# Mean and Scale-Factor Modeling of Under- and Overdispersed Count Data

**David M. Smith**
IBM Watson Health

**Malcolm J. Faddy**
Queensland University of Technology

### Abstract

This vignette describes an updated version of the R package **CountsEPPM** and its use in determining maximum likelihood estimates of the parameters of extended Poisson process models. These provide a Poisson process based family of flexible models that can handle both underdispersion and overdispersion in observed count data, with the negative binomial and Poisson distributions being special cases. Within **CountsEPPM** models with mean and scale-factor related to covariates are constructed to match a generalized linear model formulation. Use of the package is illustrated by application to several published datasets.

## 1. Introduction

This vignette, which relates to version 3.0 of **CountsEPPM**, is a revised and added-to version of Smith and Faddy (2016a) which related to version 2.1. The important differences between previous versions and 3.0 are a focus on mean and scale-factor models with variance models dropped, the addition of generic (S3) methods, using only `optim` for optimization i.e., no use of `nlm`, and offsets included in the formulae. Readers should refer to Smith and Faddy (2016a) for details of the actual modelling which involves Markov birth processes.

The models using extended Poisson process models (EPPMs) were originally developed in Faddy (1997), where the construction of discrete probability distributions having very general dispersion properties was described. The Poisson and negative binomial distributions are special cases of this modeling which includes both underdispersion and overdispersion relative to the Poisson, with the negative binomial having the most extreme level of overdispersion within the EPPM family. Faddy and Smith (2008) incorporated covariate dependence in the mean via a reparameterization using an approximate form of the mean; Faddy and Smith (2011) extended this to incorporate covariate dependence in the dispersion, this being achieved by a reparameterization using an approximate form of the variance. The supplementary material for Faddy and Smith (2011) contained R code illustrating fitting these models. This R code has been extended and generalized to have inputs and outputs more akin to those of a generalized linear model (GLM) as in the R function `glm` and the R function `betareg` (Cribari-Neto and Zeileis 2010, Grün, Kosmidis, and Zeileis 2012). Both Hilbe (2011) and Hilbe (2014) have commented about a software package for EPPMs being developed in the

R system (R Core Team 2016); the package **CountsEPPM** Smith and Faddy (2016b) whose use is described in this vignette is that software.

# 2. Description of the functions

The main function of the package, also named `CountsEPPM`, is focused on models with two covariate dependences linked to the mean and scale-factor. The input into the function is a formula involving a single response variable and one or two formulae related to the mean and scale-factor models. Although the input formula involves a single response variable, the actual model fitting has a list of frequency distributions `list.counts` in place of the response variable, which is either input or constructed from the input data according to whether a `list` or a `data.frame` is input. For all models the GLM link function between the response variable (mean, scale-factor) and linear predictor of covariates is log; the log of parameter $b$ of Equation 2 of Smith and Faddy (2016a) is also used but the parameter $c$ of Equation 2 of Smith and Faddy (2016a) is untransformed. The full three parameter version of Equation 2 of Smith and Faddy (2016a) has been labeled the Faddy distribution as in Grunwald, Bruce, Jiang, Strand, and Rabinovitch (2011). Because of possible issues with the parameter $b$, variants of the models where $b$ is fixed have been included in the lists of models. This enables profile log-likelihoods to be produced for this parameter.

Nash (2014) is a recent reference on optimization using R functions and contains information on, and insights into, the methods used. All models are fitted to the data using maximum likelihood, the optimization method used being the R function `optim` (options used being the simplex method of Nelder and Mead (1967) or `BFGS` using numerical derivatives). A facility to change options for `optim` through use of the argument `control` is included. The elements of this argument are the options for `optim` as described in R Core Team (2016). The default values set within `CountsEPPM` are `fnscale = -1`, `trace = 0`, `maxit = 1000` for `optim`. Although for most data sets the two options of `optim` give similar results in terms of log-likelihood and parameter estimates, etc., some results may be a little different depending on particular features of the data set. The simplex method, is robust to discontinuities in the log-likelihood surface. However, it is slow to converge. In contrast the function `BFGS` makes use of derivatives, in this case numerical derivatives, resulting in faster convergence, but there is a reliance on the log-likelihood surface being smooth i.e., no sudden changes in derivative values. Only `BFGS` makes use of derivatives in the actual model fitting, but both options calculate a hessian matrix from the derivatives to produce standard errors for the parameter estimates. The calculation of the numerical derivatives and hessians use the functions `grad` and `hessian` from the package **numDeriv** (Gilbert and Varadhan 2015). The derivatives are more accurately calculated using **numDeriv** (Gilbert and Varadhan 2015) than using alternative central difference approximations, resulting in better model fitting and better conditioned hessians. However, as stated in Nash (2014, p. 131), a longer time is taken. The occurrence of `NA` in the vector of standard errors is an indication of problems with the model fitting, possibly caused by an inappropriate model. This is particularly so when all estimates of parameter standard errors are `NA`, which results when the hessian matrix can not be be inverted due to its determinant being zero or it being otherwise ill-conditioned. Having derivatives available means that they can be reviewed, together with the hessian, at the conclusion of parameter estimation to evaluate whether maximum likelihood estimates have been attained.

The code for the main analysis function is

```
CountsEPPM(formula, data, subset = NULL, na.action = NULL, weights = NULL,
  model.type = "mean and scale-factor", model.name = "general", link = "log",
  initial = NULL, ltvalue = NA, utvalue = NA, method = "Nelder-Mead",
  control = NULL, fixed.b=NA)
```

with details of the arguments given in Table 1 together with defaults if any.

| Argument | Description | Default |
|---|---|---|
| `formula` | a single response variable & paired formulae Zeileis and Croissant (2010) | |
| `data` | `data.frame` or `list` | |
| `subset` | subsetting commands | NULL |
| `na.action` | action taken for NAs in data | NULL |
| `weights` | vector if `data` is a `data.frame` <br> a `list` if `data` is a `list` <br> attributes `normalization`, `norm.to.n` | vector of ones <br> list of lists of ones <br> both NULL |
| `model.type` | `"mean only"` <br> `"mean and scale-factor"` | `"mean and scale-factor"` |
| *if* `model.type = "mean only"` *(only a in Equation 2 of Smith and Faddy (2016a) modeled)* | | |
| `model.name` | `"Poisson"` <br> `"negative binomial"` <br> `"negative binomial fixed b"` <br> `"Faddy distribution"` Equation 2 of Smith and Faddy (2016a) <br> `"Faddy distribution fixed b"` | |
| *if* `model.type = "mean and scale-factor"` *(both modeled)* | | |
| `model.name` | `"general"` as Equations 3 and 4 <br> `"general fixed b"` <br> `"limiting"` as Equations 9 and 10 of Smith and Faddy (2016a) | `"general"` |
| `link` | the glm link function for mean count only log allowed | `"log"` |
| `initial` | vector of initial values for parameters, means first followed by the variances and/or parameters $c$, $\log(b)$ | Poisson `glm` output augmented by 0's for other parameters |
| `ltvalue` | lower truncation value (excluded) | NA |
| `utvalue` | upper truncation value (excluded) | NA |
| `method` | `"Nelder-Mead"` <br> `"BFGS"` attribute `"grad.method"` <br> which is `"simple"` or `"Richardson"` | `"Nelder-Mead"` <br> attribute `"simple"` |
| `control` | list of control parameters `optim` | see text for more details |
| `fixed.b` | value $b$ is fixed at | NA |

Table 1: Arguments of `CountsEPPM`.

