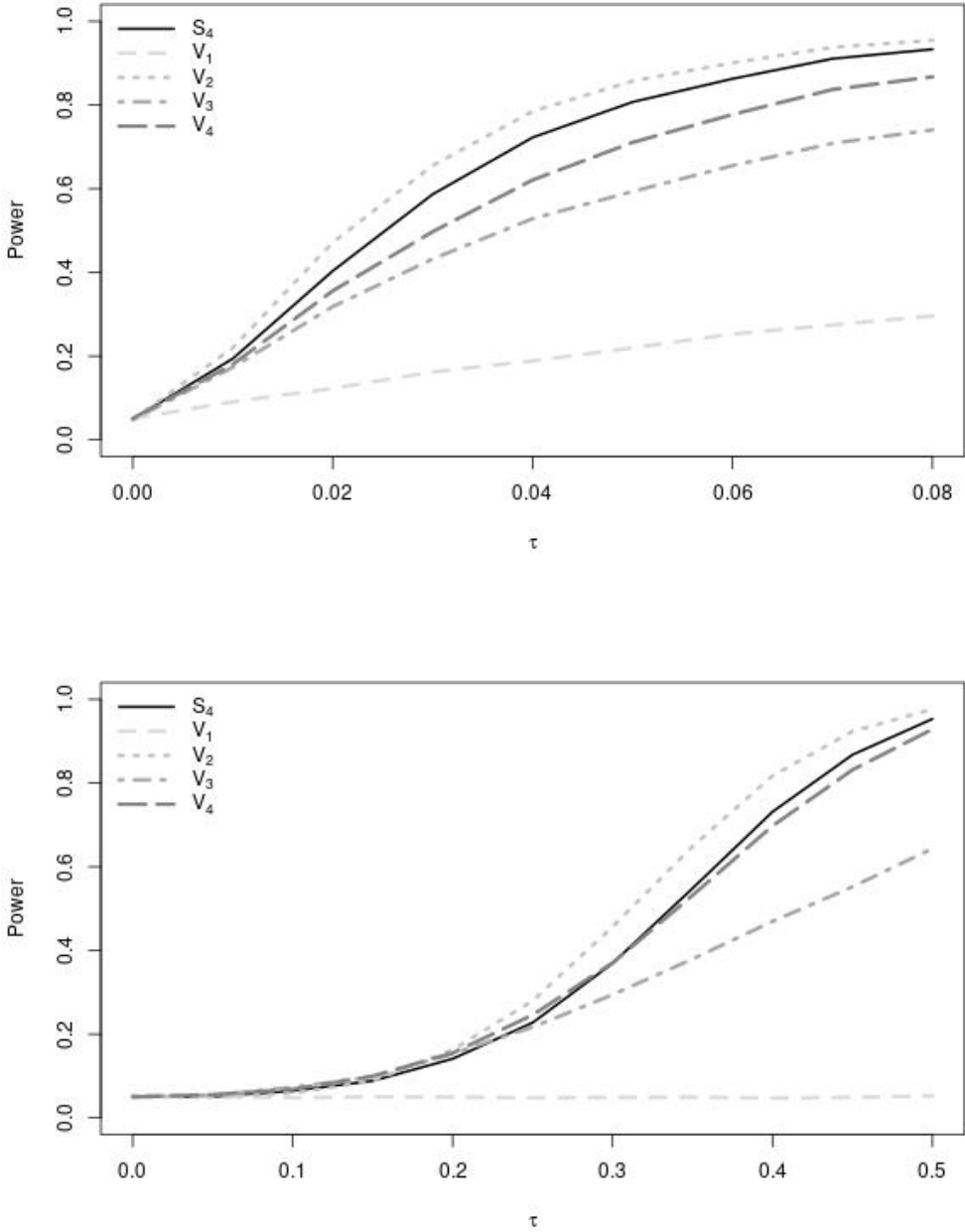


Figure 1. Powers of the smooth test and its component tests for the Poisson regression model. All results are based on 2,000 Monte Carlo simulation runs and 200 parametric bootstrap runs. Top: data generated with negative binomial regression model with overdispersion parameter τ . Bottom: data generated with Poisson regression model with a misspecified mean model that includes an additional term ζx_2 .

