

Case-Control with Gestational Age Matching

Christof Seiler

12 February 2020

Contents

1	Goal	2
2	Prerequisites	2
3	Import Data	2
4	Quality Control	7
5	Preprocessing.	8
6	Data Exploration	9
7	Differential Expression Analyses	11
8	Patient Variability	15
9	Power Analysis	18
10	Gene Set Enrichment Analysis	20
	Session Info	20

1 Goal

1. Load and normalize data using `oligo`
2. Differential analysis using `limma`

2 Prerequisites

Install necessary packages from bioconductor repository. Run this code only once to install packages.

```
pkgs_needed = c("oligo", "limma", "hta20transcriptcluster.db", "pd.hta.2.0",
  "affycoretools", "genefilter", "MatchIt", "optmatch", "ggfortify",
  "magrittr", "statmod", "readr", "dplyr", "readxl", "stringr",
  "tibble", "ggrepel", "tidyr", "locfdr")
letsinstall = setdiff(pkgs_needed, installed.packages())
if (length(letsinstall) > 0) {
  source("http://bioconductor.org/biocLite.R")
  biocLite(letsinstall)
}
```

Load packages.

```
library("oligo")
library("limma")
library("hta20transcriptcluster.db")
library("affycoretools")
library("genefilter")
library("MatchIt")
library("ggfortify")
library("magrittr")
library("statmod")
library("readr")
library("dplyr")
library("readxl")
library("stringr")
library("tibble")
library("ggrepel")
library("tidyr")
library("locfdr")
```

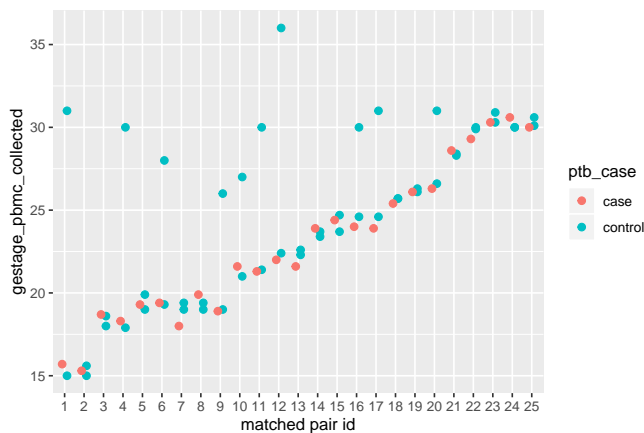
3 Import Data

Read sample tables. Merge with sample information table.

```
sample_table = read_csv("sample_table_from_word.csv")
## Parsed with column specification:
## cols(
##   ptidno = col_double(),
##   visitcode = col_double(),
##   gestage_delivery = col_double(),
```

Case-Control with Gestational Age Matching

```
## gestage_pbmc_collected = col_double(),
## ptb_case = col_character(),
## ptb_casenum = col_double(),
## ptb_controlnum = col_double()
## )
new_matched = read_excel("MSS Case-Control Matches (Masked)_new_matched.xlsx")
sample_table %<>% dplyr::left_join(new_matched, by = "ptidno",
                                suffix = c("", "2"))
ggplot(sample_table, aes(as.factor(ptb_casenum), gestage_pbmc_collected,
                        color = ptb_case)) +
  geom_jitter(position=position_dodge(width = 0.5), size = 2) +
  xlab("matched pair id")
```



```
sample_table %<>% dplyr::select(ptidno,
                                ptb_case,
                                gestage_pbmc_collected,
                                gestage_delivery,
                                sample_id = SampleID)
```

Find CEL files in current folder.

```
file_names_cel = list.files("./", pattern = "CEL")
params$treatment
## [1] "stim"
if(params$treatment == "stim") {
  file_names_cel = file_names_cel[str_detect(
    string = file_names_cel,
    pattern = "Bayless_H1N1|\\+\\-\\(HTA)"]
} else {
  file_names_cel = file_names_cel[str_detect(
    string = file_names_cel,
    pattern = "Bayless_US|\\-\\-\\(HTA|\\-\\-\\(HTA|\\-\\-\\(HTA)"]
}
tb_file_name = lapply(sample_table$sample_id, function(id) {
  pattern = paste0("[ _]", id)
  name = file_names_cel[which(str_detect(file_names_cel, pattern))]
  if(length(name)==0) name = NA
})
```

