

1 parameters (p) cannot be used for genomic-enabled prediction where the number of parameters
2 (p) is larger than the sample size (n). Here we propose a Bayesian mixed negative binomial
3 (BMNB) genomic regression model for counts that takes into account genotype by environment
4 ($G \times E$) interaction. We also provide all the full conditional distributions to implement a Gibbs
5 sampler. We evaluated the proposed model using a simulated data set and a real wheat data set
6 from the International Maize and Wheat Improvement Center (CIMMYT) and collaborators.
7 Results indicate that our BMNB model is a viable alternative for analyzing count data.

8

9 **Keyword:** Bayesian model; Count data; Genome enabled prediction; Gibbs sampler.

10

11

Introduction

12 A phenotype is the result of genotype (G), environment (E) and the genotype by environment
13 interactions ($G \times E$) in most living organisms. Garrod (1902) observed that the effect of genes
14 on phenotype could be modified by the environment (E). Similarly, Turesson (1922)
15 demonstrated that the development of a plant is often influenced by its surroundings. He
16 postulated the existence of a close relationship between crop plant varieties and their
17 environment, and stressed that the presence of a particular variety in a given locality is not just a
18 chance occurrence; rather, there is a genetic component that helps the individual adapt to that
19 area.

20 For these reasons, today the consensus is that $G \times E$ is useful for understanding genetic
21 heterogeneity under different environmental exposures (Kraft *et al.*, 2007; Van Os and Rutten,
22 2009) and for identifying high-risk or productive subgroups in a population (Murcay *et al.*,

1 2009); it also provides insight into the biological mechanisms of complex traits such as disease
2 resistance and yield (Thomas, 2011), and improves the ability to discover resistance genes that
3 interact with other factors that have little marginal effects (Thomas, 2011). However, finding
4 significant $G \times E$ interactions is challenging. Model misspecification, inconsistent definition of
5 environmental variables, and insufficient sample sizes are just a few of the issues that often lead
6 to low-power and non-reproducible findings in $G \times E$ studies (Jiao *et al.*, 2013; Winham and
7 Biernacka, 2013).

8 Genomics and its breeding applications are developing very quickly with the goal of
9 predicting yet-to-be observed phenotypes or unobserved genetic values for complex traits and
10 inferring the underlying genetic architecture utilizing large collections of markers (Goddard and
11 Hayes, 2009; Zhang *et al.*, 2014). Also, genomics is useful when dealing with complex traits
12 that are multi-genic in nature and have major environmental influence (Perez-de-Castro *et al.*,
13 2012). For these reasons, the use of whole genome prediction models continues to increase. In
14 genomic prediction, all marker effects are fitted simultaneously on a model and simulation
15 studies promote the use of this methodology to increase genetic progress in less time. For
16 continuous phenotypes, models have been developed to regress phenotypes on all available
17 markers using a linear model (Goddard and Hayes, 2009; de los Campos *et al.*, 2013). However,
18 in plant breeding, the response variable in many traits is a count ($y=0,1,2,\dots$), for example,
19 number of panicle per plant, number of seed per panicle, weed count per plot, etc. Count data
20 are discrete, non-negative, integer-valued, and typically have right-skewed distributions
21 (Yaacob *et al.*, 2010).

22 Poisson regression and negative binomial regression are often used to deal with count
23 data. These models have a number of advantages over an ordinary linear regression model,

1 including a skewed, discrete distribution $(0,1,2,3,\dots)$ and the restriction of predicted values for
2 phenotypes to non-negative numbers (Yaacob *et al.*, 2010). These models are different from an
3 ordinary linear regression model. First, they do not assume that counts follow a normal
4 distribution. Second, rather than modeling y as a linear function of the regression coefficients,
5 they model a function of the response mean as a linear function of the coefficients (Cameron
6 and Trivedi, 1986). Regression models for counts are usually nonlinear and have to take into
7 consideration the specific properties of counts, including discreteness and non-negativity, and
8 are often characterized by overdispersion (variance greater than the mean) (Zhou *et al.*, 2012).

9 However, in the context of genomic selection, it is still common practice to apply linear
10 regression models to these data or to transformed data (Montesinos-López *et al.*, 2015a,b). This
11 does not take into account that: (a) many distributions of count data are positively skewed, many
12 observations in the data set have a value of 0, and the high number of 0's in the data set does not
13 allow a skewed distribution to be transformed into a normal one (Yaacob *et al.*, 2010); and (b) it
14 is quite likely that the regression model will produce negative predicted values, which are
15 theoretically impossible (Yaacob *et al.*, 2010; Stroup, 2015). When transformation is used, it is
16 not always possible to have normally distributed data and many times transformations not only
17 do not help, they are counterproductive. There is also mounting evidence that transformations
18 do more harm than good for the models required by the vast majority of contemporary plant and
19 soil science researchers (Stroup, 2015). To the best of our knowledge, only the paper of
20 Montesinos-López *et al.* (2015c) is appropriate for genomic prediction for count data under a
21 Bayesian framework; however it does not take into account $G \times E$ interaction.

22 In this paper, we extend the NB regression model for counts proposed by Montesinos-
23 López *et al.* (2015c) to take into account $G \times E$ by using a data augmentation approach. A

1 Gibbs sampler was derived since all full conditional distributions were obtained, which allows
2 drawing samples from them to estimate the required parameters. In addition, we provide all the
3 details of the efficient derived Gibbs sampler so it can be easily implemented by most plant and
4 animal scientists. We illustrate our proposed methods with a simulated data set and a real data
5 set on wheat Fusarium head blight. We compare our proposed models (NB and Poisson) with
6 the Normal and Log-Normal models that are commonly implemented for analyzing count data.
7 We also provide R code for implementing the proposed models.

8

9

Materials and Methods

10 The data used in this study were taken from a Ph.D. thesis (Falconi-Castillo, 2014) aimed at
11 identifying sources of resistance to Fusarium head blight (FHB), caused by *Fusarium*
12 *graminearum* and identify genomic regions and molecular markers linked to FHB resistance
13 through association analysis.

14

15

Experimental data

16 Phenotypic data

17 A total of 297 spring wheat lines developed by the International Maize and Wheat Improvement
18 Center (CIMMYT) was assembled and evaluated for resistance to *F. graminearum* in México over
19 two years (2012 and 2014) and Ecuador for one year (2014). In this paper we used only 182 spring
20 wheat lines since only for these lines we have complete marker information.

21 Genotypic data

1 DNA samples were genotyped using an Illumina 9K SNP chip with 8,632 SNPs (Cavanagh *et al.*,
2 2013). SNP markers with unexpected genotype AB (heterozygous) were recoded as either AA or
3 BB based on the graphical interface visualization tool of the software GenomeStudio® (Illumina).
4 SNP markers that did not show clear clustering patterns were excluded. In addition, 66 simple
5 sequence repeats (SSR) markers were screened. After filtering the markers for the minor allele
6 frequency (MAF) of 0.05 and deleting markers with more than 10% of no calls, the final set of SNPs
7 was of 1,635 SNP.

8

9 **Data and software availability**

10 The phenotypic (FHB) and genotypic (marker) data used in this study as well as basic R codes
11 (R Core Team, 2015) for fitting the models can be directly downloaded from the repository at
12 <http://hdl.handle.net/11529/10575>

13

14

Statistical Models

15 We used y_{ijt} to represent the count response for the t th replication of the j th line in the
16 i th environment with $i = 1, \dots, I$; $j = 1, 2, \dots, J$, $t = 1, 2, \dots, n_{ij}$ and we propose the following
17 linear predictor that takes into account $G \times E$:

$$18 \quad \eta_{ij} = E_i + g_j + gE_{ij} \quad (1)$$

19 where E_i represents the environment i , g_j is the marker effect of genotype j , and gE_{ij} is the
20 interaction between markers and environment; $I = 3$, since we have three environments (Batan
21 2012, Batan 2014, and Chunchi 2014), $J = 182$, since it is the number of lines under study,
22 and n_{ij} represents the number of replicates of each line in each environment (the minimum and
23 maximum n_{ij} found per line were 10 and 20). The number of observations in each environment

1 i is $n_i = \sum_{j=1}^J n_{ij}$, while the total number of observations is $n = \sum_{i=1}^I n_i$. IJ is the product of
2 the number of environments and number of lines. Four models were implemented using the
3 linear predictor given in expression (1).

