

Zero-inflated generalized Dirichlet multinomial regression model for microbiome compositional data analysis

Emily Palmer

Oregon State University

palmerem@oregonstate.edu

February 10, 2022

Biostatistics (2019) **20**, 4, pp. 698–713
doi:10.1093/biostatistics/kxy025
Advance Access publication on June 24, 2018

Zero-inflated generalized Dirichlet multinomial regression model for microbiome compositional data analysis

ZHENG-ZHENG TANG*

Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison, Madison, WI 53715, USA and Wisconsin Institute for Discovery, Madison, WI 53715, USA

tang@biostat.wisc.edu

GUANHUA CHEN

Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison, Madison, WI 53715, USA

- Microbial data contains an abundance of zeros
- To properly model these data, zeros must be taken into account
- Using a distribution that more appropriately approximates the true zero abundances gives statistical tests more power.
- Types of zeros:
 - Structural zeros: True zeros - taxa are absent
 - Sampling zeros: Observed zeros, undersampling

Dirichlet Multinomial

- Frequently used to model microbial counts
- Imposes negative correlation
- Only one dispersion parameter, so modeling different dispersion parameters and zero-inflation levels among multiple taxa is not considered

Generalized Dirichlet Distribution

- $K + 1$ taxa
- N total sequencing reads
- $\mathbf{Y} = (Y_1, \dots, Y_K)$, $Y_{K+1} = N - \sum_{j=1}^K Y_j$ read counts
- Unobserved underlying proportions
 $\mathbf{P} = (P_1, \dots, P_K)$, $P_{K+1} = 1 - \sum_{j=1}^K P_j$
- The Generalized Dirichlet constructs the proportions \mathbf{P} out of mutually independent Beta variables: $\mathbf{Z} = (Z_1, \dots, Z_K)$,
 $Z_j \sim \text{Beta}(a_j, b_j)$,

$$P_1 = Z_1, \quad P_j = Z_j \prod_{i=1}^{j-1} (1 - Z_i), j = 2, \dots, K$$

Generalized Dirichlet Density

- Density function for unobserved proportions \mathbf{P}

$$f(\mathbf{P}) = \prod_{j=1}^K \frac{1}{\mathcal{B}(a_j, b_j)} P_j^{a_j-1} (1 - P_1 - \dots - P_j)^{c_j}$$

where $c_j = b_j - a_{j+1} - b_{j+1}$ for $j = 1, \dots, K-1$ and $c_K = b_K - 1$

- $\mathbf{P} \sim GD(\mathbf{a}, \mathbf{b})$
- Can also obtain mutually independent Beta RVs from a GD distribution:

$$Z_1 = P_1, Z_j = \frac{P_j}{1 - \sum_{i=1}^{j-1} P_i}, j = 2, \dots, K$$

Generalized Dirichlet Multinomial (GDM)

- The GDM distribution is created by using the GD as a prior for the multinomial distribution
- \mathbf{Y} is multinomially distributed with $GD(\mathbf{a}, \mathbf{b})$ prior on proportion parameters \mathbf{P}
- Posterior probability of $(\mathbf{P}|\mathbf{Y})$ is $GD(\mathbf{a}^*, \mathbf{b}^*)$ with
 $\mathbf{a}^* = (a_1^*, \dots, a_K^*), \mathbf{b}^* = (b_1^*, \dots, b_K^*),$
 $a_j^* = a_j + Y_j, b_j^* = b_j + Y_{J+1} + \dots + Y_{K+1}, j = 1, \dots, K$
- GD distribution assumes all taxa have positive proportions (every taxa present in the sample)
- Observed zeros are sampling zeros using GD distribution

Zero-inflated Generalized Dirichlet distribution

- Introduce zero-inflation to model structural zeros (absent taxa)
- Assume Z_j follows a zero-inflated Beta (ZIB) distribution
 $Z_j \sim ZIB(\pi_j, a_j, b_j)$
 - π_j is the probability $Z_j = 0$
- Create the ZIGD distribution by replacing Beta distributed Z_j with ZIB Z_j in GD

$$P_1 = Z_1, \quad P_j = Z_j \prod_{i=1}^{j-1} (1 - Z_i), j = 2, \dots, K$$

- $\mathbf{P} \sim ZIGD(\boldsymbol{\pi}, \mathbf{a}, \mathbf{b})$

- $Z_j = 0$ implies $P_j = 0$
- $\Delta_j = I(Z_j = 0) = I(P_j = 0) \sim \text{Bern}(\pi_j)$
- If we assume $\Delta_1 = \dots = \Delta_K = 0$, the ZIGD is the GD
- Assume there are L taxa present in a sample, $U = (u_1, \dots, u_L)$ indices for present taxa. \bar{U} indices of absent taxa
- M observed zero count taxa with indices $V = (v_1, \dots, v_M)$
- $U \cap V$ taxa present in the sample but have zero counts due to undersampling

