

Batch Effects Correction for Microbiome Data with Dirichlet-multinomial Regression

DAI ZHENWEI

23 Oct, 2018

1. Introduction

Metagenomic sequencing techniques enable quantitative analyses of the microbiome. However, combining the microbial data from these experiments is challenging due to the variations between experiments. The existing methods for correcting batch effects do not consider the interactions between variables—microbial taxa in microbial studies—and the overdispersion of the microbiome data. Therefore, they are not applicable to microbiome data.

We develop a new method, Bayesian Dirichlet-multinomial regression meta-analysis (BDMMA), to simultaneously model the batch effects and detect the microbial taxa associated with phenotypes. BDMMA automatically models the dependence among microbial taxa and is robust to the high dimensionality of the microbiome and their association sparsity.

The package `BDMMAcorrect` includes functions to perform meta-analysis of the metagenomic compositional data and select taxa significantly associated with the covariates. BDMMA is based on the following assumptions:

1. Only a small proportion of taxa are significantly associated with the covariates.
2. The taxonomic read counts follows a Dirichlet-Multinomial (DM) distribution.
3. The batch effects are independent of the covariates' effects.
4. The batch information for each sample is known.

Brief introduction to the BDMMA model

Suppose the taxonomic read counts of a sample y_{ij} (the j -th sample in i -th batch) follows a DM distribution parameterized by $\gamma_{ij} = (\gamma_{ij1}, \gamma_{ij2}, \dots, \gamma_{ijG})$,

$$f_{DM}(y_{ij}|\gamma_{ij}) = \frac{\Gamma(\gamma_{ij+})\Gamma(y_{ij+} + 1)}{\Gamma(y_{ij+} + \gamma_{ij+})} \times \prod_{g=1}^G \frac{\Gamma(y_{ijg} + \gamma_{ijg})}{\Gamma(\gamma_{ijg})\Gamma(y_{ijg} + 1)},$$

where $y_{ij+} = \sum_{g=1}^G \gamma_{ijg}$ and G encodes the number of all the taxa included in the analysis. We model the parameter γ_{ijg} with,

$$\gamma_{ijg} = \alpha_g \times \exp\left(\sum_{p=1}^P X_{ijp}\beta_{pg} + \delta_{ig}\right),$$

where g means the g -th taxon; $X = (X_{ijp})_{N \times P}$ is the covariate matrix (N is the total number of samples and P is the number of covariate variables); $\delta = (\delta_{ig})_{I \times G}$ is the batch effects matrix satisfying $\sum_{i=1}^I n_i \delta_{ig} = 0$ (n_i is the sample size of the i -th batch). α_g and β_{pg} encode the intercept and covariate coefficient respectively.

We adopt the Bayesian approach and provide proper prior distributions for the parameters. To select the taxa significantly associated with the variable of interest, we impose a spike-and-slab prior on the corresponding

coefficients and estimate their posterior inclusion probability (PIP). The variable selection is conducted by thresholding PIP. In the next section, we provide an example to show the usage of functions in our package.

2. Analysis Example

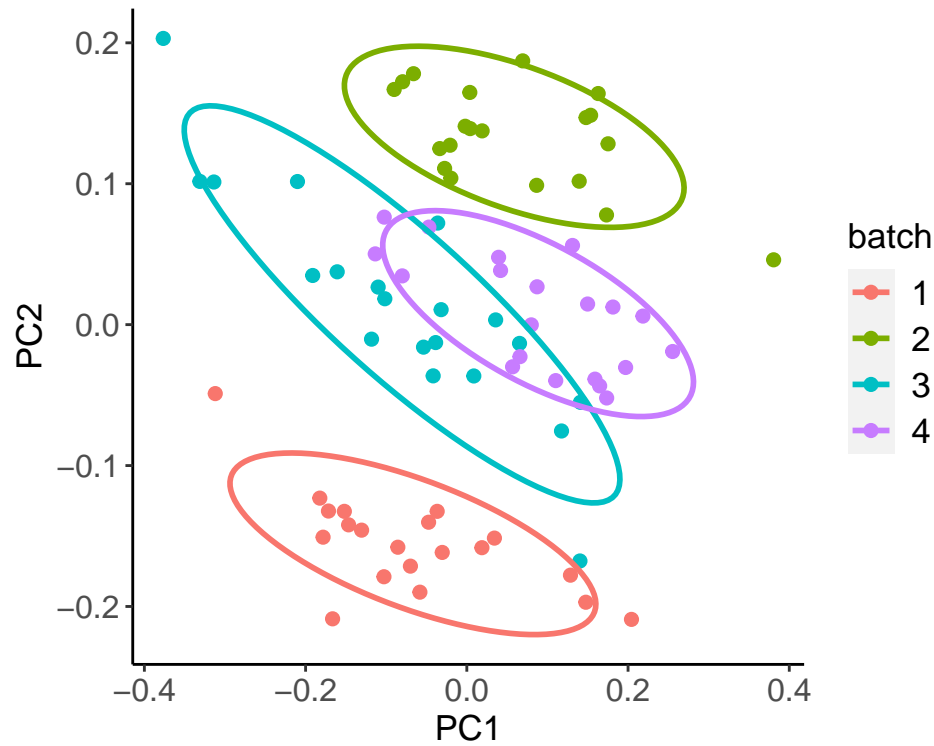
We simulated a sample data set, named `Microbiome_dat`, including 80 samples and 40 taxa. The read counts were simulated from the DM distribution according to the BDMMA model assumptions. We also simulated a main effect variable (case/control status) and a confounding variable. `Microbiome_dat` is a `SummarizedExperiment` object and be loaded directly.