4

5 **Model NB**

6 Distributions: $y_{ijt}|g_j, gE_{ij} \sim \text{NB}(\mu_{ij}, r)$, with r being the scale parameter, $\mu_{ij} = \exp(\eta_{ij})$, $\mathbf{g} =$
7 $(g_1, \dots, g_J)^T \sim N(\mathbf{0}, \mathbf{G}_1 \sigma_g^2)$, $\mathbf{gE}_i = (gE_{i1}, \dots, gE_{iJ})^T \sim N(\mathbf{0}, \mathbf{G}_2 \sigma_{gE}^2)$. Note that the NB
8 distribution has expected value μ_{ij} and is smaller than the variance $\mu_{ij} + \frac{\mu_{ij}^2}{r}$. \mathbf{G}_1 and \mathbf{G}_2 were
9 assumed known, with \mathbf{G}_1 computed from marker \mathbf{X} data (for $k = 1, \dots, p$ markers) as
10 $\mathbf{G}_1 = \frac{\mathbf{X}\mathbf{X}^T}{p}$; this matrix is called the Genomic Relationship Matrix (GRM) (VanRaden, 2008).

11 While \mathbf{G}_2 is computed as $\mathbf{G}_2 = \mathbf{I}_I \otimes \mathbf{G}_1$ of order $IJ \times IJ$ and \otimes denotes the Kronecker product,
12 \mathbf{I}_I means that we assume independence between environments.

13

14 **Model Pois**

15 This model is the same as **Model NB**, except that $y_{ijt}|g_j, gE_{ij} \sim \text{Poisson}(\mu_{ij})$. Since according
16 to Zhou *et al.* (2012) and Teerapabolarn and Jaioun (2014) the $\lim_{r \rightarrow \infty} \text{NB}(\mu_{ij}, r) =$
17 $\text{Pois}(\mu_{ij})$, **Model Pois** was implemented using the same method as **Model NB**, but fixing r to
18 a large value, depending on the mean count. We used $r = 1000$, which is a good choice when
19 the mean count is less than 100.

20

21 **Model Normal**

22 Model Normal is similar to **Model NB**, except that $y_{ijt}|g_j, gE_{ij} \sim N(\eta_{ij}, \sigma_e^2)$ with identity link

1 function.

2

3 **Model Log-Normal**

4 Model Log-Normal is similar to **Model NB**, except that $\log(y_{ijt} + 1) | g_j, gE_{ij} \sim N(\eta_{ij}, \sigma_e^2)$ with
5 identity link function.

6

7 When $p > n$, implementing **Models NB and Pois** is challenging. For this reason, we propose a
8 Bayesian method for dealing with situations when $p > n$. The **Models Normal and Log-**
9 **Normal** were implemented in the package BGLR of de los campos *et al.* (2014).

10

11 *Bayesian mixed negative binomial regression*

12 Rewriting the linear predictor (1) as $\eta_{ij} = \mathbf{x}_i^T \boldsymbol{\beta} + \sum_{h=1}^2 b_{hij}$, with $\mathbf{x}_i^T = [x_{i1}, x_{i2}, x_{i3}]$,

13 where x_{ik} is an indicator variable that takes the value of 1 if it is observed in environment i and

14 0 otherwise, for $k = 1, 2, 3$; $\boldsymbol{\beta}^T = [\beta_1, \beta_2, \beta_3]$, since three is the number of environments under

15 study, $b_{1ij} = g_j$ and $b_{2ij} = gE_{ij}$. Note that in a sequence of independent Bernoulli (π_{ij}) trials,

16 the random variable y_{ijt} denotes the number of successes before the r th failure occurs. Then

$$17 \Pr(Y_{ijt} = y_{ijt} | g_j, gE_{ij}) = \binom{y_{ijt} + r - 1}{y_{ijt}} \left(1 - \frac{\mu_{ij}}{r + \mu_{ij}}\right)^r \left(\frac{\mu_{ij}}{r + \mu_{ij}}\right)^{y_{ijt}} \text{ for } y_{ijt} = 0, 1, 2, \dots$$

$$18 = \frac{\Gamma(y_{ijt} + r)}{y_{ijt}! \Gamma(r)} \frac{[\exp(\eta_{ij}^*)]^{y_{ijt}}}{[1 + \exp(\eta_{ij}^*)]^{y_{ijt} + r}}, \quad y_{ijt} = 0, 1, 2, \dots \quad (2)$$

19 Since $\pi_{ij} = \frac{\mu_{ij}}{r + \mu_{ij}} = \frac{r\mu_{ij}}{1 + r\mu_{ij}} = \frac{r\exp(\eta_{ij})}{1 + r\exp(\eta_{ij})} = \frac{\exp(\eta_{ij}^*)}{1 + \exp(\eta_{ij}^*)}$, where $\eta_{ij}^* = \mathbf{x}_i^T \boldsymbol{\beta}^* + \sum_{h=1}^2 b_{hij}$, $\boldsymbol{\beta}^* =$

20 $[\beta_1^*, \beta_2^*, \beta_3^*]$, with $\beta_i^* = \beta_i - \log(r)$ since \mathbf{x}_i^T is composed of three indicator variables. We

21 can rewrite (Eq 2) as:

$$\Pr(Y_{ijt} = y_{ijt} | g_j, gE_{ij}) = \frac{\Gamma(y_{ijt}+r)}{y_{ijt}!\Gamma(r)} 2^{-y_{ijt}-r} \exp\left(\frac{y_{ijt}-r}{2} \eta_{ij}^*\right) \int_0^\infty \exp\left[-\frac{\omega_{ijt}(\eta_{ij}^*)^2}{2}\right] f(\omega_{ijt}, y_{ijt} + r, 0) d\omega_{ijt} \quad (3)$$

Expression (3) was obtained using the equality given by Polson *et al.* (2013): $\frac{(e^\psi)^a}{(1+e^\psi)^b} = 2^{-b} e^{\kappa\psi} \int_0^\infty e^{-\frac{\omega\psi^2}{2}} f(\omega; b, 0) d\omega$, where $\kappa = a - b/2$ and $f(\cdot, b, 0)$ denotes the density of $PG(b, c = 0)$, the PG Pólya-Gamma distribution with parameters b and $c = 0$ (see Definition 1 in Polson *et al.*, 2013).

From here, conditional on $\omega_{ijt} \sim PG(y_{ijt} + r, c = 0)$,

$$\Pr(Y_{ijt} = y_{ijt} | g_j, gE_{ij}, \omega_{ijt}) = \frac{\Gamma(y_{ijt}+r)}{y_{ijt}!\Gamma(r)} 2^{-y_{ijt}-r} \exp\left(\frac{y_{ijt}-r}{2} \eta_{ij}^*\right) \exp\left[-\omega_{ijt}(\eta_{ij}^*)^2/2\right] \quad (4)$$

To be able to get the full conditional distributions, we provide the prior distributions, $f(\boldsymbol{\theta})$, for all the unknown model parameters $\boldsymbol{\theta} = (\boldsymbol{\beta}^*, \sigma_\beta^2, \mathbf{b}_1, \sigma_{b1}^2, \mathbf{b}_2, \sigma_{b2}^2, r)$. We assume prior independence between the parameters, that is,

$$f(\boldsymbol{\theta}) = f(\boldsymbol{\beta}^*)f(\sigma_\beta^2)f(\mathbf{b}_1)f(\sigma_{b1}^2)f(\mathbf{b}_2)f(\sigma_{b2}^2)f(r).$$

We assign conditionally conjugate but weakly informative prior distributions to the parameters because we have no prior information. Prior specification in terms of $\boldsymbol{\beta}^*$ instead of $\boldsymbol{\beta}$ is for convenience. We adopt proper priors with known hyper-parameters whose values we specify in model implementation to guarantee proper posteriors. We assume that $\boldsymbol{\beta}^* | \sigma_\beta^2 \sim N_p(\boldsymbol{\beta}_0, \Sigma_0 \sigma_\beta^2)$, $\sigma_\beta^2 \sim \chi^{-2}(v_\beta, S_\beta)$ where $\chi^{-2}(v_\beta, S_\beta)$ denotes a scaled inverse chi-square distribution with shape v_β and scale S_β parameters, $\mathbf{b}_1 | \sigma_{b1}^2 \sim N_{nb1}(\mathbf{0}, \mathbf{G}_1 \sigma_{b1}^2)$, $\sigma_{b1}^2 \sim \chi^{-2}(v_{b1}, S_{b1})$, $\mathbf{b}_2 | \sigma_{b2}^2 \sim N_{nb2}(\mathbf{0}, \mathbf{G}_2 \sigma_{b2}^2)$, $\sigma_{b2}^2 \sim \chi^{-2}(v_{b2}, S_{b2})$ and $r \sim G(a_0, 1/b_0)$. Next we combine (Eq 4)

1 using all data with priors to get the full conditional distribution for parameters $\boldsymbol{\beta}^*$, σ_β^2 , \mathbf{b}_1 , $\sigma_{b_1}^2$,
 2 \mathbf{b}_2 , $\sigma_{b_2}^2$ and r .