ZIGD conjugate prior for multinomial

- \mathbf{Y} follows a multinomial with $\text{ZIGD}(\boldsymbol{\pi}, \mathbf{a}, \mathbf{b})$ prior on proportion parameters \mathbf{P}
- $\boldsymbol{\Delta} = (\Delta_1, \dots, \Delta_K)$
- $\Delta_j = I(P_j = 0)$, so $f(\mathbf{P}, \boldsymbol{\Delta}) = f(\mathbf{P})$
- Posterior probability of proportions given observed counts:

$$f(\mathbf{P}|\mathbf{Y}) = f(\mathbf{P}, \boldsymbol{\Delta}|\mathbf{Y}) = f(\mathbf{P}|\boldsymbol{\Delta}, \mathbf{Y})f(\boldsymbol{\Delta}|\mathbf{Y})$$

ZIGD conjugate prior for multinomial

- Since $P_j = 0$ when $\Delta_j = 1$ (taxon has proportion of 0 if it is absent from a sample)

$$f(\mathbf{P}|\mathbf{\Delta}) = I(\mathbf{P}_{\bar{U}} = 0)f(\mathbf{P}_U|\mathbf{\Delta}_U = 0, \mathbf{\Delta}_{\bar{U}} = 1)$$

- Since $f(\mathbf{P}_U|\mathbf{\Delta}_U = 0, \mathbf{\Delta}_{\bar{U}} = 1) \sim GD(\mathbf{a}_U, \mathbf{b}_U)$, and the GD is conjugate prior for the multinomial, the posterior probability is

$$f(\mathbf{P}_U|\mathbf{\Delta}_U = 0, \mathbf{\Delta}_{\bar{U}} = 1, \mathbf{Y}) \sim GD(\mathbf{a}_U^*, \mathbf{b}_U^*)$$

- $\Delta_j = 0$ when $Y_j > 0$ (a taxon is present in the sample if the count is positive)

$$f(\mathbf{\Delta}, |\mathbf{Y}) = I(\mathbf{\Delta}_{\bar{V}} = 0)f(\mathbf{\Delta}_V|\mathbf{Y}_V = 0, Y_{\bar{V}} > 0)$$

- The mass function for Δ_V given $Y_V = 0, Y_{\bar{V}} > 0$ is then

$$\begin{aligned}
 & f(\Delta_V \mid Y_V = 0, Y_{\bar{V}} > 0) \\
 & \propto f(\Delta_V) f(Y_V = 0, Y_{\bar{V}} > 0 \mid \Delta_V, \Delta_{\bar{V}} = 0) \\
 & = f(\Delta_V) \int_{\mathbf{P}} f(Y_V = 0, Y_{\bar{V}} > 0 \mid \mathbf{P}, \Delta_V, \Delta_{\bar{V}} = 0) f(\mathbf{P} \mid \Delta_V, \Delta_{\bar{V}} = 0) d\mathbf{P} \\
 & \propto \prod_{i \in \mathcal{V}} \left\{ \pi_j^{\Delta_j} (1 - \pi_j)^{(1 - \Delta_j)} \right\} \prod_{i \in I \cap \mathcal{V}} \left\{ \frac{\mathcal{B}(a_j^*, b_j^*)}{\mathcal{B}(a_j, b_j)} \right\} \\
 & = \prod_{j \in \mathcal{V}} \left\{ \pi_j^{\Delta_j} \left[(1 - \pi_j) \frac{\mathcal{B}(a_j^*, b_j^*)}{\mathcal{B}(a_j, b_j)} \right]^{(1 - \Delta_j)} \right\}.
 \end{aligned}$$

- Thus the posterior probability $f(\mathbf{P} \mid \mathbf{Y})$ follows a ZIGD with the zero-inflation on the taxa having observed zero counts
- Probability of an observed zero being structural is

$$\frac{\pi_j}{\pi_j + (1 - \pi_j) \frac{\mathcal{B}(a_j^*, b_j^*)}{\mathcal{B}(a_j, b_j)}} \quad (j \in \mathcal{V})$$