Case-Control with Gestational Age Matching

```
tibble(sample_id = id, file_name = name)
}) %>% bind_rows()
sample_table %<>% dplyr::left_join(tb_file_name, by = "sample_id")
sample_table %<>% na.omit
table(sample_table$ptb_case)
##
##      case control
##      19      37
```

Match samples.

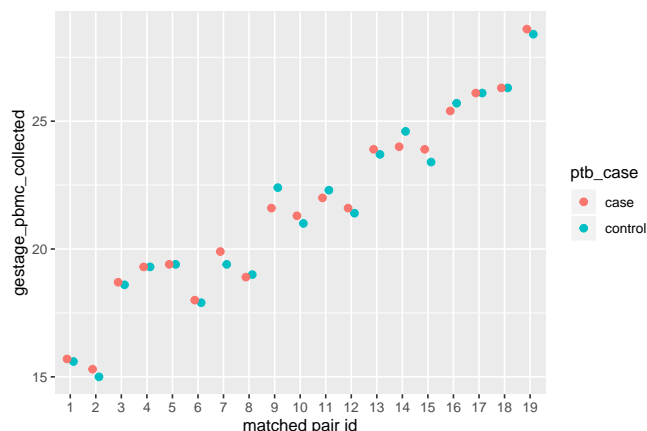
```
set.seed(0xdada2)
sample_table$group = sample_table$ptb_case == "case"
mout = matchit(group ~ gestage_pbmc_collected,
               data = sample_table,
               method = "optimal",
               ratio = 1)
## Warning in optmatch::fullmatch(d, min.controls = ratio, max.controls = ratio, : Without 'data' argument the
##      to be the same as your original data.
summary(mout)
##
## Call:
## matchit(formula = group ~ gestage_pbmc_collected, data = sample_table,
##      method = "optimal", ratio = 1)
##
## Summary of balance for all data:
##               Means Treated Means Control SD Control Mean Diff
## distance                0.3695         0.3238    0.1123    0.0457
## gestage_pbmc_collected    21.5737        23.7649    5.1037   -2.1912
##               eQQ Med eQQ Mean eQQ Max
## distance                0.0501    0.0466  0.1038
## gestage_pbmc_collected  2.2000    2.2895  7.4000
##
##
## Summary of balance for matched data:
##               Means Treated Means Control SD Control Mean Diff
## distance                0.3695         0.3701    0.0898   -0.0006
## gestage_pbmc_collected    21.5737        21.5526    3.7082    0.0211
##               eQQ Med eQQ Mean eQQ Max
## distance                0.0046    0.0058  0.0171
## gestage_pbmc_collected  0.2000    0.2421  0.7000
##
## Percent Balance Improvement:
##               Mean Diff. eQQ Med eQQ Mean eQQ Max
## distance                98.6463  90.8686  87.5475  83.5288
## gestage_pbmc_collected   99.0392  90.9091  89.4253  90.5405
##
## Sample sizes:
##               Control Treated
## All                37      19
## Matched            19      19
```

