

Test with genus abundance data

Weijia Xiong

5/14/2020

Real data

From <https://github.com/chvlyl/PLEASE>

Inflammation, Antibiotics, and Diet as Environmental Stressors of the Gut Microbiome in Pediatric Crohn's Disease. Cell Host & Microbe. 2015

Load the raw data

The “unclassified” taxa were removed and the total relative abundance in each sample were normalized to be one. “P”, “F”, “G”, “S” at the beginning of each file indicate taxonomic levels “phylum”, “family”, “genus”, “species”.

We use genus abundance data.

```
PLEASE.raw =  
  read_excel("./real data/G_Remove_unclassified_Renormalized_Merge_Rel_MetaPhlAn_Result.xlsx",  
             col_names = T) %>%  
  as.data.frame()  
row.names(PLEASE.raw) = PLEASE.raw[,1]  
PLEASE.raw = PLEASE.raw[,-1]  
colname = as.numeric(colnames(PLEASE.raw))  
samplepoint = as.Date(colname,origin = "1900-01-01")  
head(PLEASE.raw[,1:5])
```

```
##           1132619    1132650    1132678    1132709    1132984  
## g__Bacteroides    0.4272231  0.2318847  0.020411137  0.005330001  72.3829938  
## g__Ruminococcus    0.0000000  0.0000000  0.006130341  0.000000000  4.9631993  
## g__Faecalibacterium 0.0000000  0.0000000  0.000000000  0.000000000  8.4710030  
## g__Akkermansia    0.0000000  0.0000000  0.000000000  0.000000000  0.0000000  
## g__Bifidobacterium 0.3756509  2.0725290  1.033117545  0.602090120  1.3608252  
## g__Escherichia    38.3172402  21.4206238  41.259008130  36.090417220  0.8564302
```

```
taxa.raw <- data.frame(t(PLEASE.raw),  
                      row.names = strtrim(samplepoint,7))
```

```
### Make sure you load the data correctly  
cat('samples', 'taxa', dim(taxa.raw), '\n')
```

```
## samples taxa 335 105
```

```
taxa.raw[1:3,1:3]
```

```
##           g__Bacteroides g__Ruminococcus g__Faecalibacterium  
## 5001-01      0.42722310      0.000000000                0  
## 5001-02      0.23188470      0.000000000                0
```

```
## 5001-03      0.02041114      0.006130341      0
```

Each row represent a sample.

Load total non-human read counts

```
human.read.file <- './real data/please_combo_human_reads.xlsx'
human.read <-
  read_excel(human.read.file, col_names = T) %>%
  as.data.frame() %>%
  mutate(
    Sample = strtrim(as.Date(Sample,origin = "1900-01-01"),7)
  )

head(human.read)
```

```
##      Sample NonHumanReads TotalReads HumanReads      HumanPer GroupFcp GroupPcdai
## 1 1910-12      17525422      19515697      1990275 10.19832959      Combo      Combo
## 2 1910-12      18089762      18185930       96168  0.52880444      Combo      Combo
## 3 1910-12      27311061      27338002       26941  0.09854781      Combo      Combo
## 4 1910-12      11051439      11092808       41369  0.37293536      Combo      Combo
## 5 1910-12       9434025       9492196       58171  0.61282980      Combo      Combo
## 6 1910-12      23327496      23634581       307085  1.29930382      Combo      Combo
```

```
## first column: sample id
```

```
### Filter low depth samples (low non human reads)
```

```
low.depth.samples <- subset(human.read,NonHumanReads<10000)
head(low.depth.samples[,1:5])
```

```
##      Sample NonHumanReads TotalReads HumanReads      HumanPer
## 56 5010-02          1014         1104          90  8.15217391
## 85 5018-03          6101          6121          20  0.32674400
## 98 5023-04          1954          2679          725 27.06233669
## 257 6007-01          1809          4249         2440 57.42527654
## 309 7001-02          9965          9971           6  0.06017451
## 310 7001-03           566           613          47  7.66721044
```

```
### Delete these samples from PLEASE data.
```

```
# row.names(taxa.raw)
```

```
# row.names(low.depth.samples)
```

```
row.names(taxa.raw)[which(row.names(taxa.raw) %in% low.depth.samples$Sample)]
```

```
## [1] "5018-03" "7001-02" "7003-02" "7003-03" "7009-03" "7010-03"
```

```
### Before deletion
```

```
dim(taxa.raw)
```

```
## [1] 335 105
```

```
### After deletion
```

```
taxa.raw <- taxa.raw[-which(rownames(taxa.raw) %in% low.depth.samples$Sample),]
dim(taxa.raw)
```

```
## [1] 329 105
```