ZIGDM Regression Model

- n subjects measured on $K + 1$ taxa
- Y_{ij}, P_{ij} observed count and true proportion for taxon j in subject i .
- \mathbf{X}_i d -dimensional vector of intercept, covariates, and confounding variables
- Assume \mathbf{Y}_i follows a ZIGDM($\boldsymbol{\pi}_i, \mathbf{a}_i, \mathbf{b}_i$)
- The model is then

$$\Delta_{ij} \sim \text{Bernoulli}(\pi_{ij}), \quad j = 1, \dots, K,$$

$$Z_{ij} = 0 \text{ if } \Delta_{ij} = 1, \quad Z_{ij} \mid \Delta_{ij} = 0 \sim \text{Beta}(a_{ij}, b_{ij}), \quad j = 1, \dots, K,$$

$$P_{i1} = Z_{i1}, \quad P_{ij} = Z_{ij} \prod_{k=1}^{j-1} (1 - Z_{ik}), \quad j = 2, \dots, K,$$

$$\mathbf{Y}_i \mid \mathbf{P}_i \sim \text{Multinomial}(\mathbf{P}_i, N_i), \text{ where } \mathbf{P}_i = (P_{i1}, \dots, P_{iK}) \text{ and } N_i = \sum_{j=1}^{K+1} Y_{ij}.$$

Link functions

- With this model, we can link π_i , \mathbf{a}_i , and \mathbf{b}_i to \mathbf{X}_i
- π_{ij} is the probability of absence
- a_{ij} , b_{ij} control abundance distribution at the presence
- Reparametrize $\mu_{ij} = \frac{a_{ij}}{a_{ij} + b_{ij}}$ and $\sigma_{ij} = \frac{1}{1 + a_{ij} + b_{ij}}$ as the mean of the Beta variable and the dispersion parameter (since the variance of the Beta variables are of the form $\mu_{ij}(1 - \mu_{ij})\sigma_{ij}$)
- Use logit link functions (μ_{ij} , μ_{ij} , σ_{ij} all between 0 and 1)

$$\pi_{ij} = \frac{e^{\boldsymbol{\gamma}_j^T \mathbf{x}_i}}{1 + e^{\boldsymbol{\gamma}_j^T \mathbf{x}_i}}, \quad \mu_{ij} = \frac{e^{\boldsymbol{\alpha}_j^T \mathbf{x}_i}}{1 + e^{\boldsymbol{\alpha}_j^T \mathbf{x}_i}}, \quad \text{and} \quad \sigma_{ij} = \frac{e^{\boldsymbol{\beta}_j^T \mathbf{x}_i}}{1 + e^{\boldsymbol{\beta}_j^T \mathbf{x}_i}}, \quad j = 1, \dots, K,$$

- $\boldsymbol{\gamma}_j = (\gamma_{1j}, \dots, \gamma_{dj})$, $\boldsymbol{\alpha}_j = (\alpha_{1j}, \dots, \alpha_{dj})$, $\boldsymbol{\beta}_j = (\beta_{1j}, \dots, \beta_{dj})$ regression coefficients for taxon j .

Estimating parameters

- Denote $\boldsymbol{\theta} = (\gamma_1, \dots, \gamma_K, \alpha_1, \dots, \alpha_K, \beta_1, \dots, \beta_K)$ the complete set of parameters
- Likelihood-based inference on $\boldsymbol{\theta}$ is difficult since the observed log-likelihood function is complicated
- Complete data log-likelihood in terms of Z (instead of \mathbf{P})

$$\begin{aligned} l(\boldsymbol{\theta}) &= \log \left[\prod_{i=1}^n \left\{ f(\mathbf{Y}_i \mid \mathbf{Z}_i) \prod_{j=1}^K f(Z_{ij}) \right\} \right] \\ &= \sum_{i=1}^n \log \{f(\mathbf{Y}_i \mid \mathbf{Z}_i)\} \\ &\quad + \sum_{j=1}^K \sum_{i=1}^n \left\{ \Delta_{ij} \log \pi_{ij} + (1 - \Delta_{ij}) \log(1 - \pi_{ij}) + \right. \\ &\quad \left. (1 - \Delta_{ij}) \left[-\log(\mathcal{B}(a_{ij}, b_{ij})) + (a_{ij} - 1) \log(Z_{ij}) + (b_{ij} - 1) \log(1 - Z_{ij}) \right] \right\} \end{aligned}$$