Case-Control with Gestational Age Matching

```
## Unmatched      18      0
## Discarded      0      0
sample_table$pair = NA
A = rownames(mout$match.matrix) %>% as.integer
B = mout$match.matrix %>% as.integer
for(i in 1:nrow(mout$match.matrix)) {
  sample_table$pair[A[i]] = i
  sample_table$pair[B[i]] = i
}
sample_table %<>% dplyr::select(-group)
sample_table %<>% na.omit
sample_table %>%
  dplyr::select(sample_id,ptidno,pair,ptb_case,gestage_pbmc_collected,
                gestage_delivery) %>%
  arrange(pair,ptb_case) %>%
  print(n = Inf)
## # A tibble: 38 x 6
##   sample_id ptidno pair ptb_case gestage_pbmc_collected gestage_delivery
##   <chr>      <dbl> <int> <chr>          <dbl>          <dbl>
## 1 1.2      14104760     1 case           15.7           21.3
## 2 2.2      14104580     1 control        15.6           37.4
## 3 2.1      14101980     2 case           15.3           33.6
## 4 2.3      14105160     2 control        15           39.3
## 5 3.3      14225030     3 case           18.7           33.4
## 6 3.2      14223910     3 control        18.6           38
## 7 5.3      14222170     4 case           19.3           30.9
## 8 6.3      14223900     4 control        19.3           39
## 9 6.1      14220060     5 case           19.4           29.6
## 10 7.3      14104410     5 control        19.4           43.6
## 11 7.2      14100460     6 case           18            32.9
## 12 4.3      14102700     6 control        17.9           40
## 13 8.3      14224240     7 case           19.9           23
## 14 8.1      14104890     7 control        19.4           41
## 15 9.2      14102840     8 case           18.9           24.1
## 16 9.1      14102770     8 control        19            38.3
## 17 10.1     14223520     9 case           21.6           27.9
## 18 12.3     14223850     9 control        22.4           41.1
## 19 11.2     14102560    10 case           21.3           31.9
## 20 10.3     14224520    10 control        21            37.6
## 21 12.1     14220350    11 case           22            27.3
## 22 13.1     14100770    11 control        22.3           39.1
## 23 13.2     14224670    12 case           21.6           26.7
## 24 11.1     14104710    12 control        21.4           39.9
## 25 14.3     14222790    13 case           23.9           31
## 26 14.2     14105140    13 control        23.7           43
## 27 16.1     14102780    14 case           24            33.7
## 28 17.3     14100120    14 control        24.6           39.6
## 29 17.2     14105630    15 case           23.9           31.3
## 30 14.1     14100400    15 control        23.4           39.6
## 31 18.3     14222340    16 case           25.4           27.6
## 32 18.1     14104350    16 control        25.7           38.4
```

Case-Control with Gestational Age Matching

```
## 33 19.1      14103660      17 case      26.1      29.1
## 34 19.3      14102370      17 control    26.1      41
## 35 20.2      14100230      18 case      26.3      26.7
## 36 19.2      14225270      18 control    26.3      39
## 37 21.3      14106020      19 case      28.6      31.7
## 38 21.2      14220100      19 control    28.4      39
write_csv(sample_table, path = paste0("sample_table_matched_", params$treatment, ".csv"))
ggplot(sample_table, aes(as.factor(pair), gestage_pbmc_collected,
                        color = ptb_case)) +
  geom_jitter(position=position_dodge(width = 0.5), size = 2) +
  xlab("matched pair id")
```



Then load Affymetrix CEL files. At this stage, Bioconductor will automatically download the necessary annotation packages and install them for us.

```
pd = as(as.data.frame(sample_table), "AnnotatedDataFrame")
rawData = read.celfiles(sample_table$file_name,
                        phenoData = pd,
                        sampleNames = sample_table$sample_id)

## Loading required package: pd.hta.2.0
## Loading required package: RSQLite
## Loading required package: DBI
## Platform design info loaded.
## Reading in : Nicholas Bayless_H1N1 1.2_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 2.1_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 2.3_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 2.2_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 3.3_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 3.2_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 4.3_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 5.3_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 6.1_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 6.3_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_7.2 +_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_7.3 +_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_8.3 +_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_8.1 +_(HTA-2_0).CEL
```