408 From the results (Figure 2) it can be seen that V_1 is not very useful for detecting the
409 considered alternatives, but the other four tests detect the alternatives well. The order 4
410 smooth test (S_4) appears most powerful and can be generally recommended. All tests control
411 the type I error rate.

412

5. Conclusion

413 Many regression models contain a distributional assumption for the outcome,
414 conditional on the regressors, allowing for maximum likelihood estimation of the regression
415 parameters. An important class of such models is formed by Generalised Linear Models
416 (GLM). We have developed smooth goodness of fit tests for testing this distributional
417 assumption. The construction starts from sets of polynomials that are orthonormal to
418 the conditional outcome distribution, and thus for each observed regressor another set of
419 orthonormal functions must be computed. The smooth test statistic is a score test statistic
420 developed in a larger class of distributions that embeds the hypothesised regression model.
421 Our methods correctly account for the estimation of nuisance parameters (i.e. regression
422 parameters and other parameters of the conditional outcome distribution). The test statistic
423 asymptotically has a chi-squared null distribution, but simulation studies have shown a slow
424 convergence to the limiting distribution. Therefore it is suggested the parametric bootstrap
425 should be used for p -value calculation. The test statistic is build up from components each
426 of which can also be used as a test statistic. These tests can be helpful in understanding in
427 what sense the data deviate from the hypothesised distribution in terms of e.g. skewness and
428 kurtosis. Even when the data analyst does not want to use formal hypothesis testing for the
429 assessment of the distributional assumption, these components are informative and may be
430 helpful in exploring potential deviations from the hypothesised distribution.

431 The first component is based on a contrast between the observed outcomes and the fitted
432 mean model. From this perspective it may seem like a statistic for assessing the correct
433 specification of the mean model, but when higher order moments are misspecified by the
434 distributional assumption, its diagnostic property may be lost; see e.g. Klar (2000) for a
435 similar issue with the one-sample smooth tests. Also, a misspecified mean model may cause
436 the higher order components to give small p -values, even when the distributional assumption
437 holds true. The ‘apparent overdispersion’ (Hilbe 2011) in the Poisson regression example is
438 an illustration of this issue. As a consequence, we always advise not blindly applying the
439 smooth tests proposed in this paper, but to always complement them with other diagnostic
440 tools for assessing the correctness of the mean model (e.g. Pearson or deviance residual plots).

441 When the regression model is a GLM or when the score functions of its nuisance
442 parameters are linear combinations of the orthonormal basis functions, the construction of the