EM algorithm for estimating parameters

- In the t -th E-step, compute the expected complete data log-likelihood:

$$\mathcal{Q}_{\theta}^* = \sum_{j=1}^K \sum_{i=1}^n \mathbb{E} \left\{ \Delta_{ij} \log \pi_{ij} + (1 - \Delta_{ij}) \log(1 - \pi_{ij}) + \right. \\ \left. (1 - \Delta_{ij}) \left[-\log(\mathcal{B}(a_{ij}, b_{ij})) + (a_{ij} - 1) \log Z_{ij} + (b_{ij} - 1) \log(1 - Z_{ij}) \right] \right\}$$

- Expectation is with respect to the posterior distributions of $(\Delta_i | \mathbf{Y}_i; \boldsymbol{\theta}^{(t-1)})$ and $(\mathbf{Z}_i | \Delta_i, \mathbf{Y}_i; \boldsymbol{\theta}^{t-1})$
- $\boldsymbol{\theta}^{(t-1)}$ the parameter estimates in the $(t - 1)$ -th M-step

$$\Delta_{ij}^* = E(\Delta_{ij} \mid \mathbf{Y}_i) = \begin{cases} 0 & \text{if } Y_{ij} > 0 \\ \frac{\pi_{ij}}{\pi_{ij} + (1 - \pi_{ij}) \frac{\mathcal{B}(a_{ij}^*, b_{ij}^*)}{\mathcal{B}(a_{ij}, b_{ij})}} & \text{if } Y_{ij} = 0 \end{cases},$$

$$A_{ij}^* = E(\log Z_{ij} \mid \mathbf{Y}_i, \Delta_{ij} = 0) = \psi(a_{ij}^*) - \psi(a_{ij}^* + b_{ij}^*),$$

$$B_{ij}^* = E(\log(1 - Z_{ij}) \mid \mathbf{Y}_i, \Delta_{ij} = 0) = \psi(b_{ij}^*) - \psi(a_{ij}^* + b_{ij}^*)$$

- $a_{ij}^* = a_{ij} + Y_{ij}$, $b_{ij}^* = b_{ij} + Y_{i(j+1)} + \cdots + Y_{i(K+1)}$, $\psi()$ is the digamma function
- Rewrite Q_θ^* as:

$$Q_\theta^* = \sum_{j=1}^K Q_{\gamma_j}^* + \sum_{j=1}^K Q_{\alpha_j, \beta_j}^*$$

- $Q_{\gamma_j}^* = \sum_{i=1}^n \{\Delta_{ij}^* \log \pi_{ij} + (1 - \Delta_{ij}^*) \log(1 - \pi_{ij})\}$ and
 $Q_{\alpha_j, \beta_j}^* = \sum_{i=1}^n (1 - \Delta_{ij}^*) \{-\log(\mathcal{B}(a_{ij}, b_{ij})) + (a_{ij} - 1)A_{ij}^* + (b_{ij} - 1)B_{ij}^*\}$

- In the t -th M-step, for each taxon j , obtain $\gamma_j^{(t)}$ from maximizing $Q_{\gamma_j}^*$, and obtain $\alpha_j^{(t)}, \beta_j^{(t)}$ by maximizing Q_{α_j, β_j}^*
- Computation for optimization is the same as a logistic and weighted Beta regression
- Parameters for individual taxa updated independently, estimation can be done in parallel for different parameters
- Results in a computationally efficient EM algorithm relying on simple posterior estimation calculations and parallel parameter updates for taxa

Association tests

- Test null hypothesis that covariates are not associated with the mean
 $H_0 : \alpha_{*1} = \dots = \alpha_{*K} = 0$
- Or covariates are not associated with dispersion
 $H_0 : \beta_{*1} = \dots = \beta_{*K} = 0$
- α_{*j}, β_{*j} subsets of α_j, β_j corresponding to covariates of interest
- Can test using score, Wald, or likelihood-ratio statistics
- This paper uses score statistics, which are computationally faster and more stable
- Use permutation techniques for p values since asymptotic approximations may not be accurate when most observations are zero
- Permute covariate of interest and calculate score test statistic in each permutation

Thank you!