Case-Control with Gestational Age Matching

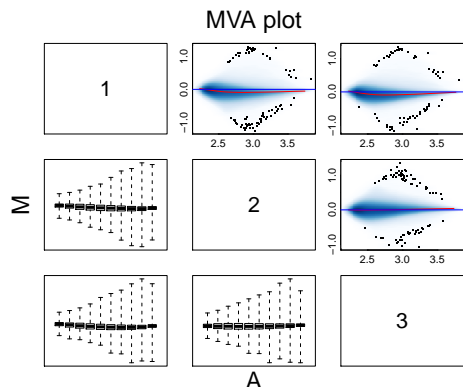
```
## Reading in : Nicholas Bayless_9.2 +_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_9.1 +_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_10.1 +_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_10.3 +_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_11.2 +_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_11.1 +_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_12.1 +_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_12.3 +_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 13.2_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 13.1_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 14.3_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 14.2_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 14.1_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 16.1_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 17.2_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 17.3_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 18.3_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 18.1_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 19.1_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 19.3_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 19.2_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 20.2_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 21.3_(HTA-2_0).CEL
## Reading in : Nicholas Bayless_H1N1 21.2_(HTA-2_0).CEL
## Warning in read.celfiles(sample_table$file_name, phenoData = pd, sampleNames
## = sample_table$sample_id): 'channel' automatically added to varMetadata in
## phenoData.
rawData
## HTAFeatureSet (storageMode: lockedEnvironment)
## assayData: 6892960 features, 38 samples
## element names: exprs
## protocolData
## rowNames: 1 2 ... 38 (38 total)
## varLabels: exprs dates
## varMetadata: labelDescription channel
## phenoData
## rowNames: 1 2 ... 38 (38 total)
## varLabels: ptidno ptb_case ... pair (7 total)
## varMetadata: labelDescription channel
## featureData: none
## experimentData: use 'experimentData(object)'
## Annotation: pd.hta.2.0
```

4 Quality Control

MA plots on the first three samples.

Case-Control with Gestational Age Matching

```
MPlot(rawData[, 1:3], pairs=TRUE)
```



5 Preprocessing

Background subtraction, normalization and summarization using median-polish.

```
eset = rma(rawData)
## Background correcting
## Normalizing
## Calculating Expression
```

Get rid of background probes and annotate using functions in `affycoretools` package.

```
dbGetQuery(db(pd.hta.2.0), "select * from type_dict;")
##   type                                     type_id
## 1     1                                           main
## 2     2               Antigenomic background control
## 3     3               control->affx->bac_spike
## 4     4               control->affx->polya_spike
## 5     5 ERCC (External RNA Controls Consortium) step control
## 6     6           Exonic normalization control (Positive Control)
## 7     7       Intronic normalization control (Negative Control)
## 8     8                               Positive Control
table(getMainProbes("pd.hta.2.0")$type)
##
##      1      2      3      4      5      6      7
## 67516   23     4     4   155   698   646
eset = getMainProbes(eset)
```

Filter probes that we cannot map to symbols.

```
e2s = toTable(hta20transcriptclusterSYMBOL)
prob_ids = rownames(exprs(eset))
keep_ids = which(prob_ids %in% e2s$probe_id)
eset = ExpressionSet(assayData = exprs(eset)[keep_ids,],
                    phenoData = phenoData(eset),
                    experimentData = experimentData(eset),
                    annotation = annotation(eset))
```


Case-Control with Gestational Age Matching

Save to file.

```
class(eset)
## [1] "ExpressionSet"
## attr(,"package")
## [1] "Biobase"
show(eset)
## ExpressionSet (storageMode: lockedEnvironment)
## assayData: 29503 features, 38 samples
## element names: exprs
## protocolData: none
## phenoData
## sampleNames: 1 2 ... 38 (38 total)
## varLabels: ptidno ptb_case ... pair (7 total)
## varMetadata: labelDescription channel
## featureData: none
## experimentData: use 'experimentData(object)'
## Annotation: pd.hta.2.0
exprs(eset)[1:10, 1:2]
##           1      2
## TC01000003.hg.1  2.122029  2.255396
## TC01000007.hg.1 10.668143 10.429572
## TC01000010.hg.1  3.908038  3.703156
## TC01000018.hg.1  6.427095  7.361067
## TC01000019.hg.1  5.635121  5.827241
## TC01000020.hg.1  5.838369  5.833540
## TC01000021.hg.1  5.611789  5.540196
## TC01000022.hg.1  5.946168  5.844652
## TC01000023.hg.1 10.394409  9.303554
## TC01000024.hg.1  6.160168  6.149659
save(eset, file = "eset.Rdata")
```

Write processed expressions to file for GEO upload.