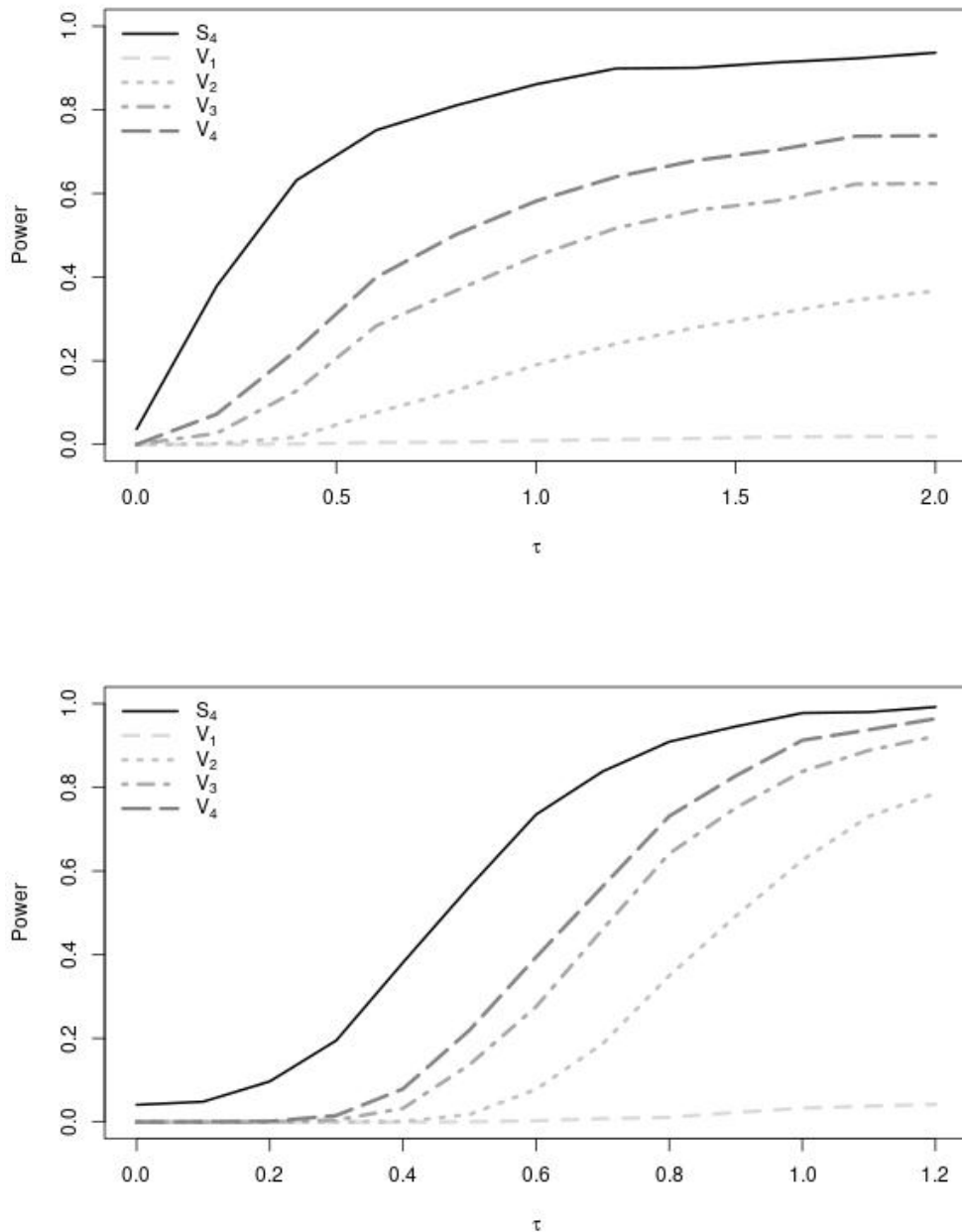


Figure 2. Powers of the smooth test and its component tests for the ZIP regression model. All results are based on 2000 Monte Carlo simulation runs and 200 parametric bootstrap runs. Top: data generated with zero inflated negative binomial regression model with overdispersion parameter τ . Bottom: data generated with ZIP regression model with a misspecified mean model that includes an additional term $\zeta_2 x_2$.

443 smooth test permits simpler expressions. We have given explicit forms of the test statistics
444 for Poisson regression, normal linear regression, zero-inflated Poisson (ZIP) regression and
445 logistic regression. The simulation study, which was restricted to Poisson and ZIP regression,
446 demonstrates that the new tests control the type I error rate when the parametric bootstrap
447 is used. In the PhD thesis of Rippon (2012) smooth tests for more regression models are
448 presented and evaluated in simulation studies. All tests have power for interesting alternative
449 hypotheses. In conclusion, the new tests are customised, focussed and comprehensive.