3

4 *Full conditional distributions*

5 The full conditional distribution of $\boldsymbol{\beta}^*$ is given as:

$$6 \quad f(\boldsymbol{\beta}^* | \mathbf{y}, ELSE) \sim N(\tilde{\boldsymbol{\beta}}_0, \tilde{\boldsymbol{\Sigma}}_0) \quad (5)$$

7 where $\tilde{\boldsymbol{\Sigma}}_0 = (\boldsymbol{\Sigma}_0^{-1} \sigma_\beta^{-2} + \mathbf{X}^T \mathbf{D}_\omega \mathbf{X})^{-1}$, $\tilde{\boldsymbol{\beta}}_0 = \tilde{\boldsymbol{\Sigma}}_0 (\boldsymbol{\Sigma}_0^{-1} \sigma_\beta^{-2} \boldsymbol{\beta}_0 - \mathbf{X}^T \mathbf{D}_\omega \sum_{h=1}^2 \mathbf{Z}_h \mathbf{b}_h + \mathbf{X}^T \boldsymbol{\kappa})$,

8 $\mathbf{y}_{ij} = [y_{ij1}, \dots, y_{ijn_{ij}}]^T$, $\mathbf{y}_i = [\mathbf{y}_{i1}^T, \dots, \mathbf{y}_{ij}^T]^T$, $\mathbf{y} = [\mathbf{y}_1^T, \dots, \mathbf{y}_I^T]^T$, $\boldsymbol{\kappa}_{ij} = \frac{1}{2} [y_{ij1} - r, \dots, y_{ijn_{ij}} -$

9 $r]^T$, $\boldsymbol{\kappa}_i = [\boldsymbol{\kappa}_{i1}^T, \dots, \boldsymbol{\kappa}_{ij}^T]^T$, $\boldsymbol{\kappa} = [\boldsymbol{\kappa}_1^T, \dots, \boldsymbol{\kappa}_I^T]^T$, $\mathbf{X}_{ij} = [\mathbf{1}_{n_{ij}}^T \otimes \mathbf{x}_i]^T$, $\mathbf{X}_i = [\mathbf{X}_{i1}^T, \dots, \mathbf{X}_{ij}^T]^T$, $\mathbf{X} =$

10 $[\mathbf{X}_1^T, \dots, \mathbf{X}_I^T]^T$, $\mathbf{D}_{\omega ij} = \text{diag}(\omega_{ij1}, \dots, \omega_{ijn_{ij}})$, $\mathbf{D}_{\omega i} = \text{diag}(\mathbf{D}_{\omega i1}, \dots, \mathbf{D}_{\omega ij})$,

11 $\mathbf{D}_\omega = \text{diag}(\mathbf{D}_{\omega 1}, \dots, \mathbf{D}_{\omega I})$, $\mathbf{b}_{hi} = [b_{hi1}, \dots, b_{hij}]^T$, $\mathbf{b}_h = [\mathbf{b}_{h1}^T, \dots, \mathbf{b}_{hI}^T]^T$,

$$12 \quad \mathbf{Z}_{1i} = \begin{bmatrix} \mathbf{1}_{n_{i1}} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n_{i2}} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{1}_{n_{ij}} \end{bmatrix}, \mathbf{Z}_1 = [\mathbf{Z}_{11}^T, \dots, \mathbf{Z}_{1I}^T]^T \text{ and } \mathbf{Z}_2 = \mathbf{Z}_1 * \sim \mathbf{X}, \text{ where } * \sim \text{ indicates}$$

13 the horizontal Kronecker product between \mathbf{Z}_1 and \mathbf{X} . The horizontal Kronecker product

14 performs a Kronecker product of \mathbf{Z}_1 and \mathbf{X} and creates a new matrix by stacking these row

15 vectors into a matrix. \mathbf{Z}_1 and \mathbf{X} must have the same number of rows, which is also the same

16 number of rows in the result matrix. The number of columns in the result matrix is equal to the

17 product of the number of columns in \mathbf{Z}_1 and \mathbf{X} . When the prior for $\boldsymbol{\beta}^* \propto \text{constant}$, the posterior

18 distribution of $\boldsymbol{\beta}^*$ is also normally distributed, $N(\tilde{\boldsymbol{\beta}}_0, \tilde{\boldsymbol{\Sigma}}_0)$, but we set the term $\boldsymbol{\Sigma}_0^{-1} \sigma_\beta^{-2}$ to zero in

19 both $\tilde{\boldsymbol{\Sigma}}_0$ and $\tilde{\boldsymbol{\beta}}_0$.

1 The fully conditional distribution of ω_{ijt} is

$$2 \quad f(\omega_{ijt}|\mathbf{y}, ELSE) \sim PG(y_{ijt} + r, \mathbf{x}_i^T \boldsymbol{\beta}^* + \sum_{h=1}^2 b_{hij}) \quad (6)$$

3 Defining $\boldsymbol{\eta}^h = \mathbf{X} \boldsymbol{\beta}^* + \mathbf{Z}_h \mathbf{b}_h$, with $h = 1, 2$, the conditional distribution of \mathbf{b}_h is given as

$$4 \quad f(\mathbf{b}_h|\mathbf{y}, ELSE) \sim N(\tilde{\mathbf{b}}_h, \mathbf{F}_h) \quad (7)$$

5 If $\boldsymbol{\eta}^1 = \mathbf{X} \boldsymbol{\beta}^* + \mathbf{Z}_2 \mathbf{b}_2$, then $\mathbf{F}_1 = (\sigma_{b_1}^{-2} \mathbf{G}_1^{-1} + \mathbf{Z}_1^T \mathbf{D}_\omega \mathbf{Z}_1)^{-1}$, $\tilde{\mathbf{b}}_1 = \mathbf{F}_1 (\mathbf{Z}_1^T \boldsymbol{\kappa} - \mathbf{Z}_1^T \mathbf{D}_\omega \boldsymbol{\eta}^1)$ and

6 then $\mathbf{b}_1|\mathbf{y}, ELSE \sim N(\tilde{\mathbf{b}}_1, \mathbf{F}_1)$. Similarly, by defining $\boldsymbol{\eta}^2 = \mathbf{X} \boldsymbol{\beta}^* + \mathbf{Z}_1 \mathbf{b}_1$, we arrive at the full

7 conditional of \mathbf{b}_2 as $\mathbf{b}_2|\mathbf{y}, ELSE \sim N(\tilde{\mathbf{b}}_2, \mathbf{F}_2)$, where $\mathbf{F}_2 = (\sigma_{b_2}^{-2} \mathbf{G}_2^{-1} + \mathbf{Z}_2^T \mathbf{D}_\omega \mathbf{Z}_2)^{-1}$, $\tilde{\mathbf{b}}_2 =$

$$8 \quad \mathbf{F}_2 (\mathbf{Z}_2^T \boldsymbol{\kappa} - \mathbf{Z}_2^T \mathbf{D}_\omega \boldsymbol{\eta}^2).$$

9 The fully conditional distribution of $\sigma_{b_h}^2$, for $h = 1, 2$, is

$$10 \quad f(\sigma_{b_h}^2|\mathbf{y}, ELSE) \sim \chi^{-2}(\tilde{\nu}_b = \nu_{b_h} + n_{b_h}, \tilde{S}_b = (\mathbf{b}_h^T \mathbf{G}_h^{-1} \mathbf{b}_h + \nu_{b_h} S_{b_h}) / \nu_{b_h} + n_{b_h}) \quad (8)$$

11 with $n_{b_1} = J$ and $n_{b_2} = IJ$.

12 The conditional distribution of $\sigma_{\beta^*}^2$ is

$$13 \quad f(\sigma_{\beta^*}^2|\mathbf{y}, ELSE) \sim \chi^{-2}(\tilde{\nu}_{\beta^*} = \nu_{\beta^*} + I, \tilde{S}_{\beta^*} = [(\boldsymbol{\beta}^* - \boldsymbol{\beta}_0)^T \boldsymbol{\Sigma}_0^{-1} (\boldsymbol{\beta}^* - \boldsymbol{\beta}_0) + \nu_{\beta^*} S_{\beta^*}] / \nu_{\beta^*} + I) \quad (9)$$

14 Taking advantage of the fact that the NB distribution can also be generated using a

15 Poisson representation (Quenouille, 1949) as $Y = \sum_{l=1}^L u_l$, where $u_l \sim \text{Log}(\pi)$, $\pi = \frac{\mu}{r+\mu}$ and is

16 independent of $L \sim \text{Pois}(-r \log(1 - \pi))$, where Log and Pois denote logarithmic and Poisson

17 distributions, respectively. Then we infer a latent count L for each $Y \sim \text{NB}(\mu, r)$ conditional on

18 Y and r . Therefore, following Zhou *et al.* (2012), we obtain the full conditional of r by

19 alternating

$$1 \quad f(r|\mathbf{y}, ELSE) \sim G(a_0 - \sum_{i=1}^I \sum_{j=1}^J \sum_{t=1}^{n_{ij}} \log(1 - \pi_{ijt}), \frac{1}{b_0 + \sum_{i=1}^I \sum_{j=1}^J \sum_{t=1}^{n_{ij}} L_{ijt}}) \quad (10)$$

$$2 \quad f(L_{ijt}|\mathbf{y}, ELSE) \sim CRT(y_{ijt}, r) \quad (11)$$

3 where $CRT(y_{ijt}, r)$ denotes a Chinese restaurant table (CRT) count random variable that can be
4 generated as $L_{ijt} = \sum_{l=1}^{y_{ijt}} d_l$, where $d_l \sim \text{Bernoulli}\left(\frac{r}{l-1+r}\right)$. For details of the CRT random
5 variable derivation, see Zhou and Carin (2012, 2015).