As in earlier versions, `data` can be either a `list` or a `data.frame`. The response variable in `formula` is a vector if a `data.frame` is input or a `list` if a `list` is input. The response variables `mean.obs` and `variance.obs` are constructed within `CountsEPPM` prior to being used to fit models. The R package **Formula** of Zeileis and Croissant (2010) is used to extract model information from the `formula` input to `CountsEPPM`. To avoid repeated extractions within subordinate functions, the extraction of model information used in the model fitting, such as `covariates.matrix.mean`, is only done once within `CountsEPPM`. In version 3.0 a set of S3 generic extractor functions for objects of class `"CountsEPPM"` has been added. The set is the same as that for **BinaryEPPM** (Smith and Faddy 2018) which is itself similar to that of Table 1 of `betareg` (Cribari-Neto and Zeileis 2010).

As iteration is involved in the model fitting, initial estimates of the parameters are needed. These can optionally be provided in the vector `initial`. Within `CountsEPPM`, if `initial` is unset, a Poisson model is fitted using `glm` and the estimates from that fit are used to provide estimates for the parameters of the mean linear predictor. If the scale-factor is also being modeled the initial estimates of the parameters of the scale-factor linear predictor are set to 1.0 recognising that for the Poisson distribution the variance equals the mean. The initial value of $\log(b)$ of Equation 2 of Smith and Faddy (2016a) is set to zero. The matrix exponential function used for calculating the probabilities of Equation 1 of Smith and Faddy (2016a) is that of the package **expm** of Goulet, Dutang, Maechler, Firth, Shapira, and Stadelmann (2015) which depends on the package **Matrix** of Bates and Maechler (2016). `CountsEPPM` returns an object of class `"CountsEPPM"` summarizing the model fit, the components of which are given in Table 2.

Table 3 gives details of a set of S3 generic extractor functions for objects of class `"CountsEPPM"`. The set is similar to that of Table 1 of Cribari-Neto and Zeileis (2010) related to package **betareg**, except there are no functions `estfun`, `bread` or `linear.hypothesis`. Also, `gleverage` and `cooks.distance` are variants of the functions `glm.diag` and `glm.diag.plots` from package **boot** Canty and Ripley (2017) rather than **betareg**. The first four blocks refer to functions specific to **CountsEPPM**. The last block contains generic functions, the default versions of which work because of the information supplied by the functions of the first four blocks. Package **lmtest** Zeileis and Hothorn (2002) needs to be loaded to use `coeftest` and `lrtest`. Function `AIC` comes from **stats** which is a default package loaded when R is started. In Table 2 both `n` and `nobs` are included, so that functions from both packages **lmtest** and **stats** can use the object returned.

As the vectors of frequency distributions are only required to be of length the maximum observed count value +1, this is how they are set up. However, the fitted models can have probability masses at counts greater than these maximum counts. A component of the output object from `CountsEPPM` i.e., `$vnmax` is a vector of the maximum observed counts. If probabilities for counts greater than these maximums are wanted, the values in `output.fn$vnmax` can be increased in value and `predict` with `type="distribution"` run to obtain these probabilities.

| Component | Description |
|---|---|
| `data.type` | `"data.frame"` or `"list"` |
| `list.data` | data as a `"list"` of frequency distributions |
| `call` | the call to **CountsEPPM** |
| `formula` | the `formula` input |
| `model.type` | `"mean only"` or `"mean and scale-factor"` |
| `model.name` | `"Poisson"`, `"negative binomial"`, `"negative binomial fixed b"`, `"Faddy distribution"` (Equation 2), `"Faddy distribution fixed b"`, `"general"` (Equations 3 & 4), `"general fixed b"`, `"limiting"` (Equations 9 & 10). Equation numbers of Smith and Faddy (2016a) |
| `link` | the glm link function for mean count |
| `covariates.matrix.mean` | matrix of covariates for the mean |
| `covariates.matrix.scalef` | matrix of covariates for the scale-factor |
| `offset.mean` | offset vector for the mean |
| `offset.scalef` | offset vector for the scale-factor |
| `coefficients` | the estimated coefficients |
| `loglik` | the final log-likelihood value |
| `vcov` | the estimated variance/covariance matrix |
| `n` needed for `lmtest` | the number of observations |
| `nobs` needed for `stats` | the number of observations |
| `df.null` | null model degrees of freedom |
| `df.residual` | residual degrees of freedom |
| `ltvalue` | lower truncation value (excluded) |
| `utvalue` | upper truncation value (excluded) |
| `fixed.b` | value $b$ is fixed at |
| `vnmax` | a vector of maximum counts in each of the grouped data vectors |
| `weights` | a vector or list of weights |
| `converged` | whether converged |
| `iterations` | number of iterations |
| `method` | `"Nelder-Mead"` or `"BFGS"` |
| `start` | initial estimates input |
| `optim` | final estimates of coefficients |
| `control` | control parameters of `optim` |
| `fitted.values` | fitted values of mean count |
| `y` | observed values of mean count |
| `terms` | model terms |

Table 2: Components of object returned by `CountsEPPM`.

| Function | Description |
|---|---|
| `print()` | a simple printed display |
| `summary()` | standard regression output (coefficient estimates, standard errors, partial Wald tests); returns an object of class `summary.BinaryEPPM` containing the relevant summary statistics (which has a `print` method) |
| `coef()` | extract coefficients of model (full, mean, or precision components), a single vector of all coefficients by default |
| `vcov()` | associated covariance matrix (with matching names) |
| `predict()` | predictions (response, linear predictor $p_s$, linear predictor scale-factor, $p_s$, scale-factor, scale-factor limits, mean, variance, distribution probabilities, distribution parameters) for existing and new data |
| `fitted()` | fitted means for observed data |
| `residuals()` | extract residuals (deviance, Pearson, response, standardized deviance, standardized Pearson residuals), defaulting to standardized Pearson residuals |
| `terms()` | extract terms of model components |
| `model.matrix()` | extract model matrix of model components |
| `model.frame()` | extract full original model frame |
| `logLik()` | extract fitted log-likelihood |
| `plot()` | diagnostic plots of residuals, predictions, leverages, etc. |
| `hatvalues()` | hat values (diagonal of hat matrix) |
| `cooks.distance()` | Cook's distance |
| `gleverage()` | generalized leverage |
| `waldtest()` | Wald tests of model parameters |
| `coeftest()` | partial Wald tests of coefficients |
| `lrtest()` | likelihood ratio tests of model parameters |
| `AIC()` | compute information criteria (AIC, BIC, . . . ) |

Table 3: Generic Functions for Use with Objects of Class `CountsEPPM`.