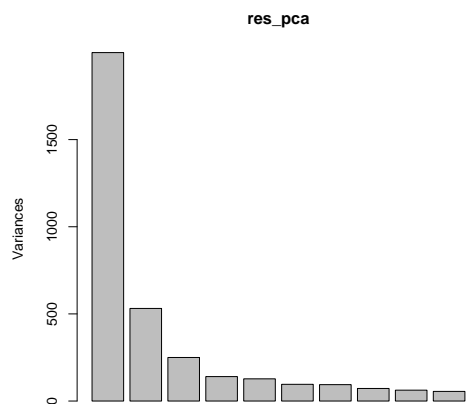
```
geo_exprs_rma = exprs(eset)
condition = ifelse(pData(eset)$ptb_case == "control",
  yes = "Term", no = "Preterm")
sample_name = paste0(pData(eset)$sample_id, "_",
  condition, "_",
  ifelse(params$treatment == "stim",
    yes = "H1N1", "Mock"))
colnames(geo_exprs_rma) = sample_name
geo_exprs_rma %<>% as_tibble(rownames = "ID_REF")
write_csv(geo_exprs_rma, path = paste0("case_control_rma_", params$treatment, ".csv"))
```

6 Data Exploration

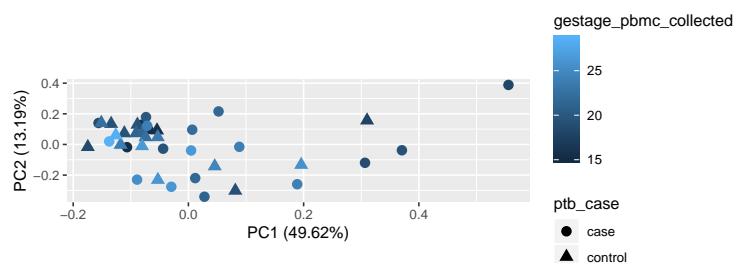
PCA plot of normalized expressions.

```
res_pca = prcomp(t(exprs(eset)), scale. = FALSE)
screeplot(res_pca)
```

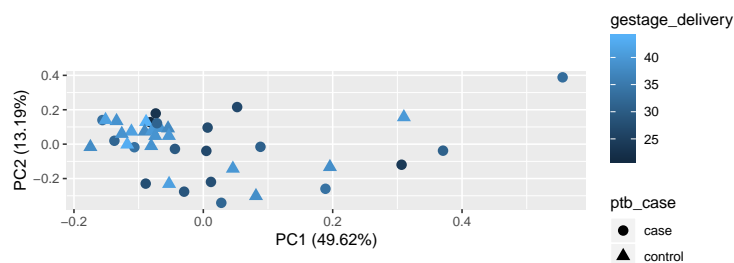
Case-Control with Gestational Age Matching



```
sample_table_annotated = eset@phenoData@data
sample_table_annotated$pair %<>% as.factor
asp_ratio = res_pca$sdev[2]^2/res_pca$sdev[1]^2
autoplot(res_pca,
  data = sample_table_annotated,
  shape = "ptb_case",
  colour = "gestage_pbmc_collected",
  size = 3,
  asp = asp_ratio)
```