6

7 **Gibbs sampler**

8 The Gibbs sampler for the latent parameters of the NB with $G \times E$ can be implemented by
9 sampling repeatedly from the following loop:

- 10 1. Sample ω_{ijt} values from the Pólya-Gamma distribution in (6).
- 11 2. Sample $L_{ijt} \sim CRT(y_{ijt}, r)$ from (11).
- 12 3. Sample the scale parameter (r) from the gamma distribution in (10).
- 13 4. Sample the location effects (β^*) from the normal distribution in (5).
- 14 5. Sample the random effects (\mathbf{b}_h) with $h = 1, 2$, from the normal distribution in (7).
- 15 6. Sample the variance effects ($\sigma_{b_h}^2$) with $h = 1, 2$, from the scaled inverted χ^2 distribution in
16 (8).
- 17 7. Sample the variance effect ($\sigma_{\beta^*}^2$) from the scaled inverted χ^2 distribution in (9).
- 18 8. Return to step 1 or terminate when chain length is adequate to meet convergence
19 diagnostics.

20

1 **Model implementation**

2 The Gibbs sampler described above for the BMNB model was implemented in R-Core
3 Team (2015). Implementation was done under a Bayesian approach using Markov Chain Monte
4 Carlo (MCMC) through the Gibbs sampler algorithm, which samples sequentially from the full
5 conditional distribution until it reaches a stationary process, converging with the joint posterior
6 distribution (Gelfand and Smith, 1990). To decrease the potential impact of MCMC errors on
7 prediction accuracy, we performed a total of 60,000 iterations with a burn-in of 30,000, so that
8 30,000 samples were used for inference. We did not apply thinning of the chains following the
9 suggestions of Geyer (1992), MacEachern and Berliner (1994) and Link and Eaton (2012), who
10 provide justification of the ban on subsampling MCMC output for approximating simple
11 features of the target distribution (e.g., means, variances, and percentiles). We implemented the
12 prior specification given in the section Bayesian mixed negative binomial regression with
13 $\boldsymbol{\beta}^* | \sigma_{\beta}^2 \sim N_p(\boldsymbol{\beta}_0 = \mathbf{0}_3^T, \mathbf{I}_3 \times 10,000)$, $\mathbf{b}_1 | \sigma_{b_1}^2 \sim N_{nb_1}(\mathbf{0}_{nb_1}^T, \mathbf{G}_1 \sigma_{b_1}^2)$, where \mathbf{G}_1 is the GRM, that is,
14 the covariance matrix of the random effects, $\sigma_{b_1}^2 \sim \chi^{-2}(v_{b_1} = 3, S_{b_1} = 0.001)$,
15 $\mathbf{b}_2 | \sigma_{b_2}^2 \sim N_{nb_2}(\mathbf{0}_{nb_2}^T, \mathbf{G}_2 \sigma_{b_2}^2)$, \mathbf{G}_2 is the covariance matrix of the random effects that belong to the
16 $G \times E$ term, $\sigma_{b_2}^2 \sim \chi^{-2}(v_{b_2} = 3, S_{b_2} = 0.001)$, and $r \sim G(a_0 = 0.01, 1/(b_0 = 0.01))$. All these
17 hyper-parameters were chosen to lead weakly informative priors. The convergence of the
18 MCMC chains was monitored using trace plots and autocorrelation functions. We also
19 conducted a sensitivity analysis on the use of the inverse gamma priors for the variance
20 components and we observed that the results are robust under different choices of priors.

21

22 **Assessing prediction accuracy**

23 We used cross-validation to compare the prediction accuracy of the proposed models

1 for count phenotypes. We implemented a 10-fold cross validation, that is, the data set was
 2 divided into 10 mutually exclusive subsets; each time we used 9 subsets for the training set and
 3 the remaining one for validation set. The training set was used to fit the model and the
 4 validation set was used to evaluate the prediction accuracy of the proposed models. To compare
 5 the prediction accuracy of the proposed models, we calculated the Spearman correlation (Cor)
 6 and the mean square error of prediction (MSEP), both calculated using the observed and
 7 predicted response variables of the validation set. Models with large absolute values of Cor
 8 indicate better prediction accuracy, while small MSEP indicate better prediction performance.
 9 The predicted observations, \hat{y}_{ij} , were calculated with M collected Gibbs samples after
 10 discarding those of the burn-in period. For **Models NB** and **Pois** the predicted values were
 11 calculated as $\hat{y}_{ij} = \frac{\sum_{s=1}^M \exp(x_{i1}\beta_1^{*(s)} + x_{i1}\beta_2^{*(s)} + x_{i1}\beta_3^{*(s)} + \log(\hat{r}^{(s)}) + \hat{g}_j^{(s)} + \widehat{gE}_{ij}^{(s)})}{s}$, where $\hat{r}^{(s)}$, $\beta_1^{*(s)}$, $\beta_2^{*(s)}$,
 12 $\beta_3^{*(s)}$ and $\hat{g}_j^{(s)}$ and $\widehat{gE}_{ij}^{(s)}$ are estimates of β_1^* , β_2^* , β_3^* , r , g_j and gE_{ij} , for line j in environment
 13 i obtained in the sth collected sample. For **Model Normal** as
 14 $\hat{y}_{ij} = \frac{\sum_{s=1}^M (x_{i1}\hat{\beta}_1^{(s)} + x_{i1}\hat{\beta}_2^{(s)} + x_{i1}\hat{\beta}_3^{(s)} + \hat{g}_j^{(s)} + \widehat{gE}_{ij}^{(s)})}{s}$ and for **Model LN** the predicted observations were
 15 calculated as $\hat{y}_{ij} = \frac{\sum_{s=1}^M \exp(x_{i1}\hat{\beta}_1^{(s)} + x_{i1}\hat{\beta}_2^{(s)} + x_{i1}\hat{\beta}_3^{(s)} + \hat{g}_j^{(s)} + \widehat{gE}_{ij}^{(s)} + \frac{\hat{\sigma}_e^2(s)}{2})}{s} - 1$, using the corresponding
 16 estimates of each model.

17

18 Simulation study

19 To show the performance of the proposed Gibbs sampler for count phenotypes that takes
 20 into account $G \times E$, we performed a simulation study under model (1) in two scenarios (S1 and
 21 S2). Scenario 1 had three environments ($I = 3$), 20 genotypes ($J = 20$), $\mathbf{G}_1 = \mathbf{I}_{60}$, $\mathbf{G}_2 = \mathbf{I}_I \otimes$
 22 \mathbf{G}_1 and $\sigma_{b_1}^2 = \sigma_{b_2}^2 = 0.5$, with four different numbers of replicates of each genotype in each

1 environment, $n_{ij} = 5, 10, 20$ and 40 . Scenario 2 is equal to scenario 1, except that $\mathbf{G}_2 =$
2 $0.7\mathbf{I}_{60} + 0.3\mathbf{J}_{60}$, where \mathbf{J}_{60} is a square matrix of ones of order 60×60 . In this second scenario,
3 we imitated the correlation between lines of real data available in genomic selection. The priors
4 used for the simulation study in both scenarios (S1 and S2) were approximately flat for all
5 parameters: for $\boldsymbol{\beta}|\sigma_{\beta}^2 \sim N(\boldsymbol{\beta}_0^T = [0,0,0], \mathbf{I}_3 \times 10000)$, for $r \sim G(0.001, 1/0.001)$, for $\sigma_{b_1}^2$ and
6 $\sigma_{b_2}^2$ a $\sim \chi^{-2}(0.50002, 4.0002)$, while for $\mathbf{b}_1|\sigma_{b_1}^2 \sim N(\mathbf{0}, \mathbf{G}_1)$, and for $\mathbf{b}_2|\sigma_{b_2}^2 \sim N(\mathbf{0}, \mathbf{G}_2)$. We
7 computed 20,000 MCMC samples; Bayes estimates were computed with 10,000 samples since
8 the first 10,000 were discarded as burning. We report average estimates obtained by using the
9 proposed Gibbs sampler along with standard deviations (SD) (Table 1). All the results in Table
10 1 are based on 50 replications.

