

genBart package Vignette

Jacob Cardenas

2018-02-14

BART (Biostatistical Analysis Reporting Tool) is a user friendly, point and click, R shiny application. With the genBart package, biostatisticians/bioinformaticians analyze high throughput biological data obtained from RNA-Seq, Microarray, Flow Cytometry, and metabolomic experiments and upload the result into BART. The tool allows the R programmer who conducted the analysis and their colleagues to efficiently examine the results. BART provides users the ability to easily view, modify and download the tables and figures generated by the app. BART offers one reporting source for all analyses in a project workflow.

- Sample information (e.g. technical, biological, clinical, and demographics) summary statistics
- Sample quality control metrics
- Unsupervised analysis (heatmap and cluster analysis)
- Differential expression (sortable genelists, venn diagrams)
- Gene set analysis
- Correlation analysis within and between biological data types

In order to effectively use BART, we introduce the genBart package. genBart is a series of functions that transforms R objects generated from statistical and modeling packages such as limma, DESeq2, and edgeR into files that are uploaded into BART. In the following sections, we illustrate how to use the functions available in genBart by walking through the analysis of a longitudinal microarray study tracking gene expression changes in 38 healthy cynomolgus macaques infected with *M. tuberculosis* (Skinner et al.). Briefly, this study consisted of data measurements taken at baseline (before infection) and at 11 additional timepoints, up to 6 months post infection. Each macaque was identified as developing symptomatic (active) TB or asymptomatic (latent) infection. For the purposes of illustration and time, a subset of the microarray data was chosen for analysis. We randomly selected 4000 probes and 10 macaques and include 3 time points (days 0, 20, and 42). In addition to microarray data, flow cytometry data consisting of 13 variables was collected on 19 of the 38 macaques. Flow cytometry data analyses results are incorporated to demonstrate the flexibility of the BART app.

Summary of genBart package

Table 1: Functions in the genBart package ordered to demonstrate a typical analysis workflow.

Function	Purpose
<code>metaData()</code>	Aligns design and expression files. Defines elements of design for further analyses.
<code>normalizeData()</code>	Normalizes data.
<code>clusterData()</code>	Clusters normalized expression data.
<code>genModScores()</code>	Generates sample level module scores for plotting.
<code>genModelResults()</code>	Processes and formats differential expression modeling results.
<code>runQgen()</code>	Runs gene set analysis using Q-Gen (generalized QuSAGE) algorithm.
<code>crossCorr()</code>	Runs pairwise cross platform correlations (e.g. metabolites vs gene expression).
<code>genFile()</code>	Generates and saves BART ready R data file.
<code>updateFile()</code>	Updates/adds analyses to existing BART file.
<code>runBart()</code>	Runs BART shiny app.

The genBart package contains 11 functions, each of which is listed in the table above. Greater detail on each of these functions, their parameters, and how they are integrated into a typical workflow is available via the

help commands in R.

Design Information

Before running any analysis, the first step to generate a BART ready file is defining the design and sample information elements. Keep in mind that not every element is necessary for the rest of the genBart functions to run, so in practice, much of this section is left to user discretion and the particular project at hand.

```
library(genBart)
library(limma)
```

Below are the first six lines of the design file from the TB microarray study described in the introduction.

```
head(tb.design)
#>      columnname monkey_id timepoint sample_id clinical_status
#> 1  AVG_Signal__7196763044_K      M1         0    M1-Pre1      Active
#> 5  AVG_Signal__6303256034_B      M1        20    M1-D20      Active
#> 7  AVG_Signal__6303256020_I      M1        42    M1-D42      Active
#> 22 AVG_Signal__7196763078_I     M11         0   M11-Pre1     Latent
#> 26 AVG_Signal__7196771011_D     M11        20   M11-D20     Latent
#> 28 AVG_Signal__7196763087_B     M11        42   M11-D42     Latent
#>      timerange      Group
#> 1           Pre  Active0
#> 5         Early Active20
#> 7       Middle Active42
#> 22          Pre  Latent0
#> 26       Early Latent20
#> 28       Middle Latent42
```