# 3. Examples

The examples illustrate various ways in which **CountsEPPM** can be used to produce informative analyses. For the first three examples `data` is a `list` where the dependent variable of counts is a `list` of vectors of frequency distributions, whereas for the last two `data` is a `data frame`. Significant time savings can be made by using the `list` form of input where applicable. The fitting of models and estimation of their parameters can be sensitive to the initial estimates and method of estimation chosen, with flatness of the log-likelihood surface possible, particularly with respect to the parameter *b*. It is recommended that analyses be run more than once using different initial estimates and optimization methods.

## 3.1. Number of young at varying effluent concentrations data

These Ceriodaphnia dubia data were used as an example by Faddy and Smith (2011). Ce-

riodaphnia dubia are water fleas used to test the impact of effluents on water quality. The data, originally from Bailer and Oris (1997), are counts of young at varying effluent concentrations, and are in `list` of frequencies and variables form. The defaults for `model.type` of `"mean and variance"`, and `model` of `"general"`, are used. The code given below is for a run using the `method="simplex"` option of `optim` followed by one using `BFGS` with `attribute="Richardson"` with this last run giving derivatives (gradients) at the final estimates. Although during these runs the estimate of $\log(b)$ changed sign, its standard error is relatively large and the estimates of the other parameters had relatively small changes in value.

```
R> data("ceriodaphnia.group")
R> output.fn <-
+   CountsEPPM(number.young ~ 1 + vdose + vdose2 | 1 + vdose + vdose2,
+   data = ceriodaphnia.group, control = list(maxit = 4000, reltol = 1.e-11))
R> names(output.fn$optim$par) <- c('mean Intercept', 'mean dose',
+   'mean dose^2', 'scale-factor Intercept', 'scale-factor dose',
+   'scale-factor dose^2', 'log(b)')
R> method <- "BFGS"
R> attr(method, which = "grad.method") <- "Richardson"
R> output.fn <- update(output.fn, initial = output.fn$optim$par,
+   method = method)
R> summary(output.fn)


 Dependent variable is a list of frequency distributions for counts

Call:
CountsEPPM(formula = number.young ~ 1 + vdose + vdose2 | 1 + vdose +
    vdose2, data = ceriodaphnia.group, initial = output.fn$optim$par,
    method = method, control = list(maxit = 4000, reltol = 1e-11))


Model type        : mean and scale-factor
Model name        : general
Link scale-factor : log

 Coefficients (model for mean with log link)


t test of coefficients:

                 Estimate Std. Error t value  Pr(>|t|)
mean Intercept  3.1406671  0.0859089 36.5581 < 2.2e-16 ***
mean dose       0.1733786  0.0302852  5.7249 9.177e-07 ***
mean dose^2    -0.0196106  0.0025092 -7.8155 8.666e-10 ***
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1


 Coefficients (model for scale-factor with log link)
```

```
t test of coefficients:

                       Estimate Std. Error t value Pr(>|t|)
scale-factor Intercept  1.297770   0.445457  2.9133 0.005653 **
scale-factor dose      -0.665782   0.219006 -3.0400 0.004017 **
scale-factor dose^2     0.047337   0.015775  3.0008 0.004469 **
log(b)                 -0.144932   2.625530 -0.0552 0.956234
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1

 Type of estimator: ML (maximum likelihood)
 Log-likelihood: -152.1356 on 7 Df
 Number of iterations: 67 of optim method BFGS gradient method Richardson

 final gradients of parameters
[1] 0.0020328358 0.0086772311 0.0392932715 0.0001186455 0.0005086497
[6] 0.0103859013 0.0019768983

 return code 0 successful
```

The above parameter estimates agree with those in Faddy and Smith (2011) to two decimal places, except for $\log(b)$ which is relatively poorly estimated. Further details of the model can be printed out such as the parameters of the Faddy distribution as well as the predicted values of the means, variances and scale-factors.

```
R> predict(output.fn, type = "distribution.parameters")

      out.va     out.vb      out.vc
1   7.661167 0.8650815  0.5179911
2  18.489020 0.8650815  0.1830147
3  56.473485 0.8650815 -0.2029801
4 414.407542 0.8650815 -0.9161889
5   6.799275 0.8650815  0.2157040


R> predictions <- data.frame(
+    mean = predict(output.fn, type = "mean"),
+    variance = predict(output.fn, type = "variance"),
+    scale.factor = predict(output.fn, type = "scale.factor"))
R> print(predictions)

       mean variance scale.factor
1 23.119285 84.64258    3.6611244
2 28.895726 41.96128    1.4521623
3 32.817612 23.81793    0.7257668
4 31.761146 11.51863    0.3626641
5  9.428539 13.67537    1.4504229
```

To illustrate use of the `newdata` argument of `predict` a new `data.frame` of the second and third rows is constructed and used with `predict`.

```
R> newdata <- data.frame(intercept = rep(1,2),
+    vdose = ceriodaphnia.group$vdose[2:3],
+    vdose2 = ceriodaphnia.group$vdose2[2:3], vnmax = c(35, 44))
R> predictions <- data.frame(
+    mean = predict(output.fn, newdata = newdata, type = "mean"),
+    variance = predict(output.fn, newdata = newdata, type = "variance"),
+    scale.factor = predict(output.fn, newdata = newdata,
+    type = "scale.factor"))
R> print(predictions)


       mean variance scale.factor
1 28.89573 41.96128    1.4521623
2 32.81761 23.81793    0.7257668
```