11

12

Results

13

14

15

16

17

18

19

20

21

22

Given in Table 1 are the results of the simulation study in both scenarios (S1 and S2). The bias when estimating the parameters is a little larger in S1 compared to S2. Also, parameter β_0 is the parameter with larger bias (underestimated). Both variances (σ_1^2, σ_2^2) are overestimated in scenario 1, but only σ_1^2 is overestimated in scenario 2. Also, with a sample size of $n_{ij} = 5$, parameter r had a larger SD; however, for larger sample sizes ($n_{ij} = 20, 40$), the SD were considerably reduced. In general, there was not a large reduction in SD when the sample size increased from 5 to 10, 20 and 40, the exception being the estimation of r in both scenarios and the estimation of β_0 in scenario 1, where there was a large reduction in SD when the sample size increased. Although estimations do not totally agree with the true values of the parameters, the proposed Gibbs sampler for count data that takes into account $G \times E$ did a good job of

1 estimating the parameters, since the estimates are close to the true values with a SD of
2 reasonable size.

3 Using the real data set, we compared four scenarios (given in Table 2) for each model.
4 Table 2 shows that in the linear predictor, scenarios 1 and 2 do not take into account interaction
5 effects, only main effects. Also, scenarios 1 and 3 do not use marker information. These four
6 scenarios were studied to investigate the gain in model fit and prediction ability taking into
7 account the interaction effects and using the marker information available.

8 The posterior means (Mean), posterior standard deviation (SD) of the scalar parameters,
9 and posterior predictive checks for each scenario of the proposed models are given in Table 3.
10 For the four models, the posterior means of the beta regression coefficients, variance
11 components, and over-dispersion parameters (r) are similar between scenarios 1 and 2 and
12 between scenarios 3 and 4. In terms of goodness of fit measured by the loglikelihood posterior
13 mean (loglink), the scenarios rank as follows: scenario 3, rank 1; scenario 4, rank 2; scenario 1,
14 rank 3; and scenario 2, rank 4, for the four proposed models, with the exception of **Model Pois**
15 where the ranking was scenario 3, rank 1; scenario 4, rank 2; scenario 2, rank 3; and scenario 1,
16 rank 4. Therefore, there is evidence that with the four proposed models in terms of goodness of
17 fit, the best scenario is S3. Of the four models under study, Table 3 shows that **Model LN**
18 reports the best fit since it has the largest Loglik.

19 In Table 4 we present the mean and standard deviation of the posterior predictive checks
20 (Cor and MSE) for each location (Batan 2012, Batan 2014 and Chunchi 2014) resulting from
21 the 10-fold cross-validation implemented for the four models and four scenarios. The predictive
22 checks given in Table 4 were calculated using the testing set. In **Model NB**, according to the
23 Spearman Correlation, the ranking of scenarios was as follows: in Batan 2012 and Batan 2014,

1 1 for scenario 4, 2 for scenario 3, 3 for scenario 1, and 4 for scenario 2. In Chunchi 2014, the
2 ranking was 1 for scenario 3, 2 for scenario 2, 3 for scenario 4, and 4 for scenario 4. With the
3 MSEP, the ranking for **Model NB** in Batan 2012 was 1 for scenario 3, 2 for scenario 4, 3 for
4 scenario 1, and 4 for scenario 2. In Batan 2014, the ranking was 1 for scenario 2, 2 for scenario
5 1, 3 for scenario 3, and 4 for scenario 4. In Chunchi 2014, the ranking in terms of MSEP was 1
6 for scenario 3, 2 for scenario 2, 3 for scenario 4, and 4 for scenario 1. Under **Model Pois**, the
7 ranking of the 4 scenarios in each locality was exactly the same as the ranking reported for
8 **Model NB**. For **Model Normal** in terms of the Spearman correlation, scenario 1 was the best in
9 prediction accuracy in Batan 2012 and Chunchi 2014, while scenario 4 was the worst in all three
10 locations. In terms of MSEP, the best scenario was 3 in Batan 2014 and Chunchi 2014, and the
11 worst was scenario 4 in Batan 2014 and Chunchi 2014. For **Model LN** in terms of the Spearman
12 correlation, the best scenarios were scenarios 1 and 2, and the worst was scenario 3 in Batan
13 2012. In Batan 2014, the best scenario was 1, then scenario 3 and the worst was scenario 4. In
14 Chunchi 2014, the best scenario was scenario 3, then scenario 2 and the worst was scenario 2. In
15 terms of MSEP for Batan 2012, the best scenario was 3, then scenario 1 and the worst was
16 scenario 4. In Batan 2014, the best scenario was 1, then 2 and the worst was scenario 4. Finally,
17 in Chunchi 2014, the best scenario was 3, then 2 and the worst was scenario 1.

18 Table 5 gives the average of the ranks of the two posterior predictive checks (Cor and
19 MSEP) that were used. Since we are comparing four scenarios for each model, the values of the
20 ranks range from 1 to 4, and the lower the values, the better the scenario. For ties we assigned
21 the average of the ranges that would have been assigned had there been no ties. Table 5 shows
22 that the best scenarios were scenarios 3 and 4 under **Model NB** and **Pois** in Batan 2012. In
23 Batan 2014 under **Models NB** and **Pois**, the best scenario was 3, while in Chunchi 2014, the

1 best scenarios were 3 and 1. Under **Model Normal**, the best scenario was scenario 3 in Batan
2 2014 and Chuchi 2014, while in Batan 2012, the best scenarios were 2 and 3. Finally, under
3 **Model LN**, the best scenario was 3 in Chunchi 2014, and scenario 1 in Batan 2012 and Batan
4 2014.

5 Results in Tables 4 and 5 indicate that the best models in terms of prediction accuracy
6 are **Models NB** and **Pois**, since they had better predictions in the validation set based on both
7 the posterior predictive checks (Cor and MSEP) implemented, although in terms of goodness of
8 fit, **Model LN** was the best. These results are in partial agreement with the findings of
9 Montesinos-Lopez *et al.* (2015), who came to the conclusion that **Models NB** and **Pois** are good
10 alternatives for modeling count data, although in this study, the best predictions were produced
11 by **Model LN**. However, this model did not take into account the $G \times E$ interaction.

12

13

Discussion

14 Developing specific methods for count data for genome-enabled prediction can help to
15 improve the selection of candidate genotypes early in time when the phenotypes are counts.
16 However, currently in genomic selection, phenotypic data (dependent variable) are not taken
17 into account before deciding on the modeling approach to be used, mainly due to the lack of
18 genome-enabled prediction models for non-normal phenotypes. The Bayesian regression models
19 proposed in this paper aim to fill this lack of genome-enabled prediction models for non-normal
20 data.

21 The first advantage of our proposed methods for count data is that they take into account
22 the nonlinear relationship between responses and consider the specific properties of counts,
23 including discreteness, non-negativity, and over-dispersion (variance greater than the mean);

1 this guarantees that the predictive response will not be negative, which makes no sense for count
2 data. In addition, our methods take into account $G \times E$, which plays a central role when
3 selecting candidates genotypes in plant breeding.

4 Another advantage of our proposed method is that the proposed Gibbs sampler has an
5 analytical solution since we were able to obtain all the full conditional distributions required
6 analytically. This was possible because we constructed our Gibbs sampler using the data
7 augmentation approach proposed by Polson *et al.*, (2013) for count data. For this reason, we
8 believe it is an attractive alternative for fitting complex multilevel data for counts because, in
9 addition to its simplicity, it can generate samples from a high dimensional probability
10 distribution.

11 Our proposed methods showed superior performance in terms of prediction accuracy
12 compared to **Models Normal** and **Log-Normal**. Also, we observed that in **Models NB** and **Pois**
13 taking into account the $G \times E$ increase considerable the prediction accuracy which is expected
14 since there is enough scientific evidence that including the $G \times E$ interaction improve prediction
15 accuracy. Finally, more research is needed to study the proposed methods using real data sets
16 and to extend the proposed genomic-enabled prediction models to deal with so many zeros in
17 count response variables and for modeling multiple traits.

18

19

Acknowledgments

20 We very much appreciate CIMMYT field collaborator, laboratory assistants, and technicians
21 who collected the phenotypic and genotypic data used in this study.