We use `colData` to retrieve the phenotype and batch information, and store them in `col_data`. Variable `main` is the main effect variable, `confounder` is the confounding variable and `batch` includes the batch labels of all the samples. The read counts can be accessed with `assay`, which will return a `matrix` object.

```
library(BDMMAcorrect)
require(SummarizedExperiment)
data(Microbiome_dat)
### Access phenotypes information
col_data <- colData(Microbiome_dat)
pheno <- data.frame(col_data$main, col_data$confounder)
batch <- col_data[,3]
### Access taxonomy read counts
counts <- t(assay(Microbiome_dat))
### Indicate whether the phenotype variables are continuous
continuous <- mcols(col_data)[1:2,]
```

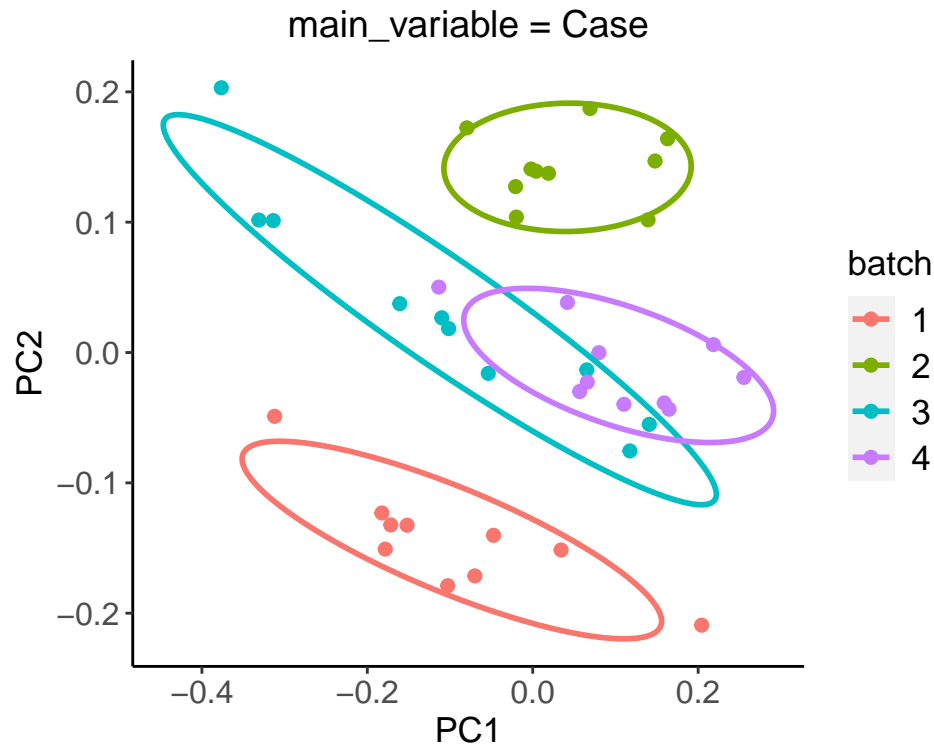
`BDMMAcorrect` provides a function to visualize the differences of the taxonomic composition across cohorts with the principal coordinate analysis. `VBatch` plots the first two principal coordinates of the corresponding samples and the 80% confidence ellipse of each batch.

```
figure = VBatch(Microbiome_dat = Microbiome_dat, method = "bray")
print(figure)
```