### 3.2. Lüning et al. data

These data are from Lüning, Sheridan, Ytterborn, and Gullberg (1966) and are in the form of a `list` of frequencies and variables. The number of trials are stated in Lüning *et al.* (1966) to be both lower and upper truncated at 4 and 11 respectively, so the data are for counts of 5 to 10. Default initial values are used for fitting the default `general` model of Equations 3 and 4 of Smith and Faddy (2016a).

```
R> data("Luningetal.litters")
R> output.fn <- CountsEPPM(number.trials ~ 0 + fdose | 0 + fdose,
+    Luningetal.litters, ltvalue = 4, utvalue = 11,
+    control = list(maxit = 2000))
R> summary(output.fn)

 Dependent variable is a list of frequency distributions for counts

 distribution truncated below at   4
 distribution truncated above at  11


Call:
CountsEPPM(formula = number.trials ~ 0 + fdose | 0 + fdose,
    data = Luningetal.litters, ltvalue = 4, utvalue = 11,
    control = list(maxit = 2000))


Model type        : mean and scale-factor
Model name        : general
Link scale-factor : log


 Coefficients (model for mean with log link)
```

```
t test of coefficients:

         Estimate Std. Error t value  Pr(>|t|)
fdose0   1.9171890  0.0022050 869.471 < 2.2e-16 ***
fdose300 1.8321149  0.0099383 184.350 < 2.2e-16 ***
fdose600 1.7239822  0.0186856  92.262 < 2.2e-16 ***
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1


 Coefficients (model for scale-factor with log link)


t test of coefficients:

          Estimate Std. Error  t value  Pr(>|t|)
fdose0   -1.446450   0.074170 -19.5018 < 2.2e-16 ***
fdose300 -1.365485   0.092402 -14.7777 < 2.2e-16 ***
fdose600 -1.209275   0.138516  -8.7302 < 2.2e-16 ***
log(b)   19.395048         NA       NA        NA
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1


 Type of estimator: ML (maximum likelihood)
 Log-likelihood: -2653.815 on 7 Df
 Number of iterations: 1068 of optim method Nelder-Mead


 return code 0 successful
Warning message:
In sqrt(diag(se)) : NaNs produced
```

The warning message In sqrt(diag(varcov)) : NaNs produced and the value of $b =$ exp(19.395048) being large suggests that the fitted model corresponds to a negative expeonential sequence of rate parameters in the underlying birth process (Equations 9 and 10 of Smith and Faddy (2016a)).

```
R> output.fn <- update(output.fn, model.name = 'limiting')
R> summary(output.fn)


 Dependent variable is a list of frequency distributions for counts

 distribution truncated below at  4
 distribution truncated above at  11

Call:
CountsEPPM(formula = number.trials ~ 0 + fdose | 0 + fdose,
    data = Luningetal.litters, model.name = "limiting",
    ltvalue = 4, utvalue = 11, control = list(maxit = 2000))
```

```
Model type        : mean and scale-factor
Model name        : limiting
Link scale-factor : log

 Coefficients (model for mean with log link)

t test of coefficients:

          Estimate Std. Error t value  Pr(>|t|)
fdose0    1.9169670  0.0077381 247.730 < 2.2e-16 ***
fdose300 1.8320442  0.0099421 184.271 < 2.2e-16 ***
fdose600 1.7242914  0.0186506  92.453 < 2.2e-16 ***
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1

 Coefficients (model for scale-factor with log link)

t test of coefficients:

          Estimate Std. Error  t value  Pr(>|t|)
fdose0    -1.444080    0.074503 -19.3828 < 2.2e-16 ***
fdose300 -1.366303    0.092760 -14.7294 < 2.2e-16 ***
fdose600 -1.210711    0.138765  -8.7249 < 2.2e-16 ***
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1

 Type of estimator: ML (maximum likelihood)
 Log-likelihood: -2653.816 on 6 Df
 Number of iterations: 533 of optim method Nelder-Mead

 return code 0 successful
```

The parameters of the limiting model can be printed out

```
R> predict(output.fn, type = "distribution.parameters")

  out.valpha  out.vbeta
1   22.99626 -0.3067979
2   18.91508 -0.3070639
3   13.95311 -0.2872457
```
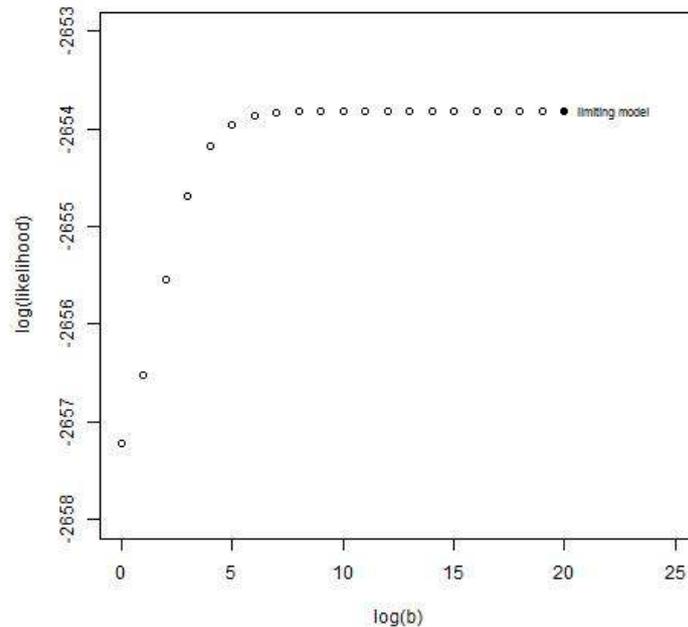
showing much the same results as the previous fit of the general model. To further explore the appropriateness of the limiting model a profile likelihood was constructed for a range of values of parameter $b$ from the version of the general fixed b model of Equations 3 and 4 of Smith and Faddy (2016a) with a plot of the resulting log(likelihoods) against log(b) being produced. In Figure 1 there is a clear trending of the log-likelihood values toward the value of the limiting model. The code used to produce Figure 1 follows.

Figure 1: log-likelihood for fixed values of parameter log(b).

```
R> vfixed.b <- c(0:19)
R> vloglikelihood <- rep(0, 20)
R> vloglikelihood <- sapply(1:20, function(i) {
R>  if (i == 1) {
R>    output.fn <- CountsEPPM(number.trials ~ 0 + fdose | 0 + fdose,
+        Luningetal.litters, model.name = 'general fixed b',
+        ltvalue = 4, utvalue = 11, fixed.b = exp(vfixed.b[i]))
R>  } else {
R>    output.fn <- CountsEPPM(number.trials ~ 0 + fdose | 0 + fdose,
+        Luningetal.litters, model.name = 'general fixed b',
+        ltvalue = 4, utvalue = 11, initial = output.fn$optim$par,
+        fixed.b = exp(vfixed.b[i]))
R>         } # end of if (i==1)
R>  vloglikelihood[i] <- output.fn$loglik } )

R> plot(vfixed.b, vloglikelihood, xlim = c(0, 25), ylim =c (-2658, -2653),
+     main = "Profile likelihood for log(b) Luning litters data",
+     xlab = "log(b)", ylab = "log(likelihood)")
R> points(20, lm.loglik, pch = 16)
R> text(20.1, lm.loglik, "limiting model", pos = 4, offset = 0.5, cex = 0.7)
```