22

1

References

- 2 Cameron, A.C., and P.K. Trivedi. (1986). Econometric models based on count data.
3 Comparisons and applications of some estimators and tests. *Journal of Applied*
4 *Econometrics* 1(1): 29-53.
- 5 Cavanagh, C.R., Chao, S., Wang, S. *et al.* (2013) Genome-wide comparative diversity uncovers
6 multiple targets of selection for improvement in hexaploid wheat landraces and cultivars.
7 *Proceedings of the National Academy of Sciences.* 110(20) 8057-8062.
- 8 de los Campos, G., A.I. Vazquez, R. Fernando, Y.C. Klimentidis, and D. Sorensen. (2013).
9 Prediction of complex human traits using the genomic best linear unbiased predictor. *PLoS*
10 *Genetics* 9: e1003608.
- 11 de los Campos G, Pataki A, Pérez P (2014). The BGLR (Bayesian Generalized Linear
12 Regression) R-Package. [<http://bglr.r-forge.r-project.org/BGLR-tutorial.pdf>
- 13 Falconi-Castillo, E. (2014). Association mapping for detecting QTLs for Fusarium head blight
14 and Yellow rust resistance in bread wheat. Ph.D. Dissertation. Michigan State University,
15 East Lansing, Michigan, USA.
- 16 Garrod, A.E. (1902). The incidence of alkatonuria: a study in chemical individuality. *Lancet*
17 160:1616-1620.
- 18 Gelfand, A.E., and A.F. Smith. (1990). Sampling-based approaches to calculating marginal
19 densities. *J. Am. Statist. Assoc.* 85: 398-409.
- 20 Geyer, C.J. (1992). Practical Markov Chain Monte Carlo. *Stat Sci* 473-483.
- 21 Goddard, M.E., and B.J. Hayes. (2009). Mapping genes for complex traits in domestic animals
22 and their use in breeding programmes. *Nat. Rev. Genet.* 10: 381-391.
- 23 Jiao, S, L. Hsu, S. Bézieau, H. Brenner, A.T. Chan, *et al.* (2013). SBERIA: Set Based Gene-
24 Environment Interaction test for rare and common variants in complex diseases. *Genet.*
25 *Epidemiol.* 37:452-64.
- 26 Kraft, P., Y.C. Yen, D.O. Stram, J. Morrison, and W.J. Gauderman. (2007). Exploiting gene
27 environment interaction to detect genetic associations. *Hum. Hered.* 63: 111-119.

- 1 Link, W.A., and M.J. Eaton (2012). On thinning of chains in MCMC. *Methods Ecol Evol* 3:
2 112-115.
- 3 MacEachern, S.N., and L.M. Berliner. (1994). Subsampling the Gibbs sampler. *Am. Stat.*
4 48:188-190.
- 5 Montesinos-López, O.A., A. Montesinos-López, P. Pérez-Rodríguez, G. de los Campos, K.M.
6 Eskridge *et al.* (2015a). Threshold models for genome-enabled prediction of ordinal
7 categorical traits in plant breeding. *G3|Genes|Genomes|Genetics* 5: 1-10.
- 8 Montesinos-López, O.A., A. Montesinos-López, J. Crossa, J. Burgueño, and K.M. Eskridge.
9 (2015b). Genomic-enabled prediction of ordinal data with Bayesian logistic ordinal
10 regression. *G3|Genes|Genomes|Genetics* 5: 2113-2126.
- 11 Montesinos-López, O.A., A. Montesinos-López, P. Pérez-Rodríguez, K.M. Eskridge, X. He, P.
12 Juliana, P. Singh, and J. Crossa. (2015c). Genomic prediction models for count data.
13 *Journal of Agricultural, Biological, and Environmental Statistics (JABES)*. 20: 533–554,
14 DOI: 10.1007/s13253-015-0223-4.
- 15 Murcray, C.E., J.P. Lewinger, and W.J. Gauderman. (2009). Gene-environment interaction in
16 genome-wide association studies. *Am. J. Epidemiol.* 169: 219-226.
- 17 Pérez-de-Castro, A.M., S. Vilanova, J. Cañizares, L. Pascual, J.M. Blanca, M.J. Diez, and B.
18 Picó. (2012). Application of genomic tools in plant breeding. *Current Genomics* 13(3): 179.
- 19 Polson, N.G., J.G. Scott, and J. Windle. (2013). Bayesian inference for logistic models using
20 Pólya–Gamma latent variables. *J. Am. Statist. Assoc.* 108: 1339-1349.
- 21 Quenouille, M.H. (1949). A relation between the logarithmic, Poisson, and negative binomial
22 series. *Biometrics* 5: 162-164.
- 23 R Core Team (2015). R: A language and environment for statistical computing. R Foundation
24 for Statistical Computing. Vienna. Austria. ISBN 3-900051-07-0. URL [http://www.R-](http://www.R-project.org/)
25 [project.org/](http://www.R-project.org/).
- 26 Stroup, W.W. (2015). Rethinking the analysis of non-Normal data in plant and soil science.
27 *Agron. J.* 107: 811-827.

- 1 Teerapabolarn, K., and K. Jaioun. (2014). An improved Poisson approximation for the Negative
2 binomial bistribution. *Applied Mathematical Sciences* 8(89): 4441-4445.
- 3 Thomas, D. (2011). Response to ‘Gene-by-environment experiments: a new approach to finding
4 the missing heritability’ by Van Ijzendoorn *et al.* *Nature Rev. Genet.* 12: 881.
- 5 Turesson, G. (1922). The genotypical response of the plant species to the habitat. *Hereditas* 3:
6 211-350.
- 7 Van Os, J., and B. Rutten (2009). Gene-environment-wide interaction studies in psychiatry. *Am*
8 *J Psychiatry* 166: 964-966.
- 9 VanRaden, P.M. (2008). Efficient methods to compute genomic predictions. *J. Dairy Sci.* 91:
10 4414-4423.
- 11 Windle, J., C.M. Carvalho, J.G. Scott, L. Sun. (2013). Polya-Gamma Data Augmentation for
12 Dynamic Models. arXiv preprint arXiv:1308.0774.
- 13 Winham, S.J., and J.M. Biernacka. (2013). Gene–environment interactions in genome-wide
14 association studies: current approaches and new directions. *J. Child Psychol. Psychiatry* 54:
15 1120-1134.
- 16 Yaacob, W.F.W., M.A. Lazim, and Y.B. Wah. (2010). A practical approach in modelling count
17 data. In *Proceedings of the Regional Conference on Statistical Sciences* (pp. 176-183).
- 18 Zhang, Z., U. Ober, M. Erbe, H. Zhang, N. Gao, *et al.* (2014). Improving the accuracy of whole
19 genome prediction for complex traits using the results of genome wide association studies.
20 *PloS One* 9: e93017.
- 21 Zhou, M., L. Li, D. Dunson, and L. Carin. (2012). Lognormal and gamma mixed negative
22 binomial regression. In *Machine Learning: Proceedings of the International Conference on*
23 *Machine Learning* (vol. 2012. p. 1343). NIH Public Access.
- 24 Zhou, M., and L. Carin. (2012). Augment-and-conquer negative binomial processes. In
25 *Advances in Neural Information Processing Systems* (pp. 2546-2554).
- 26 Zhou, M., and L. Carin. (2015). Negative binomial process count and mixture modeling. *Pattern*
27 *Analysis and Machine Intelligence, IEEE Transactions on*, 37(2): 307-320.

1

Appendix A

2 Derivation of full conditional distribution for all parameters.

3 **Full conditional for β^***

$$\begin{aligned}
 f(\beta^* | \mathbf{y}, \text{ELSE}) &= \prod_{i=1}^I \prod_{j=1}^J \prod_{t=1}^{n_{ij}} \Pr(Y_{ijt} = y_{ijt} | \mathbf{x}_i^T, r, \omega_{ijt}, b_{1i}, b_{2ij}) f(\beta^*) \\
 &\propto \exp \left(\boldsymbol{\kappa}^T \mathbf{X} \beta^* + \boldsymbol{\kappa}^T \sum_{h=1}^2 \mathbf{Z}_h \mathbf{b}_h - \frac{1}{2} \left(\mathbf{X} \beta^* + \sum_{h=1}^2 \mathbf{Z}_h \mathbf{b}_h \right)^T \mathbf{D}_\omega \left(\mathbf{X} \beta^* + \sum_{h=1}^2 \mathbf{Z}_h \mathbf{b}_h \right) \right. \\
 &\quad \left. - \frac{1}{2} (\beta^* - \beta_0)^T \boldsymbol{\Sigma}_0^{-1} \sigma_\beta^{-2} (\beta^* - \beta_0) \right) \\
 &\propto \exp \left(-\frac{1}{2} [\beta^{*T} (\boldsymbol{\Sigma}_0^{-1} \sigma_\beta^{-2} + \mathbf{X}^T \mathbf{D}_\omega \mathbf{X}) \beta^* - 2 \left(\boldsymbol{\Sigma}_0^{-1} \sigma_\beta^{-2} \beta_0 - \mathbf{X}^T \mathbf{D}_\omega \sum_{h=1}^2 \mathbf{Z}_h \mathbf{b}_h + \mathbf{X}^T \boldsymbol{\kappa} \right)^T \beta^*] \right) \\
 &\propto \exp \left(-\frac{1}{2} [(\beta^* - \tilde{\beta}_0)^T \tilde{\boldsymbol{\Sigma}}_0^{-1} (\beta^* - \tilde{\beta}_0)] \right) \propto N(\tilde{\beta}_0, \tilde{\boldsymbol{\Sigma}}_0)
 \end{aligned}$$

4 where $\tilde{\boldsymbol{\Sigma}}_0 = (\boldsymbol{\Sigma}_0^{-1} \sigma_\beta^{-2} + \mathbf{X}^T \mathbf{D}_\omega \mathbf{X})^{-1}$, $\tilde{\beta}_0 = \tilde{\boldsymbol{\Sigma}}_0 (\boldsymbol{\Sigma}_0^{-1} \sigma_\beta^{-2} \beta_0 - \mathbf{X}^T \mathbf{D}_\omega \sum_{h=1}^2 \mathbf{Z}_h \mathbf{b}_h + \mathbf{X}^T \boldsymbol{\kappa})$.