Filter low abundant bacterial data

```
### Filter low abundant bacterial data
```

```
filter.index1 <- apply(taxa.raw,2,function(X){sum(X>0)>0.4*length(X)})
filter.index2 <- apply(taxa.raw,2,function(X){quantile(X,0.9)>1})
taxa.filter <- taxa.raw[,filter.index1 & filter.index2]
taxa.filter <- 100*sweep(taxa.filter, 1, rowSums(taxa.filter), FUN="/")
cat('after filter:', 'samples', 'taxa', dim(taxa.filter), '\n')
```

```
## after filter: samples taxa 329 18
```

```
cat(colnames(taxa.filter), '\n')
```

```
## g__Bacteroides g__Ruminococcus g__Faecalibacterium g__Bifidobacterium g__Escherichia g__Clostridium
```

```
head(rowSums(taxa.filter))
```

```
## 5001-01 5001-02 5001-03 5001-04 5002-01 5002-02
##      100      100      100      100      100      100
```

After filter, there remains 18 bacteria in the taxa table.

```
###
```

```
taxa.data <- taxa.filter
dim(taxa.data)
```

```
## [1] 329 18
```

Load sample information

```
head(sample.info)
```

```
## Sample Subject Species.Cluster Cluster Treatment FCPResponse Type Time
## 1 4000 4000 cluster 1 cluster 1 NA NA COMBO NA
## 2 4001 4001 cluster 1 cluster 1 NA NA COMBO NA
## 3 4002 4002 cluster 1 cluster 1 NA NA COMBO NA
## 4 4004 4004 cluster 1 cluster 1 NA NA COMBO NA
## 5 4005 4005 cluster 1 cluster 1 NA NA COMBO NA
## 6 4006 4006 cluster 1 cluster 1 NA NA COMBO NA
## BristolScore FCP PCDAI PUCAI log.FCP Group Response Antibiotics.visit
## 1 NA NA NA NA NA COMBO NA Not.Use
## 2 NA NA NA NA NA COMBO NA Not.Use
## 3 NA NA NA NA NA COMBO NA Not.Use
## 4 NA NA NA NA NA COMBO NA Not.Use
## 5 NA NA NA NA NA COMBO NA Not.Use
## 6 NA NA NA NA NA COMBO NA Not.Use
## Steroids Treatment.Specific Disease NonHumanReads Human.Per
## 1 NA NA Control 17525422 10.19832959
## 2 NA NA Control 18089762 0.52880444000000004
## 3 NA NA Control 27311061 9.85478050000000002E-2
## 4 NA NA Control 11051439 0.372935363000000002
## 5 NA NA Control 9434025 0.61282980399999998
## 6 NA NA Control 23327496 1.299303817
## Fungi.Per Distance Bact.Div
## 1 1.02707940499236E-4 0.452380952380952 133.36064650458701
## 2 7.5180646378874396E-3 0.42857142857142899 173.757418771479
## 3 8.2750355249838195E-4 0.40476190476190499 90.479491742827904
## 4 3.5651465840783299E-3 0.35714285714285698 103.213875049946
## 5 1.5984693701786901E-2 0.5 149.26361597807599
```

```
## 6 6.1986935931743403E-3 0.28571428571428598 112.488017364215
## Species.Distance
## 1 0.40449438202247201
## 2 0.33707865168539303
## 3 0.33707865168539303
## 4 0.426966292134831
## 5 0.43820224719101097
## 6 0.235955056179775
```

create covariates, Time, Treatment(antiTNF+EEN)

```
complete_subject =
  sample.info %>%
  filter(Sample %in% rownames(taxa.data)) %>%
  filter(Treatment.Specific!='PEN') %>%
  dplyr::select(Sample,Time,Subject,Response,Treatment.Specific) %>%
  group_by(Subject) %>%
  summarise(count = n()) %>%
  filter(count==4)

reg.cov =
  sample.info %>%
  filter(Subject %in% complete_subject$Subject) %>%
  mutate(Treat=ifelse(Treatment.Specific=='antiTNF',1,0)) %>%
  dplyr::mutate(Subject=paste('S',Subject,sep='')) %>%
  dplyr::mutate(Time=ifelse(Time=='1',0,ifelse(Time=='2',1,ifelse(Time=='3',4,ifelse(Time=='4',8,NA)))) %>%
  dplyr::mutate(Time.X.Treatment=Time*Treat) %>%
  dplyr::select(Sample,Subject,Time,Response,Treat,Time.X.Treatment,everything())
```

take out first time point

```
reg.cov.t1 <- subset(reg.cov,Time==0)
rownames(reg.cov.t1) <- reg.cov.t1$Subject

reg.cov.t234 <- subset(reg.cov,Time!=0)
reg.cov.t234 <- data.frame(
  baseline.sample=reg.cov.t1[reg.cov.t234$Subject,'Sample'],
  baseline.subject=reg.cov.t1[reg.cov.t234$Subject,'Subject'],
  reg.cov.t234,
  stringsAsFactors = FALSE)

head(reg.cov.t234)
```