The `metaData` function serves two general purposes. The first is to ensure that the design and expression data align by running through a series of checks (i.e. equal number of samples and same sample names). If any of the checks fail, the user is notified with a warning or message and a hint at where the issue is. The `columnname` argument is the key to running these checks since it specifies the column name in the design that contains the column names of the expression. If `columnname` is not specified or misspecified, then `metaData` has no point of reference to run the checks. The second purpose is to store a list of user defined design and sample information parameters that can be used for downstream analyses and in the BART app itself. In our example, since the study design is longitudinal, we can specify `long = TRUE` and declare the column in the design that contains the baseline value (`baseline.var = "timepoint"`, `baseline.val = 0`). The complete function call would be:

```
meta <- metaData(y = tb.expr, design = tb.design, data.type = "microarray",
               columnname = "columnname", long = TRUE, sample.id = "sample_id",
               subject.id = "monkey_id", time.var = "timepoint",
               baseline.var = "timepoint", baseline.val = 0)
```

We proceed similarly for the flow data:

```
meta.flow <- metaData(y = tb.flow, design = tb.flow.des, data.type = "flow",
                    columnname = "columnname", long = TRUE, sample.id = "sample_id",
                    subject.id = "monkey_id", time.var = "timepoint",
                    baseline.var = "timepoint", baseline.val = 0)
```

Please refer to the R help for `metaData` for more detail on function parameters.

Expression Normalization and Clustering

Since BART's primary function is to serve as a reporting tool for analyzed data, we normalize and cluster the

data beforehand so that users can efficiently sort through various heat maps without having to wait for the computationally intensive task of clustering thousands of genes. In `genBart`, normalization and hierarchical clustering are completed by two simple functions. First, the object generated by `metaData` containing the expression and design information is input into `normalizeData`. This function normalizes the data in various ways, depending on the study design specified in `metaData`. In our TB microarray example with repeated measurements, `normalizeData` returns three normalized datasets: all samples normalized to the mean or median (specified by `norm.method` argument), baseline samples normalized to the mean or median, and all samples normalized to baseline. For baseline normalization, instead of subtracting by the mean or median of the baseline samples, each sample's own baseline is subtracted. The list returned by `normalizeData` is then input into `clusterData`, which performs hierarchical clustering on each of the normalized datasets and returns a list of dendrograms. For more detail on the various normalization and hierarchical clustering methods, please refer to the R help for `normalizeData` and `clusterData`.

```
norm.data <- normalizeData(meta = meta, norm.method = "mean")
cluster.data <- clusterData(norm.data = norm.data, dist.method = "euclidean",
                             agg.method = "complete")
#> [1] "clustering genes from baseline samples normalized to mean..."
#> [1] "clustering genes from all samples normalized to mean..."
#> [1] "clustering genes from all samples normalized to baseline..."
```

Generate Gene Set Figures

It is often of interest to examine the behavior of genes that have been grouped together based on similar biological function or other metric. The `genModScores` function calculates the percentage of genes within a gene set that are up or down regulated with respect to baseline and/or a set of controls. These *up* or *down* percentages, referred to as *module scores*, are calculated for every gene set and sample in the study and are plotted within BART as a heatmap. As an example, we form 10 gene sets defined by clusters from hierarchical clustering and compute their sample level module scores.

```
# Form clusters
tb.expr.scale <- data.frame(t(scale(t(tb.expr)))) # center and scale probes
hc <- fastcluster::hclust(dist(tb.expr.scale))
cls <- cutree(hc, 10)
clusters <- list()
for(i in 1:10){
  clusters[[i]] <- names(cls)[which(cls == i)]
  names(clusters)[i] <- paste0("C",i)
}

# Generate module scores
mod.scores <- genModScores(meta = meta, gene.sets = clusters, sd.lim = 2)
```

Refer to the R help for `genModScores` for more detail on function parameters and how the module scores are calculated.