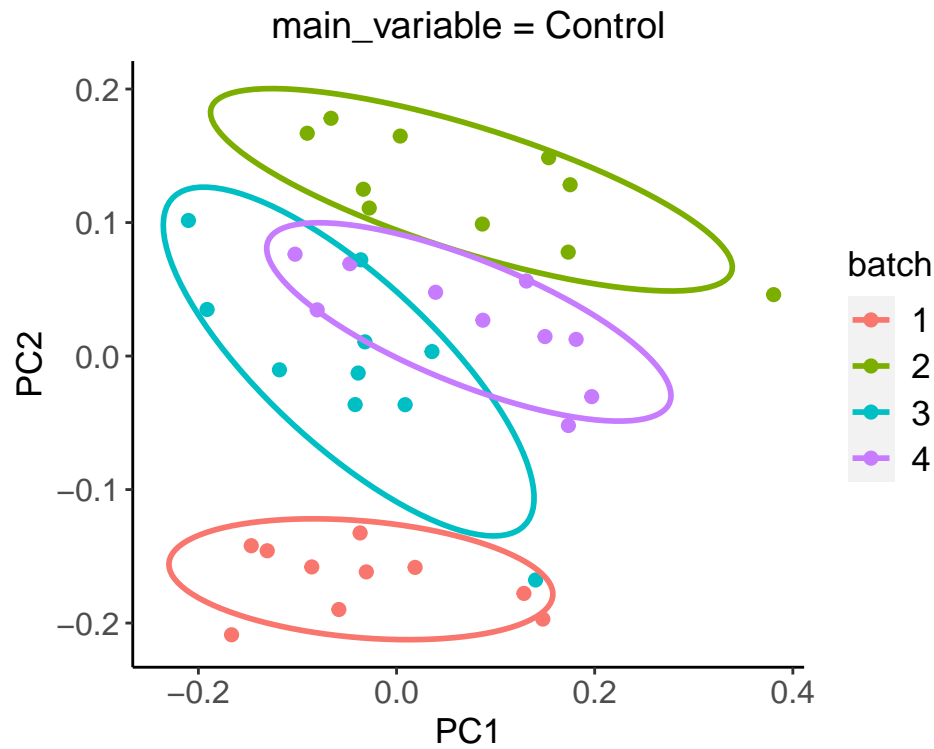


For a case/control study, `VBatch` can also visualize the batch effect of the case and control samples respectively.

```
main_variable <- pheno[,1]
main_variable[main_variable == 0] <- "Control"
main_variable[main_variable == 1] <- "Case"
figure <- VBatch(Microbiome_dat = Microbiome_dat, main_variable = main_variable, method = "bray")
print(figure[[1]])
```



```
print(figure[[2]])
```



Then, the main function `BDMMA` can work directly on `Microbiome_dat`. `BDMMAcorrect` provides the freedom to ignore the low abundant taxa. `BDMMAcorrect` runs a Markov chain Monte Carlo algorithm to sample from the posterior distribution. The users can set the lengths of the burn-in period and the sampling period for