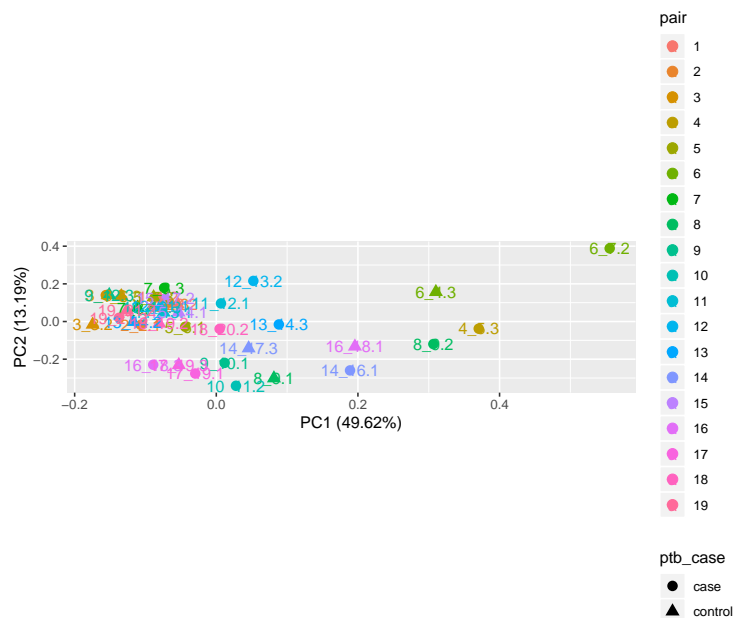
```
autoplot(res_pca,
  data = sample_table_annotated,
  shape = "ptb_case",
  colour = "gestage_delivery",
  size = 3,
  asp = asp_ratio)
```



```
rownames(sample_table_annotated) = paste0(sample_table_annotated$pair,
  "_",
  sample_table_annotated$sample_id)
autoplot(res_pca,
```

Case-Control with Gestational Age Matching

```
data = sample_table_annotated,  
shape = "ptb_case",  
colour = "pair",  
size = 3,  
label = TRUE,  
asp = asp_ratio)
```



7 Differential Expression Analyses

Use `limma` for linear models to assess difference in expression. Paired analysis as described in Section 9.4.1 on page 42 in the [limma vignette](#).

Automatic independent filtering as described in [DESeq2 doc](#):

1. Filter genes based on mean expression
2. Fit linear model
3. Compute moderated t -tests
4. Count number of rejections at FDR of 10%

Pick the threshold that maximizes the number of discoveries.

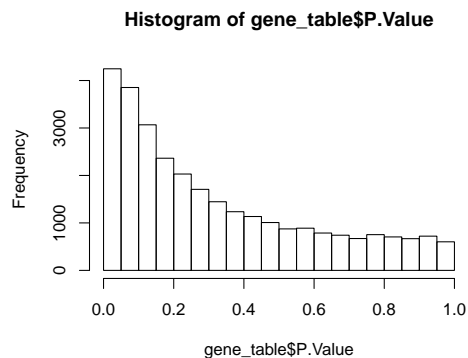
```
mean_expr = rowMeans(exprs(eset))  
thres_candidates = seq(min(mean_expr), quantile(mean_expr, probs = 0.95), 1)  
fit_list = lapply(thres_candidates, function(thres) {  
  cat("Automatic independent filtering: thres = ", thres, "\n")  
  # threshold  
  eset_thres = ExpressionSet(assayData = exprs(eset)[mean_expr >= thres,],  
                             phenoData = phenoData(eset),  
                             experimentData = experimentData(eset),  
                             annotation = annotation(eset))  
  
  # fit model
```

Case-Control with Gestational Age Matching

```
targets = eset@phenoData@data
pair = factor(targets$pair)
treat = factor(targets$ptb_case, levels = c("control", "case"))
design = model.matrix(~ pair + treat)
fit = lmFit(eset_thres, design)
eBayes(fit)
})
## Automatic independent filtering: thres = 1.397895
## Automatic independent filtering: thres = 2.397895
## Automatic independent filtering: thres = 3.397895
## Automatic independent filtering: thres = 4.397895
## Automatic independent filtering: thres = 5.397895
## Automatic independent filtering: thres = 6.397895
num_sig = sapply(fit_list, function(fit) {
  gene_table = topTable(fit, coef = "treatcase", adjust = "BH",
                        number = nrow(fit))
  gene_table %<>% dplyr::filter(adj.P.Val < 0.1)
  nrow(gene_table)
})
num_sig
## [1] 0 0 0 0 0 0
fit = fit_list[[which.max(num_sig)]]
```