Differential Analysis

Differential expression analysis can contain numerous comparisons, making it time consuming to sort through the results. BART provides a solution through an interface in which users can efficiently sort through differential results across all comparisons. The `genModelResults` function acts as a bridge between BART and common differential analysis pipelines found in R packages `limma`, `DESeq2`, and `edgeR`. The function takes result objects returned by any one of these pipelines and creates a data frame of results properly formatted for BART. There are a few things to note regarding `genModelResults` arguments: 1) `method` must be correctly specified as one of “limma”, “deseq2”, or “edgeR” in order for the function to work properly 2) `object` argument generated by `DESeq2` or `edgeR` must be wrapped in a list (not necessary if using `limma`) 3)

lm.Fit argument only necessary when `method = "limma"` and `data.type = "microarray"` or `"rnaseq"`. For more details about linear modeling in limma, DESeq2, and edgeR, please refer to the following user guides: 1) limma: Linear Models for Microarray and RNA-Seq Data 2) Analyzing RNA-seq data with DESeq2 3) edgeR: differential expression analysis of digital gene expression data.

Sample Differential Expression Analysis with limma

Now we demonstrate a quick walkthrough of differential expression analysis of microarray and flow cytometry data using limma.

Microarray

Create a grouping factor by combining clinical status and time point.

```
tb.design$Group <- paste(tb.design$clinical_status, tb.design$timepoint, sep = "")
grp <- factor(tb.design$Group)
```

Create a design matrix that includes separate coefficients for each level within the grouping factor so that the desired differences can be extracted as contrasts.

```
design2 <- model.matrix(~0+grp)
colnames(design2) <- levels(grp)
```

Since there are repeated measurements on each monkey, we treat monkey_id as a random effect and estimate the correlation between measurements on the same subject. In limma, this is done using duplicateCorrelation.

```
dupcor <- duplicateCorrelation(tb.expr, design2, block = tb.design$monkey_id)
```

Fit model and extract contrasts of interest.

```
fit <- lmFit(tb.expr, design2, block = tb.design$monkey_id,
            correlation = dupcor$consensus.correlation)
contrasts <- makeContrasts(A_20vsPre = Active20-Active0, A_42vsPre = Active42-Active0,
                          L_20vsPre = Latent20-Latent0, L_42vsPre = Latent42-Latent0,
                          levels=design2)

fit2 <- contrasts.fit(fit, contrasts)
fit2 <- eBayes(fit2, trend = FALSE)
```

Flow Cytometry

Analysis of flow data in limma proceeds similarly.

```
tb.flow.des$Group <- paste(tb.flow.des$clinical_status, tb.flow.des$timepoint, sep = "")
grp <- factor(tb.flow.des$Group)

design2 <- model.matrix(~0+grp)
colnames(design2) <- levels(grp)
tb.flow.l2 <- log2(tb.flow + 1) # log2 transform
dupcor <- duplicateCorrelation(tb.flow.l2, design2, block = tb.flow.des$monkey_id)
fit.flow <- lmFit(tb.flow.l2, design2, block = tb.flow.des$monkey_id,
                 correlation = dupcor$consensus.correlation)

# Additional contrasts since all timepoints available.
contrasts <- makeContrasts(A_20vsPre = Active20-Active0, A_30vsPre = Active30-Active0,
                          A_42vsPre = Active42-Active0, A_56vsPre = Active56-Active0,
                          L_20vsPre = Latent20-Latent0, L_30vsPre = Latent30-Latent0,
```

```

L_42vsPre = Latent42-Latent0, L_46vsPre = Latent56-Latent0,
levels=design2)

fit2.flow <- contrasts.fit(fit.flow, contrasts)
fit2.flow <- eBayes(fit2.flow, trend = FALSE)