5 **Full conditional for ω_{ijt}**

$$\begin{aligned}
 f(\omega_{ijt} | \mathbf{y}, \text{ELSE}) &\propto \exp \left[-\frac{\omega_{ijt} (\mathbf{x}_i^T \beta^* + \sum_{h=1}^2 b_{hij})^2}{2} \right] f(\omega_{ijt}; y_{ijt} + r, 0) \\
 &\propto \exp \left[-\frac{\omega_{ijt} (\mathbf{x}_i^T \beta^* + \sum_{h=1}^2 b_{hij})^2}{2} \right] f(\omega_{ijt}; y_{ijt} + r, 0) \propto \text{PG}(y_{ijt} + r, \mathbf{x}_i^T \beta^* + \sum_{h=1}^2 b_{hij})
 \end{aligned}$$

6 **Full conditional for \mathbf{b}_1**

7 Defining $\boldsymbol{\eta}^1 = \mathbf{X} \beta^* + \mathbf{Z}_2 \mathbf{b}_2$ the conditional distribution of \mathbf{b}_1 is given as

$$f(\mathbf{b}_1 | \mathbf{y}, \text{ELSE}) \propto \exp \left(\boldsymbol{\kappa}^T \mathbf{Z}_1 \mathbf{b}_1 - \frac{1}{2} (\mathbf{Z}_1 \mathbf{b}_1 + \boldsymbol{\eta}^1)^T \mathbf{D}_\omega (\mathbf{Z}_1 \mathbf{b}_1 + \boldsymbol{\eta}^1) \right) f(\mathbf{b}_1 | \sigma_{b_1}^2)$$

$$\propto \exp \left\{ -\frac{1}{2} [\mathbf{b}_1^T (\sigma_{b_1}^{-2} \mathbf{G}_1^{-1} + \mathbf{Z}_1^T \mathbf{D}_\omega \mathbf{Z}_1) \mathbf{u} - 2 (\mathbf{Z}_1^T \boldsymbol{\kappa} - \mathbf{Z}_1^T \mathbf{D}_\omega \boldsymbol{\eta}^1)^T \mathbf{b}_1] \right\}$$

1 $\propto \exp \left\{ -\frac{1}{2} (\mathbf{b}_1 - \tilde{\mathbf{b}}_1)^T \mathbf{F}_1^{-1} (\mathbf{b}_1 - \tilde{\mathbf{b}}_1) \right\} \propto f(\mathbf{b}_1 | \text{ELSE}) \sim N(\tilde{\mathbf{b}}_1, \mathbf{F}_1)$

2 where $\mathbf{F}_1 = (\sigma_{b_1}^{-2} \mathbf{G}_1^{-1} + \mathbf{Z}_1^T \mathbf{D}_\omega \mathbf{Z}_1)^{-1}$ and $\tilde{\mathbf{b}}_1 = \mathbf{F}_1 (\mathbf{Z}_1^T \boldsymbol{\kappa} - \mathbf{Z}_1^T \mathbf{D}_\omega \boldsymbol{\eta}^1)$.

3 **Full conditional for $\sigma_{b_h}^2$**

$$f(\sigma_{b_h}^2 | \mathbf{y}, \mathbf{b}) \propto \frac{1}{(\sigma_{b_h}^2)^{\frac{v_{b_h} + n_{b_h} + 1}{2}}} \exp \left(-\frac{\mathbf{b}_h^T \mathbf{G}_h^{-1} \mathbf{b}_h + v_{b_h} S_{b_h}}{2\sigma_{b_h}^2} \right)$$

4 $\propto \chi^{-2}(\tilde{v}_b = v_{b_h} + n_{b_h}, \tilde{S}_b = (\mathbf{b}_h^T \mathbf{G}_h^{-1} \mathbf{b}_h + v_{b_h} S_{b_h}) / (v_{b_h} + n_{b_h}))$

5 with $n_{b_1} = J$ and $n_{b_2} = IJ$.

6 **Full conditional for σ_{β}^2**

$$f(\sigma_{\beta}^2 | \mathbf{y}, \text{ELSE}) \propto \frac{1}{(\sigma_{\beta}^2)^{\frac{v_{\beta} + 3}{2} + 1}} \exp \left(-\frac{(\boldsymbol{\beta}^* - \boldsymbol{\beta}_0)^T \boldsymbol{\Sigma}_0^{-1} (\boldsymbol{\beta}^* - \boldsymbol{\beta}_0) + v_{\beta} S_{\beta}^*}{2\sigma_{\beta}^2} \right)$$

7 $\propto \chi^{-2}(\tilde{v}_{\beta} = v_{\beta} + 1, \tilde{S}_{\beta} = [(\boldsymbol{\beta}^* - \boldsymbol{\beta}_0)^T \boldsymbol{\Sigma}_0^{-1} (\boldsymbol{\beta}^* - \boldsymbol{\beta}_0) + v_{\beta} S_{\beta}^*] / (v_{\beta} + 1))$

8

9 **Full conditional for r**

10 To make the inference of r , we first place a gamma prior on it as $r \sim G(a_0, 1/b_0)$. Then we infer
 11 a latent count L for each count conditional on Y and r . To derive the full conditional of r , we use
 12 the following parameterization of the NB distribution: $Y \sim \text{NB}(\pi, r)$ with $\pi = \frac{\mu}{r + \mu}$. Since
 13 $L \sim \text{Pois}(-r \log(1 - \pi))$, by construction we can use the Gamma-Poisson conjugacy to update
 14 r . Therefore,

$$\begin{aligned}
 f(\mathbf{r}|\mathbf{y}, \text{ELSE}) &\propto f(r) \prod_{i=1}^I \prod_{j=1}^J \prod_{t=1}^{n_{ij}} f(y_{ijt}|L_{ijt}) f(L_{ijt}) \\
 &\propto r^{a_0-1} \exp(-rb_0) \prod_{i=1}^I \prod_{j=1}^J \prod_{t=1}^{n_{ij}} (-r \log(1 - \pi_{ij}))^{L_{ijt}} \exp(r \log(1 - \pi_{ij})) \\
 &\propto r^{a_0 + \sum_{i=1}^I \sum_{j=1}^J \sum_{t=1}^{n_{ij}} L_{ijt} - 1} \exp[-(b_0 - \sum_{i=1}^I \sum_{j=1}^J \sum_{t=1}^{n_{ij}} \log(1 - \pi_{ij})) r] \\
 &\propto G(a_0 - \sum_{i=1}^I \sum_{j=1}^J \sum_{t=1}^{n_{ij}} \log(1 - \pi_{ij}), \frac{1}{b_0 + \sum_{i=1}^I \sum_{j=1}^J \sum_{t=1}^{n_{ij}} L_{ijt}}) \quad (\text{A5})
 \end{aligned}$$

1

2

According to Zhou *et al.* (2012), the conditional posterior distribution of L_{ijt} is a Chinese

3

restaurant table (CRT) count random variable. That is, $L_{ijt} \sim CRT(y_{ijt}, r)$ and we can sample it

4

as $L_{ijt} = \sum_{l=1}^{y_{ijt}} d_l$, where $d_l \sim \text{Bernoulli}(\frac{r}{l-1+r})$. For details of the CRT random variable

5

derivation, see Zhou and Carin (2012, 2015).

6

7

8

9

10

11

12

13

14

1 **Table 1.** Posterior mean (Mean) and posterior standard deviation (SD) of the Bayesian method
 2 with four sample sizes (n_{ij}) for Model NB. S denotes scenario.

S	Parameter	True	$n_{ij} = 5$		$n_{ij} = 10$		$n_{ij} = 20$		$n_{ij} = 40$	
			Mean	SD	Mean	SD	Mean	SD	Mean	SD
1	β_0	1.5	1.484	0.357	1.488	0.269	1.542	0.233	1.549	0.213
	β_1	-1	-0.981	0.256	-0.994	0.247	-1.075	0.250	-1.016	0.190
	β_2	1	0.997	0.270	0.985	0.223	0.994	0.268	0.949	0.223
	r	5	5.079	0.916	5.078	0.519	5.017	0.471	5.027	0.330
	σ_1^2	0.5	0.542	0.196	0.594	0.176	0.582	0.180	0.590	0.216
	σ_2^2	0.5	0.503	0.134	0.524	0.136	0.531	0.110	0.512	0.114
2	β_0	1.5	1.4808	0.5009	1.4596	0.5041	1.5611	0.6108	1.4723	0.4979
	β_1	-1	-1.0631	0.2348	-0.9975	0.2040	-1.008	0.2226	-1.025	0.1908
	β_2	1	0.9504	0.2356	1.0294	0.2167	0.9925	0.1954	0.9685	0.2018
	r	5	5.1030	0.8060	4.9901	0.5928	5.0367	0.3485	5.0275	0.2033
	σ_1^2	0.5	0.5422	0.1827	0.5650	0.2199	0.5785	0.1872	0.5296	0.1837
	σ_2^2	0.5	0.4987	0.1155	0.5084	0.1423	0.5302	0.1301	0.5123	0.1047

3

4

- 1 **Table 2.** Scenarios proposed to fit the real data set with **Models NB, Pois, Normal and LN.** E
2 stands for Environment, L for lines, G for lines taking into account markers; EL and EG are
3 interaction effects of E and L and E and G.