### 3.3. Number of attempts at feeding of herons

These data are originally from Zhu, Eickhoff, and Kaiser (2003) and are in the form of a
list of frequencies and variables. Faddy and Smith (2005) described an alternative modeling

approach to that of Zhu *et al.* (2003) constructing a bivariate EPPM for both count (number of attempts) and grouped (number of successful attempts) data. Here a univariate EPPM for the numbers of trials (attempts at foraging) of 20 adult and 20 immature green-backed herons is considered. The first model fitted was a negative binomial using the default initial values.

```
R> data("herons.group")
R> output.fn.one  <- CountsEPPM(number.attempts ~ 0 + group, herons.group,
+     model.type = 'mean only', model.name = 'negative binomial')
R> names(output.fn.one$optim$par) <- c('Adult mean',
+     'Immature mean', 'log(b)')
R> print(summary(output.fn.one))


 Dependent variable is a list of frequency distributions for counts

Call:
CountsEPPM(formula = number.attempts ~ 0 + group, data = herons.group,
    model.type = "mean only", model.name = "negative binomial")

Model type         : mean only
Model name         : negative binomial

 Coefficients (model for mean with log link)

t test of coefficients:

              Estimate Std. Error t value  Pr(>|t|)
group Adult    0.56230    0.15500  3.6279 0.0008573 ***
group Immature 0.47586    0.16439  2.8947 0.0063315 **
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1

 Coefficients (model for scale-factor with log link)

t test of coefficients:

        Estimate Std. Error t value Pr(>|t|)
log(b)  0.50825    0.26793   1.897  0.06566 .
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1

 Type of estimator: ML (maximum likelihood)
 Log-likelihood: -120.2042 on 3 Df
 Number of iterations: 108 of optim method Nelder-Mead

 return code 0 successful
```

The second model fitted was a more general Faddy distribution again using the default initial values.

```
R> output.fn.two <- update(output.fn.one, model.name = 'Faddy distribution')
R> names(output.fn.two$optim$par) <- c('Adult mean', 'Immature mean',
+     'c', 'log(b)')
R> print(summary(output.fn.two))

 Dependent variable is a list of frequency distributions for counts

Call:
CountsEPPM(formula = number.attempts ~ 0 + group, data = herons.group,
    model.type = "mean only", model.name = "Faddy distribution")

Model type        : mean only
Model name        : Faddy distribution

 Coefficients (model for mean with log link)

t test of coefficients:

               Estimate Std. Error t value  Pr(>|t|)
group Adult     0.56282    0.15496  3.6320 0.0008688 ***
group Immature  0.47652    0.16436  2.8993 0.0063351 **
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1

 Coefficients (model for scale-factor with log link)

t test of coefficients:

        Estimate Std. Error    t value Pr(>|t|)
c     1.0000e+00 1.5724e-05 63596.0001   <2e-16 ***
log(b) 5.0703e-01 2.6792e-01     1.8924   0.0665 .
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1

 Type of estimator: ML (maximum likelihood)
 Log-likelihood: -120.2042 on 4 Df
 Number of iterations: 331 of optim method Nelder-Mead

 return code 0 successful
```

The estimate of the parameter $c$ here is essentially one, the upper limit of permitted values of this parameter corresponding to a negative binomial model, and its standard error is largely noise; Wald and (log-)likelihood tests on this parameter are therefore invalid. Print outs of the parameters of the Faddy distribution and the associated distribution for the first group can be produced.

```
R> predict(output.fn.two, type = "distribution.parameters")
R> predict(output.fn.two, type = "distribution")[1]


    out.va   out.vb    out.vc
1 1.755613 1.66035 0.9999997
2 1.610463 1.66035 0.9999997

[[1]]
 [1] 0.054207702 0.074451022 0.081919897 0.082680037 0.079683626 0.074619482
 [7] 0.068518504 0.062025215 0.055542261 0.049315532 0.043487554 0.038132360
[13] 0.033278821 0.028926520 0.025056608 0.021639253 0.018638723 0.016016821
[19] 0.013735148 0.011756573 0.010046116 0.008571446 0.007303109 0.006214567
[25] 0.005282124
```

By increasing the maximum values for the grouped counts using the following code, the probabilities for the next ten counts in the sequence for the first group can be obtained.

```
R> wks <- output.fn.two$vnmax [1] + 2
R> wke <- output.fn.two$vnmax[1] + 11
R> output.fn.two$vnmax[1] <- output.fn.two$vnmax[1] + 10
R>  predict(output.fn.two, type = "distribution")[[1]][wks:wke]

[[1]]
 [1] 0.0044847776 0.0038040231 0.0032236418 0.0027294804 0.0023092342
 [6] 0.0019522416 0.0016492910 0.0013924445 0.0011748767 0.0009907317
```

A weighted analysis using the reciprocal of the predicted variances can be performed.

```
R> herons.group$weights <- herons.group$number.attempts
R> weights <- 1 / predict(output.fn.one, type = "variance")
R> herons.group$weights <- lapply(1:length(herons.group$weights), function(i) {
+   herons.group$weights[[i]] <- rep(weights[i],
+   length(herons.group$weights[[i]])) } ) # end of lapply
R> attr(herons.group$weights, which = "normalize") <- TRUE
R> output.fn  <- CountsEPPM(number.attempts ~ 0 + group, herons.group,
+   model.type = 'mean only', model.name = 'Poisson',
+   weights = herons.group$weights)
R> names(output.fn$optim$par) <- c('Adult mean', 'Immature mean')
R> summary(output.fn)


 Dependent variable is a list of frequency distributions for counts

Call:
CountsEPPM(formula = number.attempts ~ 0 + group, data = herons.group,
    weights = herons.group$weights, model.type = "mean only",
    model.name = "Poisson")
```

```
Model type         : mean only
Model name         : Poisson

 Coefficients (model for mean with log link)

t test of coefficients:

               Estimate Std. Error t value  Pr(>|t|)
group Adult     2.073172   0.086557  23.951 < 2.2e-16 ***
group Immature 1.894617   0.080490  23.538 < 2.2e-16 ***
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1

 Maximum weighted likelihood regression.
 List of weights used.

 Type of estimator: ML (maximum likelihood)
 Log-likelihood: -168.3474 on 2 Df
 Number of iterations: 45 of optim method Nelder-Mead

 return code 0 successful
```

The same data as `herons.group`, but in `data.frame` form, has been included in the package as `herons.case`. Running the same code with `herons.group` replaced by `herons.case` will produce essentially the same outputs.