```

Generate Model Results File

The next step in the process is generating a model results file that is formatted specifically for BART. This can be done using `genModelResults`, which takes as input the expression data used for modeling, as well as model result output from `limma`, `DESeq2`, or `edgeR`. In our example, we ran linear models using `limma` for both the microarray and flow data.

```

data(gene.symbols) # gene symbols for Illumina probes
mod.results <- genModelResults(y = tb.expr, data.type = "microarray", object = fit2,
                              lm.Fit = fit, method = "limma", var.symbols = gene.symbols)
mod.results.flow <- genModelResults(y = tb.flow, data.type = "flow", object = fit2.flow,
                                   lm.Fit = fit.flow, method = "limma")

```

Gene Set Analysis: runQgen

Next we run gene set analysis by incorporating functions available in the `qusage` package, which tests whether the average log2 fold change of the genes within a gene set are different than 0. `runQgen` simplifies the process by requiring two input parameters, the model results object produced by `genModelResults` and a list of gene sets.

```

qus <- runQgen(model.results = mod.results, gene.sets = modules)
#> Aggregating gene data for gene sets.Done.
#> Calculating VIF's on residual matrix.
#> Q-Gen analysis complete.Aggregating gene data for gene sets.Done.
#> Calculating VIF's on residual matrix.
#> Q-Gen analysis complete.Aggregating gene data for gene sets.Done.
#> Calculating VIF's on residual matrix.
#> Q-Gen analysis complete.Aggregating gene data for gene sets.Done.
#> Calculating VIF's on residual matrix.
#> Q-Gen analysis complete.

```

Correlations

In addition to providing tools for standard analysis, BART also offers a tool that allows users to visualize and sort through potentially hundreds of thousands pairwise correlations across multiple platforms (e.g. gene expression vs metabolites) or versus clinical outcomes. `crossCorr` takes two numeric data frames and outputs a long format file of all pairwise correlations. The function has an option to run correlations by a grouping factor and various additional parameters for labeling in BART. We walk through an example in which we wish to correlate gene expression with flow variables. Since this is a longitudinal setting, we correlate by time.

```

# Format expression data to align with flow data
gene.data <- data.frame(t(tb.expr))
rownames(gene.data) <- paste0(tb.design$monkey_id, "_", tb.design$timepoint)
flow.data <- data.frame(t(tb.flow))
flow.data <- flow.data[match(rownames(gene.data), rownames(flow.data), nomatch = 0), ]
gene.data <- gene.data[match(rownames(flow.data), rownames(gene.data), nomatch = 0), ]

# Create time variable

```

```
time <- tb.flow.des$timepoint[match(rownames(flow.data),tb.flow.des$columnname,nomatch = 0)]

# Run correlations formatted for BART
corrs <- crossCorr(x = gene.data, y = flow.data, by = time, by.name = "days",
  description = "Genes vs Flow", x.var = "Genes",
  y.var = "Flow", method = "spearman")
```

Generate BART file

The last step in the process is to generate the final file that can be uploaded into BART. This is done through the function `genFile`, which takes as its arguments the objects generated from the previous functions. Not every object is required to generate a BART file. For example, if gene set analysis and correlations are not run, `metaData`, `clusterData`, `genModScores`, and `genModelResults` objects can still be run through `genFile`. BART will only show the results that are input.

```
genFile(meta = list(meta, meta.flow), module.scores = mod.scores,
  dendrograms = cluster.data, model.results = list(mod.results, mod.results.flow),
  project.name = "BART example")
```

In our example, we did run gene set analysis and correlations. Adding these results into BART is made simple through the `updateFile` function. The user must simply provide the path to the BART file that needs to be updated and add the new additions in the same way that `genFile` takes them.

```
path <- paste0(getwd(), "/", "BART example")
updateFile(load.path = path, qusage.results = qus, corr.results = list(corrs))
```

Run BART App

Now that a BART file has been created, the application can be run with a call to `runBart`. Please note that the default parameters for figures and tables (e.g. height, width, labeling, significance thresholds, etc.) are not optimized for any given project. As a result, it is likely that initial visualizations are far from expected or desired. However, default parameters can be easily adjusted via various widgets provided in the tool.

```
runBart()
```