Scenario	Main effects			Interaction effects	
	E	L	G	EL	EG
1	X	X			
2	X		X		
3	X	X		X	
4	X		X		X

4

5

1 **Table 3.** Estimated beta coefficients, variance components, and posterior predictive checks for
 2 the four scenarios (S1, S2, S3, S4) for each proposed model (**Model NB, Model Pois, Model**
 3 **Normal and Model LN**). Mean stands for posterior mean and SD for posterior standard
 4 deviation.

	Model NB							
	S1		S2		S3		S4	
Parameter	Mean	SD	Mean	SD	Mean	SD	Mean	SD
β_1^*	-0.933	0.600	-1.046	0.611	-2.521	0.711	-2.383	0.992
β_2^*	-0.826	0.710	-1.158	0.661	-2.273	0.577	-2.725	1.001
β_3^*	-0.026	0.480	-0.152	0.564	-1.688	0.851	-1.961	0.777
σ_1^2	0.425	0.048	1.374	0.167	0.341	0.050	1.033	0.153
σ_2^2					0.376	0.031	1.035	0.096
r	2.802	0.117	2.813	0.116	11.866	1.115	11.549	1.170
Loglik	-1526.649		1526.882		1268.827		-1275.253	
Cor	0.694		0.694		0.899		0.891	
MSEP	2.130		2.116		0.750		0.767	
	Model Pois							
β_1^*	-7.135	0.217	-7.211	0.388	-6.693	0.111	-6.802	0.327
β_2^*	-7.075	0.132	-7.166	0.108	-7.072	0.161	-7.266	0.185
β_3^*	-5.969	0.431	-6.463	0.293	-5.879	0.163	-6.658	0.276
σ_1^2	0.443	0.049	1.457	0.172	0.348	0.047	1.027	0.144
σ_2^2	-	-	-	-	0.381	0.031	1.045	
r	1000		1000		1000		1000	
Loglik	-1477.634		-1477.52		-1228.73		-1234.973	
Cor	0.662		0.662		0.899		0.891	
MSEP	1.866		1.860		0.743		0.758	
	Model Normal							
β_1	-12.3	5.86	7.9	4.36	13.7	3.69	9.22	3.11
β_2	-12.2	5.8	7.93	4.41	13.6	3.73	9.11	3.16
β_3	-10.4	5.87	9.66	4.36	15.5	3.69	10.94	3.1
σ_1^2	0.957	0.161	1.42	0.345	0.722	0.175	1.58	0.403
σ_2^2	-	-	-	-	1.33	0.182	1.13	0.343
r	2.75	0.136	2.91	0.147	1.67	0.109	2.23	0.172
Loglik	-1918		-1957		-1542		-1747	
Cor	0.595		0.557		0.831		0.705	
MSEP	2.405		2.600		1.073		1.679	
	Model LN							
β_1	-3.950	0.505	-6.340	3.330	1.410	0.481	3.320	1.310
β_2	-3.950	0.483	-6.330	3.320	1.410	0.487	3.320	1.290
β_3	-3.510	0.485	-5.850	3.330	1.860	0.494	3.790	1.310

σ_1^2	0.085	0.013	0.146	0.028	0.069	0.013	0.157	0.030
σ_2^2	-	-	-	-	0.081	0.011	0.053	0.018
r	0.172	0.009	0.181	0.009	0.107	0.007	0.147	0.011
Loglik	-484		-518		-125		-354	
Cor	0.709		0.679		0.882		0.789	
MSEP	2.500		2.626		1.254		1.974	

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

1 **Table 4.** Estimated posterior predictive checks with cross validation for **Models NB, Pois,**
 2 **Normal and LN.** () denotes the ranking of the four scenarios for each posterior predictive
 3 check. Each average was obtained as the mean of the rankings of the four posterior predictive
 4 checks for each scenario.

		Batan 2012		Batan 2014		Chunchi 2014	
Model NB							
Scenario		Cor	MSEP	Cor	MSEP	Cor	MSEP
S1	Mean	0.426 (3)	0.977 (3)	0.427 (4)	1.388 (2)	0.182 (3)	11.733 (4)
	SD	0.331	0.723	0.327	1.351	0.401	9.471
S2	Mean	0.423 (4)	0.980 (4)	0.432 (3)	1.383 (1)	0.204 (2)	11.222 (2)
	SD	0.327	0.717	0.325	1.356	0.373	8.614
S3	Mean	0.539 (2)	0.497 (1)	0.522 (2)	1.480 (3)	0.224 (1)	8.645 (1)
	SD	0.283	0.376	0.292	2.318	0.386	5.688
S4	Mean	0.557 (1)	0.607 (2)	0.564 (1)	1.850 (4)	0.122 (4)	11.343 (3)
	SD	0.243	0.438	0.222	2.684	0.407	8.154
Model Pois							
S1	Mean	0.426 (3)	0.977 (3)	0.427 (4)	1.388 (2)	0.182 (3)	11.733 (4)
	SD	0.331	0.723	0.327	1.351	0.401	9.471
S2	Mean	0.423 (4)	0.980 (4)	0.432 (3)	1.383 (1)	0.204 (2)	11.222 (2)
	SD	0.327	0.717	0.325	1.356	0.373	8.614
S3	Mean	0.539 (2)	0.497 (1)	0.522 (2)	1.480 (3)	0.224 (1)	8.645 (1)
	SD	0.283	0.376	0.292	2.318	0.386	5.688
S4	Mean	0.557 (1)	0.607 (2)	0.564 (1)	1.850 (4)	0.122 (4)	11.343 (3)
	SD	0.243	0.438	0.222	2.684	0.407	8.154
Model Normal							
S1	Mean	0.358 (1)	1.096 (4)	0.367 (2)	1.788 (1)	0.148 (1)	7.425 (2)
	SD	0.280	0.883	0.397	1.701	0.318	4.151
S2	Mean	0.344 (2)	0.988 (2)	0.334 (3)	2.010 (3)	0.074 (3)	7.454 (3)
	SD	0.326	0.652	0.440	2.462	0.330	4.339
S3	Mean	0.330 (3)	0.806 (1)	0.371 (1)	1.963 (2)	0.146 (2)	7.318 (1)
	SD	0.300	0.495	0.400	2.986	0.287	4.159
S4	Mean	0.267 (4)	1.029 (3)	0.237 (4)	2.373 (4)	0.039 (4)	8.482 (4)
	SD	0.338	0.731	0.445	3.420	0.238	4.326
Model LN							
S1	Mean	0.510 (1.5)	0.661 (2)	0.455 (1)	1.601 (1)	0.149 (2)	8.099 (4)
	SD	0.208	0.419	0.307	2.348	0.379	5.113
S2	Mean	0.510 (1.5)	0.663 (3)	0.433 (3)	1.778 (2)	0.086 (4)	7.819 (2)
	SD	0.224	0.392	0.353	2.820	0.459	5.311
S3	Mean	0.505 (3)	0.639 (1)	0.449 (2)	1.871 (3)	0.153 (1)	7.759 (1)
	SD	0.208	0.451	0.313	3.162	0.371	5.209
S4	Mean	0.428 (4)	0.721 (4)	0.427 (4)	1.951 (4)	0.087 (3)	8.038 (3)
	SD	0.246	0.415	0.327	3.148	0.413	5.187

1 **Table 5.** Rank averages for the four scenarios for each model (**Models NB, Pois, Normal and**
2 **LN**) resulting from the 10-fold cross-validation implemented. Each average was obtained as the
3 mean of the rankings given in Table 4 for the two posterior predictive checks (Cor and MSEP)
4 in each scenario.

Scenario	Batan 2012	Batan 2014	Chunchi 2014	Batan 2012	Batan 2014	Chunchi 2014
	Model NB			Model Normal		
S1	3	3	3.5	2.5	1.5	1.5
S2	4	2	2	2	3	3
S3	1.5	2.5	1	2	1.5	1.5
S4	1.5	2.5	3.5	3.5	4	4
	Model Pois			Model LN		
S1	3	3	3.5	1.75	1	3
S2	4	2	2	2.25	2.5	3
S3	1.5	2.5	1	2	2.5	1
S4	1.5	2.5	3.5	4	4	3

5

6

7