### 3.4. Titanic survivors

To illustrate the inclusion of offsets, data of passenger survival from the 1912 sinking of the Titanic are used. The data are in data frame form as given in Table 9.37 of Hilbe (2011) i.e., the numbers surviving out of the number of cases (passengers) within different age, sex, and class categories. The individual data for all 1316 passengers is available from package **msme** Hilbe and Robinson (2014). Hilbe (2011, p. 265–268) analyzes the numbers surviving as count data with an offset of the log of the number of cases for the mean, and fits a negative binomial with variance function $v = m + \alpha m^2$ which equates to the variance function of Equation 4 of Smith and Faddy (2016a) with $\alpha = \frac{1}{b}$. Both mean and scale-factor need to be offset by the log of the number of cases. A series of models was fitted: a negative binomial with the parameter $b$ fixed at the value from Hilbe (2011) of $b = 9.615385$; a negative binomial with $b$ unspecified; a more general Faddy distribution; and a general mean and constant scale-factor model. The last of these models was found to have the largest log-likelihood. Details of it and its fitting follow.

```
R> data("Titanic.survivors.case")
R> lncases <- log(Titanic.survivors.case$cases)
R> output.fn <- CountsEPPM(survive ~ age + sex + class + offset(lncases) |
+   1 + offset(lncases), Titanic.survivors.case, control = list(maxit = 2000))
R> names(output.fn$optim$par) <- c('Intercept mean', 'age adult', 'sex male',
```

```
+   'class 2nd class', 'class 3rd class', 'Intercept scale', 'log(b)')
R> output.fn <- update(output.fn, initial = output.fn$optim$par,
+    method = 'BFGS')
R> summary(output.fn)


 Dependent variable a vector of counts.

Call:
CountsEPPM(formula = survive ~ age + sex + class + offset(lncases) |
    1 + offset(lncases), data = Titanic.survivors.case,
    initial = output.fn$optim$par, method = "BFGS",
    control = list(maxit = 2000))


Model type         : mean and scale-factor
Model name         : general
Link scale-factor : log
non zero offsets in linear predictors

 Coefficients (model for mean with log link)


t test of coefficients:

                  Estimate Std. Error t value Pr(>|t|)
Intercept mean    0.047471   0.083542  0.5682 0.594447
age adult        -0.043300   0.133024 -0.3255 0.757981
sex male         -0.177376   0.129474 -1.3700 0.229012
class 2nd class  -0.031658   0.113242 -0.2796 0.791010
class 3rd class  -0.900905   0.144787 -6.2223 0.001568 **
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1


 Coefficients (model for scale-factor with log link)


t test of coefficients:

                  Estimate Std. Error t value  Pr(>|t|)
Intercept scale  -4.04624    0.54017 -7.4907 0.0006701 ***
log(b)           -7.32930    3.64626 -2.0101 0.1006346
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1


 Type of estimator: ML (maximum likelihood)
 Log-likelihood: -39.39972 on 7 Df
 Number of iterations: 60 of optim method BFGS gradient method simple

 final gradients of parameters
[1] -0.0203696866 -0.0008164389  0.0041572081  0.0023752619  0.0068211516
```

```
[6]   0.0018952763   0.0007120647
```

```
 return code 0 successful
```

The parameters of the Faddy distribution can now be printed out.

```
R> predict(outpu.fn, type = "distribution.parameters")
```

```
            out.va           out.vb           out.vc
1        0.1393159  0.0006560301  -28.10881810
2       45.9529402  0.0006560301    0.30300954
3     1594.9387938  0.0006560301   -5.21904057
4       40.8724192  0.0006560301    0.33977379
5      392.9331778  0.0006560301   -1.69941164
6       47.0261619  0.0006560301    0.19278209
7      330.3710419  0.0006560301   -2.09935775
8       39.9238096  0.0006560301    0.33263580
9       27.6100183  0.0006560301   -0.42233391
10      25.0226409  0.0006560301    0.33016005
11      20.5073026  0.0006560301   -0.09566396
12      28.2101400  0.0006560301    0.46262359
```

The fit of the general mean and scale-factor model is better than that of Hilbe (2011)(page 268), the log-likelihood values being $-39.400$ and $-43.719$ respectively; although the former has one extra parameter it would be preferred according to AIC.

### 3.5. Take over bids

These data, originally from Cameron and Johansson (1997), are used as example data in Cameron and Trivedi (2013) as well as in Sáez-Castillo and Conde-Sánchez (2013). The `takeover.bids.case` came from the website associated with Cameron and Trivedi (2013). The dependent variable NUMBIDS is the number of bids received by the firm targeted for takeover after the initial bid. As both variables CASE, CONSTANT are equal to 1 throughout they have not been included in package data set. In Smith and Faddy (2016a), related to version 2.1 of **CountsEPPM**, the continuous variables were scaled to have zero mean and unit standard deviation prior to analysis, as it was found that the scaling of the continuous variables improved the model fitting. The changes and additions made to version 3.0 make this unnecessary.

```
method <- "BFGS"
attr(method,which = "grad.method") <- "Richardson"
output.fn  <- CountsEPPM(NUMBIDS ~ LEGLREST + REALREST + FINREST +
 WHTKNGHT + BIDPREM + INSTHOLD + SIZE + SIZESQ + REGULATN | LEGLREST +
 REALREST + FINREST + WHTKNGHT + BIDPREM + INSTHOLD + SIZE + SIZESQ +
 REGULATN, data = takeover.bids.case, method = method)
summary(output.fn)
```

```
 Dependent variable a vector of counts.

Call:
CountsEPPM(formula = NUMBIDS ~ LEGLREST + REALREST + FINREST + WHTKNGHT +
    BIDPREM + INSTHOLD + SIZE + SIZESQ + REGULATN | LEGLREST + REALREST +
    FINREST + WHTKNGHT + BIDPREM + INSTHOLD + SIZE + SIZESQ + REGULATN,
    data = takeover.bids.case, method = method)

Model type        : mean and scale-factor
Model name        : general
Link scale-factor : log

 Coefficients (model for mean with log link)

t test of coefficients:

                Estimate  Std. Error  t value  Pr(>|t|)
(Intercept)   1.66980728  0.69448315    2.4044 0.0179509 *
LEGLREST      0.26775754  0.18027398    1.4853 0.1404656
REALREST     -0.16655260  0.29485396   -0.5649 0.5733706
FINREST       0.67982048  0.34621124    1.9636 0.0522207 .
WHTKNGHT      0.81090836  0.23384727    3.4677 0.0007622 ***
BIDPREM      -1.66354592  0.54936971   -3.0281 0.0030970 **
INSTHOLD     -0.81832858  0.53231273   -1.5373 0.1272260
SIZE          0.30109055  0.01639493   18.3649 < 2.2e-16 ***
SIZESQ       -0.01432358  0.00076058  -18.8325 < 2.2e-16 ***
REGULATN      0.24525971  0.22094221    1.1101 0.2695077
---
 Signif. codes: "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1

 Coefficients (model for scale-factor with log link)

t test of coefficients:

              Estimate Std. Error t value Pr(>|t|)
(Intercept)   0.443199   0.719919  0.6156   0.5395
LEGLREST     -0.104376   0.186026 -0.5611   0.5759
REALREST      0.214242   0.212209  1.0096   0.3150
FINREST       0.533311   0.513642  1.0383   0.3015
WHTKNGHT      0.267462   0.328230  0.8149   0.4170
BIDPREM      -0.840780   0.919408 -0.9145   0.3626
INSTHOLD     -0.027377   0.366394 -0.0747   0.9406
SIZE          0.145692   0.130238  1.1187   0.2658
SIZESQ       -0.015642   0.013306 -1.1756   0.2424
REGULATN      0.263825   0.366727  0.7194   0.4735
log(b)       -5.339202   6.111304 -0.8737   0.3843
```

```
Type of estimator: ML (maximum likelihood)
Log-likelihood: -156.0563 on 21 Df
Number of iterations: 244 of optim method BFGS gradient method Richardson

 final gradients of parameters
 [1]  0.0011504814  0.0007030922  0.0013814654  0.0003727108 -0.0005461059
 [6]  0.0009490258  0.0012592771  0.0046084397  0.0020728355  0.0022824005
[11] -0.0007277559  0.0017569373 -0.0024497681  0.0009933780 -0.0001212738
[16] -0.0022219230 -0.0028663406 -0.0052951310  0.0457381566 -0.0008509123
[21] -0.0001037344


 return code 0 successful
```

