

La loi négative binomiale pour le traitement des données des microbiotes

SAHBANE Abdesstar

Montpellier University

November 8, 2020



Table of contents

- 1 Présentation des données
- 2 Loi négative binomiale
- 3 Fonction de lien
- 4 IWLS

Presentation des données

Les données typiques sur le microbiome générées par le séquençage du gène, comprennent les composants suivants :

- C_{ij} : Counts.
- X_i : Host factors.
- T_i : Total sequence read.
- Z_i : Sample variables.

Loi négative binomiale

L'objectif est de détecter les associations entre les caractéristiques du microbiome C_{ij} et la variable X_i .

$\forall j$ on note $y_i = C_{ij}$ et on suppose que y_i suit la loi négative binomial :

$$y_i \sim NB(y_i \mid \mu_i, \theta) = \frac{\Gamma(y_i + \theta)}{\Gamma(\theta)y_i!} \cdot \left(\frac{\theta}{\mu_i + \theta}\right)^\theta \cdot \left(\frac{\mu_i}{\mu_i + \theta}\right)^{y_i} \quad (1)$$

où μ_i et θ sont les paramètres de forme, et Γ et la fonction gamma.

Fonction de lien

Le modèle négative binomiale mixte relie la moyenne μ_i aux variables T_i , X_i et Z_i par le lien :

$$\log(\mu_i) = \log(T_i) + X_i\beta + Z_ib \quad (2)$$

avec $b \sim N_K(0, \tau^2 I)$

Famille exponentielle

Pour tout paramètre θ fixe, la loi négative binomiale appartient à la famille exponentielle, et s'écrit sous la forme :

$$NB(y_i | \mu_i, \theta) = \exp \left\{ \frac{y_i \vartheta_i - b(\vartheta_i)}{\phi} + c(y_i, \phi) \right\} \quad (3)$$

où $\vartheta_i = \log \frac{\mu_i}{\mu_i + \theta}$,

$b(\vartheta_i) = -\theta \log \left(1 - e^{\log \frac{\mu_i}{\mu_i + \theta}} \right) = -\theta \log (1 - e^{\vartheta_i})$, et

$c(y_i, \phi) = \log \left(\frac{\Gamma(y_i + \theta) \theta^\theta}{\Gamma(\theta) y_i!} \right)$

IWLS

- Le modèle négative binomiale est donc un cas spéciale des modèles linéaires généralisés pour tout θ fixe.
- lorsque θ est inconnu, le modèle négative binomiale n'est pas un GLM. Or, les NBMMs peuvent être ajustés en mettant à jour de manière itérative les paramètres (β, b, τ^2) et θ .
- Conditionnellement à θ , le NBMM est un GLMM spécial et donc les paramètres (β, b, τ^2) peuvent être mis à jour en utilisant la procédure GLMMs.
- Conditionnellement à (β, b) , le paramètre de forme θ peut être mis à jour en maximisant la vraisemblance NB en utilisant l'algorithme standard de Newton – Raphson.

IWLS

- Conditionnellement à θ , nous mettons à jour les paramètres (β, b, τ^2) en étendant l'algorithme IWLS.
- L'algorithme IWLS procède à une approximation de la vraisemblance du modèle linéaire généralisé par une vraisemblance normale pondérée, puis met à jour les paramètres à partir du modèle normal pondéré.
- Conditionnellement aux paramètres θ , β , et b , la vraisemblance binomiale négative $NB(y_i | \mu_i, \theta)$ peut être approximée par la vraisemblance normale pondérée :

$$NB(y_i | \mu_i, \theta) \approx N(t_i | \eta_i, w_i^{-1})$$

où $\eta_i = \log(T_i) + X_i\beta + Z_ib$, les «données de réponse normales» t_i et les «poids» w_i sont appelés respectivement pseudo-réponse et pseudo-poids.

IWLS

La pseudo-réponse t_i et les pseudo-poids w_i sont calculés par:

$$t_i = \hat{\eta}_i - \frac{L' (y_i | \hat{\eta}_i, \hat{\theta})}{L'' (y_i | \hat{\eta}_i, \hat{\theta})}, \text{ et } w_i = -L'' (y_i | \hat{\eta}_i, \hat{\theta})$$

Où $\hat{\eta}_i = \log (T_i) + X_i \hat{\beta} + Z_i \hat{b}$, $L (y_i | \hat{\eta}_i, \hat{\theta}) = \log NB (y_i | \hat{\mu}_i, \hat{\theta})$,

$L' (y_i | \eta_i, \theta) = dL (y_i | \eta_i, \theta) / d\eta_i$, $L'' (y_i | \eta_i, \theta) =$

$d^2 L (y_i | \eta_i, \theta) / d\eta_i^2$,

et $(\hat{\beta}, \hat{b})$ et $\hat{\theta}$ sont les estimations actuelles de (β, b) et θ , respectivement.

IWLS

Par conséquent, les NBMM peuvent être approximés par le modèle mixte linéaire avec w_i comme poids:

$$t_i = \log(T_i) + X_i\beta + Z_ib + w_i^{-1/2}e_i, b \sim N_K(0, \tau^2), e \sim N_n(0, \sigma^2 I)$$

Les paramètres $(\beta, b, \tau^2, \sigma^2)$ sont ensuite mis à jour à partir de ce modèle mixte linéaire en utilisant l'algorithme standard d'ajustement des LMM.

Résumé IWLS

En résumé, l'IWLS pour l'ajustement des NBMM est un algorithme itératif et se déroule comme suit:

- ❶ Initialiser β , b et θ quelques valeurs plausibles;
- ❷ $j = 1, 2, \dots$
 - ❶ Sur la base des valeurs actuelles $(\beta^{(j-1)}, b^{(j-1)}, \theta^{(j-1)})$, calculer la pseudo-réponse $t_i^{(j)}$ et les pseudo-poids $w^{(j)}$;
 - ❷ Mettre à jour $(\beta, b, \tau^2, \sigma^2)$ en ajustant le LMM;
 - ❸ Mettre à jour θ par l'algorithme standard de Newton – Raphson.
- ❸ Répétez l'étape 2) jusqu'à convergence.

Nous utilisons le critère $(\eta^{(j)} - \eta^{(j-1)})^2 < \varepsilon (\eta^{(j)})^2$, pour évaluer la convergence, où $\eta^{(j)} = \sum_{i=1}^n (\log(T_i) + X_i\beta^{(j)} + Z_ib^{(j)})$, et ε est une petite valeur (disant 10^{-5}).

Frame Title

Merci