the Markov chain. In this example, the lengths of burn-in and sampling period are set to 4000.

```
output <- BDMMA(Microbiome_dat = Microbiome_dat, burn_in = 4000, sample_period = 4000)
```

```
## ##### Start MCMC #####
##
## Iteration = 50
## Iteration = 100
## Iteration = 150
## Iteration = 200
## Iteration = 250
## Iteration = 300
## Iteration = 350
## Iteration = 400
## Iteration = 450
## Iteration = 500
## Iteration = 550
## Iteration = 600
## Iteration = 650
## Iteration = 700
## Iteration = 750
## Iteration = 800
## Iteration = 850
## Iteration = 900
## Iteration = 950
## Iteration = 1000
## Iteration = 1050
## Iteration = 1100
## Iteration = 1150
## Iteration = 1200
## Iteration = 1250
## Iteration = 1300
## Iteration = 1350
## Iteration = 1400
## Iteration = 1450
## Iteration = 1500
## Iteration = 1550
## Iteration = 1600
## Iteration = 1650
## Iteration = 1700
## Iteration = 1750
## Iteration = 1800
## Iteration = 1850
## Iteration = 1900
## Iteration = 1950
## Iteration = 2000
## Iteration = 2050
## Iteration = 2100
## Iteration = 2150
## Iteration = 2200
## Iteration = 2250
## Iteration = 2300
## Iteration = 2350
## Iteration = 2400
## Iteration = 2450
```

```
## Iteration = 2500
## Iteration = 2550
## Iteration = 2600
## Iteration = 2650
## Iteration = 2700
## Iteration = 2750
## Iteration = 2800
## Iteration = 2850
## Iteration = 2900
## Iteration = 2950
## Iteration = 3000
## Iteration = 3050
## Iteration = 3100
## Iteration = 3150
## Iteration = 3200
## Iteration = 3250
## Iteration = 3300
## Iteration = 3350
## Iteration = 3400
## Iteration = 3450
## Iteration = 3500
## Iteration = 3550
## Iteration = 3600
## Iteration = 3650
## Iteration = 3700
## Iteration = 3750
## Iteration = 3800
## Iteration = 3850
## Iteration = 3900
## Iteration = 3950
## Iteration = 4000
## Iteration = 4050
## Iteration = 4100
## Iteration = 4150
## Iteration = 4200
## Iteration = 4250
## Iteration = 4300
## Iteration = 4350
## Iteration = 4400
## Iteration = 4450
## Iteration = 4500
## Iteration = 4550
## Iteration = 4600
## Iteration = 4650
## Iteration = 4700
## Iteration = 4750
## Iteration = 4800
## Iteration = 4850
## Iteration = 4900
## Iteration = 4950
## Iteration = 5000
## Iteration = 5050
## Iteration = 5100
## Iteration = 5150
```

```
## Iteration = 5200
## Iteration = 5250
## Iteration = 5300
## Iteration = 5350
## Iteration = 5400
## Iteration = 5450
## Iteration = 5500
## Iteration = 5550
## Iteration = 5600
## Iteration = 5650
## Iteration = 5700
## Iteration = 5750
## Iteration = 5800
## Iteration = 5850
## Iteration = 5900
## Iteration = 5950
## Iteration = 6000
## Iteration = 6050
## Iteration = 6100
## Iteration = 6150
## Iteration = 6200
## Iteration = 6250
## Iteration = 6300
## Iteration = 6350
## Iteration = 6400
## Iteration = 6450
## Iteration = 6500
## Iteration = 6550
## Iteration = 6600
## Iteration = 6650
## Iteration = 6700
## Iteration = 6750
## Iteration = 6800
## Iteration = 6850
## Iteration = 6900
## Iteration = 6950
## Iteration = 7000
## Iteration = 7050
## Iteration = 7100
## Iteration = 7150
## Iteration = 7200
## Iteration = 7250
## Iteration = 7300
## Iteration = 7350
## Iteration = 7400
## Iteration = 7450
## Iteration = 7500
## Iteration = 7550
## Iteration = 7600
## Iteration = 7650
## Iteration = 7700
## Iteration = 7750
## Iteration = 7800
## Iteration = 7850
```

```

## Iteration = 7900
## Iteration = 7950
## Iteration = 8000

print(output$selected.taxa)

## list()

head(output$parameter_summary)

##           mean      X2.5.      X25.      X50.      X75.      X97.5.
## alpha_1 9.202631 8.346548 8.851741 9.170373 9.507178 10.327741
## alpha_2 3.151589 2.647339 3.006299 3.150018 3.301027 3.620281
## alpha_3 4.713369 4.087989 4.530339 4.705070 4.901011 5.341901
## alpha_4 1.403778 1.101862 1.318810 1.410495 1.501926 1.665869
## alpha_5 2.741184 2.242398 2.613441 2.752701 2.897369 3.155614
## alpha_6 4.288900 3.701325 4.107043 4.300437 4.480630 4.851532

print(output$PIP)

## [1] 0.9655086 0.9780055

print(output$bFDR)

## [1] 0.02824294