The Bayesian Information Criterion (BIC) can also be calculated using `BIC(output.fn)` resulting in a value of 413.6745. Convergence to the maximum likelihood estimates was slow due to the large number (21) of parameters being estimated, the flatness of the log-likelihood surface and a small estimate of the (nuisance) parameter $b$ (negative value of $\log(b)$). The estimate of $b$ being close to 0 and the underdispersion (scale-factor $< 1$) corresponding to $c < 0$ in Equation 2 of Smith and Faddy (2016a) means that the zero-probability can be very small (Equation 1 of Smith and Faddy (2016a)). This model with a log-likelihood value of $-156.06$ and 21 parameters fits better than that from Sáez-Castillo and Conde-Sánchez (2013) with $-157.86$ and 15 parameters. But the six extra parameters associated with a relatively small increase in log-likelihood means that the BIC value of 413.67 is larger than those of the models in Sáez-Castillo and Conde-Sánchez (2013): 398.1, 393.5 and 388.3. Sáez-Castillo and Conde-Sánchez (2013) did not report details of fitting a model with the full set of 10 variables in both linear predictors, suggesting that they could not achieve convergence of their fitting algorithm for this model. The variables Sáez-Castillo and Conde-Sánchez (2013) do not include in their dispersion model as they were claimed to be not significant are `LEGLREST`, `WHTKNGHT`, `INSTHOLD`, `SIZE`, `SIZESQ`. There is reasonable agreement between the results reported above and those of Sáez-Castillo and Conde-Sánchez (2013) about the significant variables in the mean model; in both, `SIZE` and `SIZESQ` have very large $t$ statistics. However, in the dispersion model the results reported above show that `SIZE` and `SIZESQ` also have very large $t$ statistics, whereas in the Sáez-Castillo and Conde-Sánchez (2013) models they were not included. Residual plots as in Cribari-Neto and Zeileis (2010) can be produced.

```
layout(matrix(c(1:6), byrow=TRUE, ncol=2))
plot(output.fn, which = 1, type = "response")
plot(output.fn, which = 2, type = "pearson")
plot(output.fn, which = 3, type = "spearson")
plot(output.fn, which = 4, type = "likelihood")
plot(output.fn, which = 5, type = "deviance")
plot(output.fn, which = 6, type = "sdeviance")
```

Examination of the estimated $\lambda_i$ sequences of (birth) rate parameters (Equation 2 of Smith and Faddy (2016a)) shows that the $\lambda_0$ values are generally very different from the $\lambda_i$ values for $i \geq 1$, as a consequence of the small estimate of the parameter $b$. A more appropriate model for these data might be one that treats the zero counts differently from the non-zero counts;
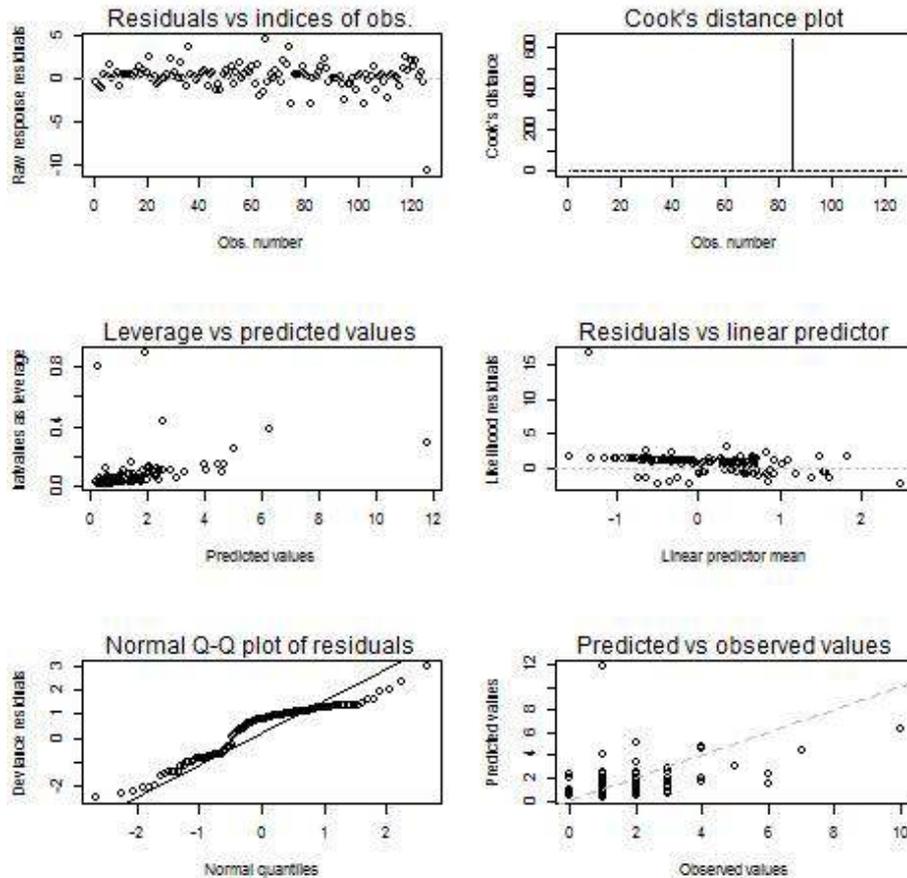
Figure 2: Residual plots.

this makes some sense as a zero count would correspond to the initial bid being accepted by the targeted firm, and different circumstances (in the form of different covariate dependence) might be operating. Such a model is not readily constructed from those considered here, and is therefore beyond the scope of the functions developed. Such modeling will be the subject of future work.

# 4. Concluding remarks

This vignette has described the use of version 3.0 of the R package **CountsEPPM** to fit EPPMs to count data that exhibit under- or over-dispersion relative to the Poisson distribution. A variety of covariate dependencies and data structures are covered in examples that illustrate the variety of ways in which the package can be used in the analysis of count data.