The `topTable` command provides us a way of ranking genes for further evaluation. In the case below, we adjust for multiple testing by FDR.

```
gene_table = topTable(fit, coef = "treatcase", adjust = "BH",
                      number = nrow(fit))
hist(gene_table$P.Value, breaks = 20)
```

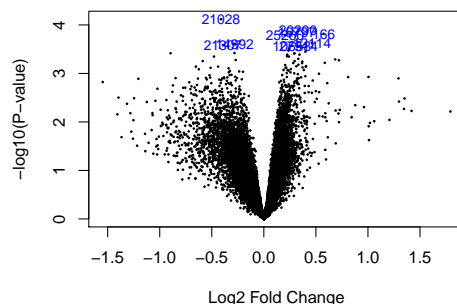


```
sum(gene_table$adj.P.Val < 0.1)
## [1] 0
```

Volcano plots for quality control.

```
volcanoplot(fit, coef = "treatcase", highlight = 10)
```

Case-Control with Gestational Age Matching



Map between manufacturer identifiers and gene symbols.

```

map_gene_symbol = function(gene_table) {
  prob_ids = rownames(gene_table)
  symbol = apply(prob_ids,function(prob_id) {
    matching_symbol = e2s$symbol[prob_id==e2s$probe_id]
    if(length(matching_symbol)==0) matching_symbol = "No_Symbol_Found"
    matching_symbol
  }) %>% unlist
  gene_table = cbind(gene_table,symbol=symbol,stringsAsFactors=FALSE)
  gene_table
}

gene_table = map_gene_symbol(gene_table)
print(head(gene_table, n = 10))

##              logFC AveExpr          t      P.Value adj.P.Val
## TC15000930.hg.1 -0.4116955 4.553749 -4.856576 7.616508e-05 0.3410067
## TC14001238.hg.1 0.3288333 4.790628 4.637631 1.295037e-04 0.3410067
## TC11001449.hg.1 0.3181510 3.176478 4.595750 1.433744e-04 0.3410067
## TC01000301.hg.1 0.5737504 6.208464 4.567539 1.535500e-04 0.3410067
## TC19000645.hg.1 0.2040581 5.445066 4.542494 1.631904e-04 0.3410067
## TC08001429.hg.1 0.4625299 3.348420 4.394225 2.340819e-04 0.3410067
## TC10000293.hg.1 -0.2801210 6.233889 -4.376315 2.445126e-04 0.3410067
## TC15001713.hg.1 -0.3917718 5.597879 -4.350500 2.603725e-04 0.3410067
## TC07001115.hg.1 0.2699228 5.593521 4.336570 2.693543e-04 0.3410067
## TC22000322.hg.1 0.3324141 3.911133 4.319523 2.807685e-04 0.3410067

##              B          symbol
## TC15000930.hg.1 0.97553299      MCTP2
## TC14001238.hg.1 0.59591834      PLEK2
## TC11001449.hg.1 0.52270786      SOX6
## TC01000301.hg.1 0.47329514      PITHD1
## TC19000645.hg.1 0.42936198      KLC3
## TC08001429.hg.1 0.16810853      MIR378D2
## TC10000293.hg.1 0.13642753      ALOX5
## TC15001713.hg.1 0.09071717      CTSH
## TC07001115.hg.1 0.06603031      ZNF890P
## TC22000322.hg.1 0.03580142      LOC100130899

```

Write to text file.

```
file_name_processed = paste0("case_control_", params$treatment, ".csv")
file_name_processed
## [1] "case_control_stim.csv"
```