```

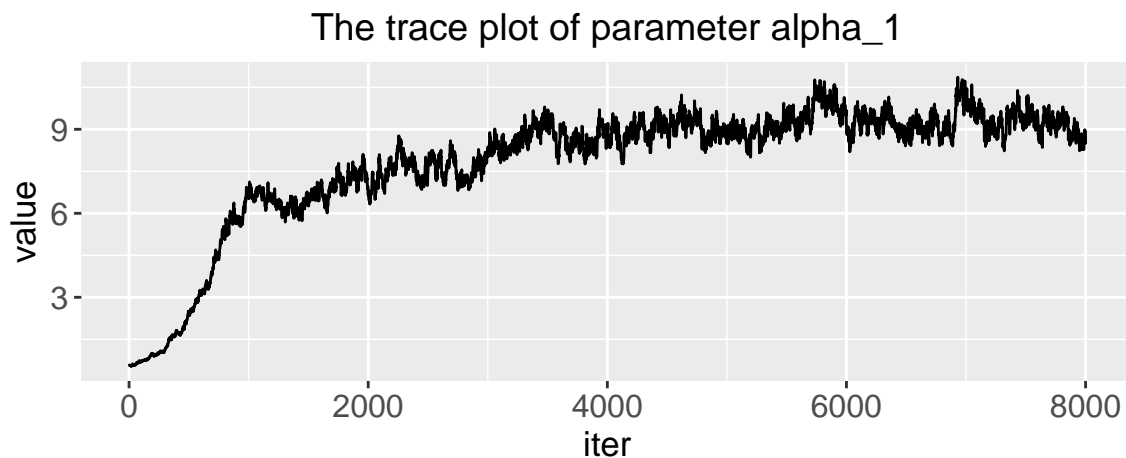
The output includes three arguments, `selected.taxa`, `parameter_summary` and `trace`. The selected taxa by thresholding PIPs and controlling Bayesian false discovery rate are listed in `output$selected.taxa`. The default bFDR level and the PIP thresholds are set to 0.1 and 0.5 respectively in the function. `BDMMAcorrect` selects V10 and V30 as significantly associated taxa. Users can check the mean and quantiles of parameters' posterior distribution in `output$parameter_summary`. `output$PIP` shows the PIPs of the selected taxa. Given the selected microbial taxa, `output$bFDR` provides the corresponding bFDR. `output$trace` includes the trace of the parameters and the function `trace_plot` can be used to check the convergence of the Markov chain.

```

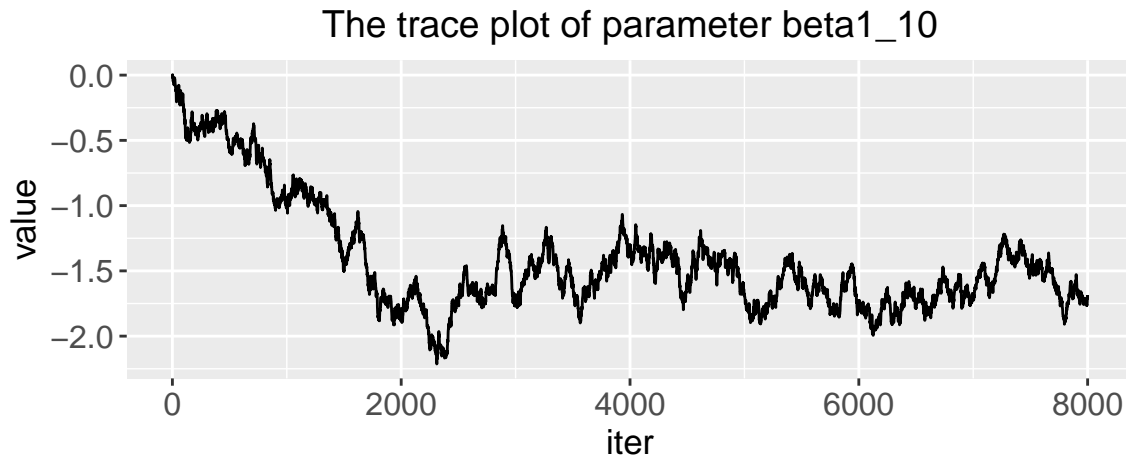
figure <- trace_plot(trace = output$trace, param = c("alpha_1", "beta1_10"))
print(figure)

```

```
## [[1]]
```



```
##
## [[2]]
```

Prepare data for analysis

To prepare develop the data for the analysis, here, the following example shows how to pack the data into a `SummarizedExperiment` object that can be used as the input of the functions.

```
### Simulate counts
counts <- rmultinom(100,10000,rep(0.02,50))
### Simulate covariates
main <- rbinom(100,1,0.5)
confounder <- rnorm(100,0,1)
### Simulate batches
batch <- c(rep(1,50),rep(2,50))

library(SummarizedExperiment)
col_data <- DataFrame(main, confounder, batch)
mcols(col_data)$continuous <- c(OL, 1L, OL)
### pack different datasets into a SummarizedExperiment object
Microbiome_dat <- SummarizedExperiment(list(counts), colData=col_data)
```