As described in Faddy and Smith (2005) and Faddy and Smith (2012), the mean and scale-factor of binary data can be modeled using EPPMs in a similar way to that described here for count data. A package **BinaryEPPM** is available in CRAN, and an article describing its use will shortly be appearing in the Journal of Statistical Software. Dependent on how the further work goes, similar functions, etc., may be developed for the model (zero and non-zero counts treated differently) mentioned in the last paragraph of the previous section.

# References

Bailer AJ, Oris JT (1997). "Estimating Inhibition Concentrations for Different Response Scales Using Generalized Linear Models." *Environmental Toxicology and Chemistry*, **16**(7), 1554–1559. `doi:10.1002/etc.5620160732`.

Bates D, Maechler M (2016). **Matrix**: *Sparse and Dense Matrix Classes and Methods*. R package version 1.2-4, URL `https://CRAN.R-project/package=Matrix`.

Cameron AC, Johansson P (1997). "Count Data Regression Using Series Expansions: With Applications." *Journal of Applied Econometrics*, **12**, 203–223. `doi:10.1002/(sici) 1099-1255(199705)12:3<203::aid-jae446>3.0.co;2-2`.

Cameron AC, Trivedi PK (2013). *Regression Analysis of Count Data*. 2nd edition. Cambridge University Press.

Canty A, Ripley BD (2017). **boot**: *Bootstrap R (S-PLUS) Functions*. R package version 1.3-20.

Cribari-Neto F, Zeileis A (2010). "Beta Regression in R." *Journal of Statistical Software*, **34**(2), 1–24. URL `http://www.jstatsoft.org/v34/i02/`.

Faddy MJ (1997). "Extended Poisson Process Modelling and Analysis of Count Data." *Biometrical Journal*, **39**(4), 431–440. `doi:10.1002/bimj.4710390405`.

Faddy MJ, Smith DM (2005). "Modelling the Dependence between the Number of Trials and the Success Probability in Binary Trials." *Biometrics*, **61**(4), 1112–1114. `doi:10.1111/j. 1541-0420.2005.00466.x`.

Faddy MJ, Smith DM (2008). "Extended Poisson Process Modelling of Dilution Series Data." *Applied Statistics*, **57**(4), 461–471. `doi:10.1111/j.1467-9876.2008.00622.x`.

Faddy MJ, Smith DM (2011). "Analysis of Count Data with Covariate Dependence in Both Mean and Variance." *Journal of Applied Statistics*, **38**(12), 2683–2694. `doi:10.1080/ 02664763.2011.567250`.

Faddy MJ, Smith DM (2012). "Extended Poisson Process Modeling and Analysis of Grouped Binary Data." *Biometrical Journal*, **54**(3), 426–435. `doi:10.1002/bimj.201100214`.

Gilbert P, Varadhan R (2015). **numDeriv**: *Accurate Numerical Derivatives*. R package version 2014.2-1, URL `https://CRAN.R-project/package=numDeriv`.

Goulet V, Dutang C, Maechler M, Firth D, Shapira M, Stadelmann M (2015). **expm**: *Matrix Exponential*. R package version 0.999.1, URL `https://CRAN.R-project/package=expm`.

Grün B, Kosmidis I, Zeileis A (2012). "Extended Beta Regression in R: Shaken, Stirred, Mixed, and Partitioned." *Journal of Statistical Software*, **48**(11), 1–25. URL `http://www. jstatsoft.org/v48/i11/`.

Grunwald GK, Bruce SL, Jiang L, Strand M, Rabinovitch N (2011). "A Statistical Model for Under- or Overdispersed Clustered and Longitudinal Count Data." *Biometrical Journal*, **53**(4), 578–594. `doi:10.1002/bimj.201000076`.

Hilbe J, Robinson A (2014). **msme**: *Functions and Datasets for 'Methods of Statistical Model Estimation'*. R package version 0.5.1, URL https://CRAN.R-project.org/package=msme.

Hilbe JM (2011). *Negative Binomial Regression*. 2nd edition. Cambridge University Press.

Hilbe JM (2014). *Modeling Count Data*. Cambridge University Press.

Lüning KG, Sheridan W, Ytterborn KH, Gullberg U (1966). "The Relationship Between the Number of Implantations and the Rate of Intra-Uterine Death in Mice." *Mutation Research*, **3**, 444–451. doi:10.1016/0027-5107(66)90054-6.

Nash JC (2014). *Nonlinear Parameter Optimization Using R Tools*. John Wiley & Sons. doi:10.1002/9781118884003.

Nelder JA, Mead R (1967). "A Simplex Method for Function Minimisation." *The Computer Journal*, **7**(4), 308–313. doi:10.1093/comjnl/7.4.308.

R Core Team (2016). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL https://www.R-project.org/.

Sáez-Castillo AJ, Conde-Sánchez A (2013). "A Hyper-Poisson Regression Model for Overdispersed and Underdispersed Count Data." *Computational Statistics & Data Analysis*, **61**, 148–157. doi:10.1016/j.csda.2012.12.009.

Smith DM, Faddy MJ (2016a). "Mean and Variance Modeling of Under- and Overdispersed Count Data." *Journal of Statistical Software*, **69**(6), 1–23. URL http://www.jstatsoft.org/v69/i06/.

Smith DM, Faddy MJ (2016b). **CountsEPPM**: *Fitting of EPPM Models to Count Data*. R package version 2.1, URL https://CRAN.R-project/package=CountsEPPM.

Smith DM, Faddy MJ (2018). **BinaryEPPM**: *Mean and Variance Modeling of Binary Data*. R package version 2.2, URL https://CRAN.R-project.org/package=BinaryEPPM.

Zeileis A, Croissant Y (2010). "Extended Model Formulas in R: Multiple Parts and Multiple Responses." *Journal of Statistical Software*, **34**(1), 1–13. doi:10.18637/jss.v034.i01.

Zeileis A, Hothorn T (2002). "Diagnostic Checking in Regression Relationships." *R News*, **2**(3), 7–10. URL http://CRAN.R-project.org/doc/Rnews/.

Zhu J, Eickhoff JC, Kaiser MS (2003). "Modelling the Dependence between Number of Trials and Success Probability in Beta-Binomial-Poisson Mixture Distributions." *Biometrics*, **59**(4), 955–961. doi:10.1111/j.0006-341x.2003.00110.x.

**Affiliation:**

David M. Smith
IBM Watson Health
75 Binney Street
Cambridge, MA 02142-1123, United States of America
E-mail: smithdm1@us.ibm.com

Malcolm J. Faddy
School of Mathematical Sciences
Queensland University of Technology
G.P.O. Box 2434
Brisbane Qld 4001, Australia
E-mail: `m.faddy@qut.edu.au`