References

450

- 451 ANSCOMBE, F.J. & TUKEY, J.W. (1963). The examination and analysis of residuals. *Technometrics* **5**,
452 141–160.
- 453 BOOS, D.D. & STEFANSKI, L.A. (2013). *Essential Statistical Inference: Theory and Methods*, vol. 120.
454 New-York: Springer Science & Business Media.
- 455 DAVISON, A. (2003). *Statistical Models*. Cambridge: Cambridge University Press.
- 456 DEAN, C. & LAWLESS, J.F. (1989). Tests for detecting overdispersion in poisson regression models. *Journal*
457 *of the American Statistical Association* **84**, 467–472.
- 458 HART, J. (2013). *Nonparametric smoothing and lack-of-fit tests*. Springer Science & Business Media.
- 459 HENZE, N. (1997). Do components of smooth tests of fit have diagnostic properties? *Metrika* **45**, 121–130.
- 460 HENZE, N. & KLAR, B. (1996). Properly rescaled components of smooth tests of fit are diagnostic.
461 *Australian Journal of Statistics* **38**, 61–74.
- 462 HILBE, J.M. (2011). *Negative Binomial Regression*. Cambridge: Cambridge University Press.
- 463 HUANG, A. (2014). Joint estimation of the mean and error distribution in generalized linear models. *Journal*
464 *of the American Statistical Association* **109**, 186–196.
- 465 HUANG, A. & RATHOUZ, P.J. (2017). Orthogonality of the mean and error distribution in generalized linear
466 models. *Communications in Statistics-Theory and Methods* **46**, 3290–3296.
- 467 KALLENBERG, W., LEDWINA, T. & RAFAJLOWICZ, E. (1997). Testing bivariate independence and
468 normality. *Sankhya, Series A* **59**, 42–59.
- 469 KAUERMAN, G. & CARROLL, R.J. (2001). A note on the efficiency of sandwich covariance matrix
470 estimation. *Journal of the American Statistical Association* **96**, 1387–1396.
- 471 KHMALADZE, E.V., KOUL, H.L. et al. (2004). Martingale transforms goodness-of-fit tests in regression
472 models. *The Annals of Statistics* **32**, 995–1034.
- 473 KLAR, B. (2000). Diagnostic smooth tests of fit. *Metrika* **52**, 237–252.
- 474 KOSTIC, A.D., GEVERS, D., SILJANDER, H., VATANEN, T., HYÖTYLÄINEN, T., HÄMÄLÄINEN, A.M.,
475 PEET, A., TILLMANN, V., PÖHÖ, P., MATTILA, I. et al. (2015). The dynamics of the human infant
476 gut microbiome in development and in progression toward type 1 diabetes. *Cell host & microbe* **17**,
477 260–273.
- 478 LOVE, M.I., HUBER, W. & ANDERS, S. (2014). Moderated estimation of fold change and dispersion for
479 rna-seq data with *deseq2*. *Genome biology* **15**, 550.
- 480 MARDIA, K. & KENT, J. (1991). Rao score tests for goodness-of-fit and independence. *Biometrika* **78**,
481 355–363.
- 482 MCMURDIE, P.J. & HOLMES, S. (2014). Waste not, want not: why rarefying microbiome data is
483 inadmissible. *PLoS computational biology* **10**, e1003531.
- 484 NELDER, J. & WEDDERBURN, R. (1972). Generalized linear models. *Journal of the Royal Statistical*
485 *Society, Series A* **135**, 370–384.
- 486 PEÑA, E.A. & SLATE, E.H. (2006). Global validation of linear model assumptions. *Journal of the American*
487 *Statistical Association* **101**, 341–354.
- 488 RAYNER, J., THAS, O. & DE BOECK, B. (2008). A generalised Emerson recurrence relation. *Australian*
489 *and New Zealand Journal of Statistics* **50**, 235240.
- 490 RAYNER, J.C.W., THAS, O. & BEST, D.J. (2009). *Smooth Tests of Goodness of Fit: Using R*. Singapore:
491 Wiley.
- 492 RIPPON, P. (2012). Application of smooth tests of goodness of fit to generalized linear models. Ph.D. thesis,
493 University of Newcastle, Newcastle, Australia.
- 494 RISSO, D., PERRAUDEAU, F., GRIBKOVA, S., DUDOIT, S. & VERT, J.P. (2017). Zinb-wave: A general and
495 flexible method for signal extraction from single-cell rna-seq data. *bioRxiv* , 125112.
- 496 SPINELLI, J.J., LOCKHART, R.A. & STEPHENS, M.A. (2002). Tests for the response distribution in a
497 poisson regression model. *Journal of statistical planning and inference* **108**, 137–154.

- 498 STUTE, W. & ZHU, L.X. (2002). Model checks for generalized linear models. *Scandinavian Journal of*
499 *Statistics* **29**, 535–545.
- 500 THAS, O. (2010). *Comparing Distributions*. New York, USA: Springer.
- 501 THAS, O. & RAYNER, J. (2005). Smooth tests for the zero-inflated Poisson distribution. *Biometrics* **61**,
502 808–815.
- 503 WEDDERBURN, R.W. (1974). Quasi-likelihood functions, generalized linear models, and the gauss-newton
504 method. *Biometrika* **61**, 439–447.
- 505 XU, L., PATERSON, A.D., TURPIN, W. & XU, W. (2015). Assessment and selection of competing models
506 for zero-inflated microbiome data. *PLoS one* **10**, e0129606.