##	baseline.sample	baseline.subject	Sample	Subject	Time	Response	Treat
## 2	5001-01	S5001	5001-02	S5001	1	Non.Response	1
## 3	5001-01	S5001	5001-03	S5001	4	Non.Response	1
## 4	5001-01	S5001	5001-04	S5001	8	Non.Response	1
## 6	5002-01	S5002	5002-02	S5002	1	Non.Response	1
## 7	5002-01	S5002	5002-03	S5002	4	Non.Response	1
## 8	5002-01	S5002	5002-04	S5002	8	Non.Response	1
##	Time.X.Treatment	Species.Cluster	Cluster	Treatment	FCPResponse	Type	
## 2	1	cluster 2	cluster 2	antiTNF	0	PLEASE-T2	
## 3	4	cluster 2	cluster 2	antiTNF	0	PLEASE-T3	
## 4	8	cluster 2	cluster 2	antiTNF	0	PLEASE-T4	

```

## 6      1      cluster 2 cluster 1      antiTNF      0 PLEASE-T2
## 7      4      cluster 2 cluster 1      antiTNF      0 PLEASE-T3
## 8      8      cluster 2 cluster 1      antiTNF      0 PLEASE-T4
## BristolScore FCP PCDAI PUCAI      log.FCP Group Antibiotics.visit
## 2      6 607      NA      25 6.4085287910595001 PLEASE      Not.Use
## 3      6 867      NA      20 6.7650389767805397 PLEASE      Not.Use
## 4      6 557      5      15 6.3225652399272798 PLEASE      Not.Use
## 6      6 950      NA      10 6.8564619845945902 PLEASE      Not.Use
## 7      6 1947     NA      50 7.5740450053722004 PLEASE      Not.Use
## 8      6 1880     35      40 7.5390270558239996 PLEASE      Not.Use
## Steroids Treatment.Specific Disease NonHumanReads      Human.Per
## 2 Not.Use      antiTNF Crohn      1350309      89.31803171
## 3 Not.Use      antiTNF Crohn      10946591 19.6891700000000001
## 4 Not.Use      antiTNF Crohn      14230882 0.851336733999999998
## 6 Use      antiTNF Crohn      12020377      17.08929796
## 7 Use      antiTNF Crohn      1910666 88.5445408399999996
## 8 Use      antiTNF Crohn      606565      89.5498726
## Fungi.Per      Distance      Bact.Div
## 2 7.0946724046125703E-2 0.69047619047619002 79.225334004595695
## 3 1.2168171808008501E-2 0.52272727272727304 66.915876889694502
## 4 1.26499538117174 0.59090909090909105 53.865413747151202
## 6 1.05154771767974E-2 0.42857142857142899 81.330542193164007
## 7 6.6678320543726605E-2 0.59523809523809501 89.862102084722594
## 8 5.6712800771557902E-2 0.66666666666666696 61.502647827475897
## Species.Distance
## 2 0.82417582417582402
## 3 0.72527472527472503
## 4 0.733333333333333295
## 6 0.59550561797752799
## 7 0.75280898876404501
## 8 0.85393258426966301

```

```

taxa_all = colnames(taxa.data)
store = function(taxa){
  # X: Baseline abundance time Treat
  # Y: Response abundance at time 1 4 8
  X <- data.frame(
    Baseline=taxa.data[reg.cov.t234$baseline.sample,taxa]/100,
    reg.cov.t234[,c('Time','Treat')])

  rownames(X) <- reg.cov.t234$Sample
  Y <- data.frame(
    Abundance=taxa.data[reg.cov.t234$Sample, taxa]/100,
    reg.cov.t234[,c('Time','Treat')])

  return(list(X = X,
             Y = Y))
}

store_results = lapply(taxa_all, store)
names(store_results) = taxa_all
#example

```

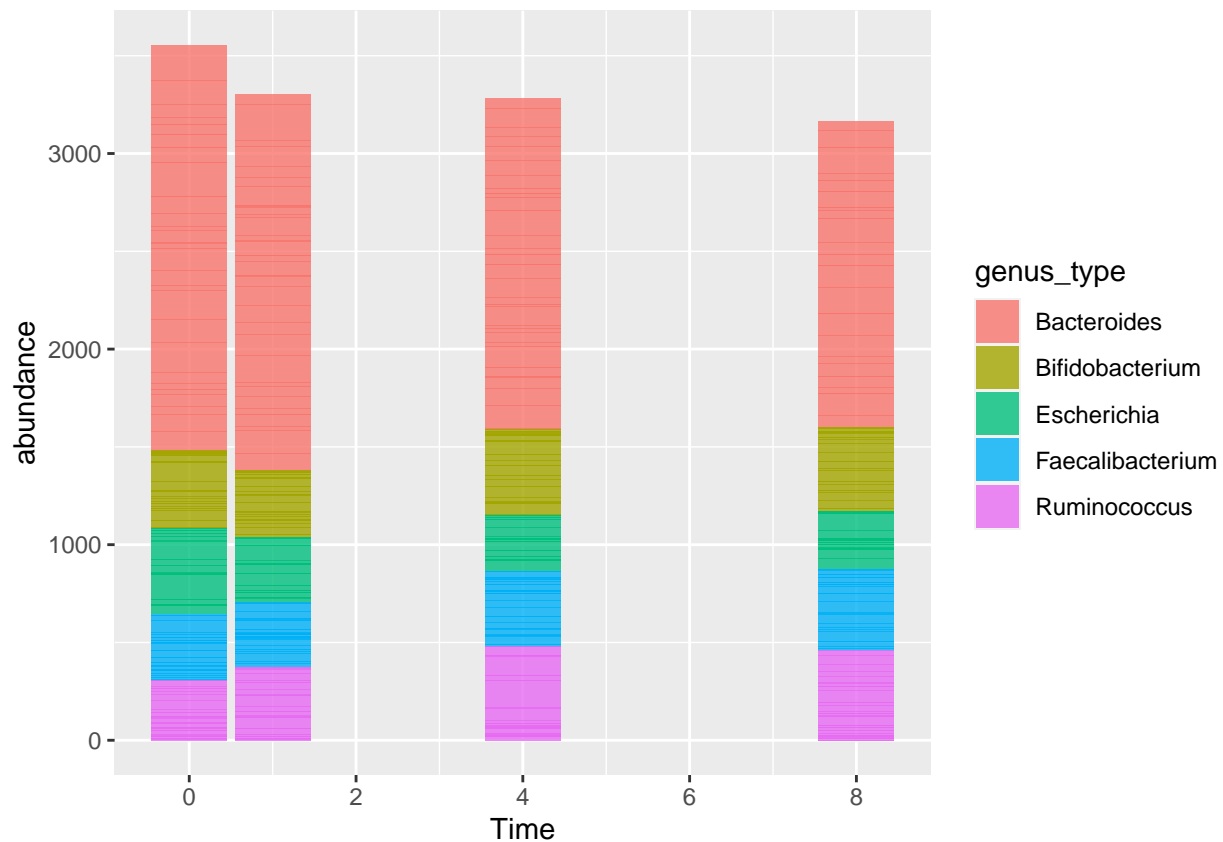
```
head(store_results$g__Bacteroides$X)
```

```
##           Baseline Time Treat
## 5001-02 0.004514608     1     1
## 5001-03 0.004514608     4     1
## 5001-04 0.004514608     8     1
## 5002-02 0.732645555     1     1
## 5002-03 0.732645555     4     1
## 5002-04 0.732645555     8     1
```

```
all_set = cbind(reg.cov, taxa.data[reg.cov$Sample,])
baseline_set = cbind(reg.cov.t1, taxa.data[reg.cov.t1$Sample,])
visit_set = cbind(reg.cov.t234, taxa.data[reg.cov.t234$Sample,])
```

```
test_genus_Eriz =
  all_set %>%
  pivot_longer(
    g__Bacteroides : g__Escherichia,
    names_to = "genus_type",
    names_prefix = "g__",
    values_to = "abundance"
  )
```

```
test_genus_Eriz %>%
  ggplot(aes(x = Time, y = abundance)) +
  geom_bar(stat="identity", aes(x = Time, y = abundance, fill = genus_type), alpha = 0.8)
```



```
save.image(file = "Explore_EricZ.RData")
```