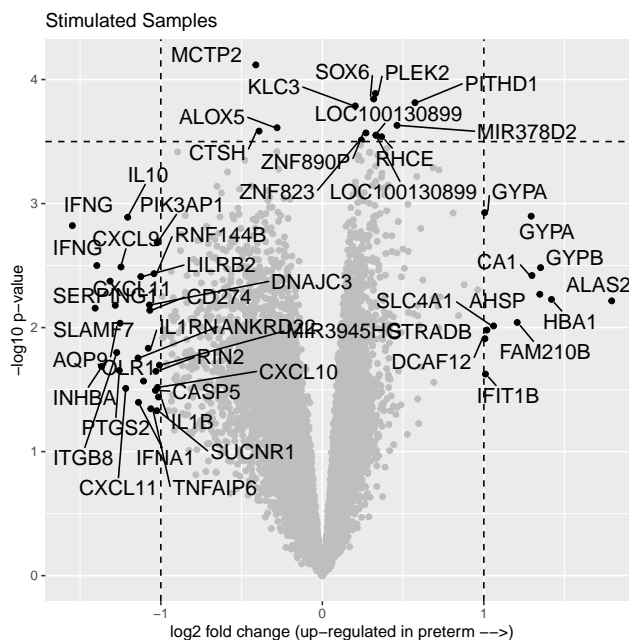
Case-Control with Gestational Age Matching

```
write_csv(gene_table, path = file_name_processed)
```

Add gene names to standard volcano plot.

```
logfc_thres = 1 # logFC threshold
pvalue_thres = 3.5 # -log10 of unadjusted p-value
gene_table %<>% mutate(
  show = ifelse(abs(logFC) > logfc_thres | pvalue_thres < -log10(P.Value),
    "yes", "no")
)
gvolcano = ggplot(gene_table, aes(logFC, -log10(P.Value), color = show)) +
  geom_point() +
  geom_vline(xintercept = c(-logfc_thres, logfc_thres), linetype = 2) +
  geom_hline(yintercept = pvalue_thres, linetype = 2) +
  geom_text_repel(
    data = dplyr::filter(gene_table, show == "yes"),
    aes(label = symbol),
    size = 5,
    box.padding = unit(0.35, "lines"),
    point.padding = unit(0.3, "lines")
  ) +
  xlab("log2 fold change (up-regulated in preterm -->)" +
  ylab("-log10 p-value" +
  theme(legend.position = "none") +
  scale_colour_manual(values = c("gray", "black")) +
  ggtitle(ifelse(params$treatment == "unstim",
    "Unstimulated Samples",
    "Stimulated Samples"))

gvolcano
```



```
save(gvolcano, file = paste0("gvolcano_", params$treatment, ".Rdata"))
```

8 Patient Variability

Visualize the pair-to-pair variability.

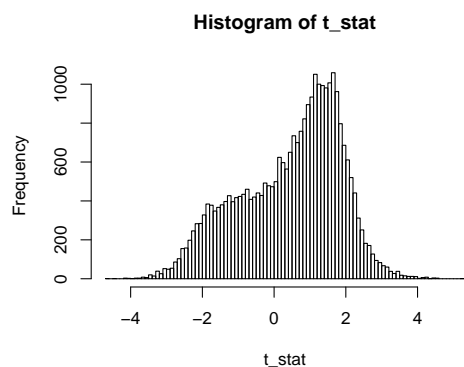
```
# prepare tables
targets = eset@phenoData@data
targets %<>% mutate(id = 1:nrow(targets))
GA_affy = read_csv("GA_affy.csv")
## Parsed with column specification:
## cols(
##   `Sample name` = col_character(),
##   title = col_character(),
##   `CEL file` = col_character(),
##   `source name` = col_character(),
##   organism = col_character(),
##   `characteristics: condition` = col_character(),
##   `characteristics: treatment` = col_character(),
##   `characteristics: rin` = col_double(),
##   `characteristics: run_day` = col_datetime(format = ""),
##   `characteristics: run_batch` = col_double(),
##   `characteristics: viable_cell_count` = col_double(),
##   `characteristics: viability` = col_double(),
##   `characteristics: ptidno` = col_double(),
##   `characteristics: gestage_delivery` = col_double(),
##   `characteristics: gestage_enroll` = col_double(),
##   molecule = col_character(),
##   label = col_character(),
##   description = col_logical(),
##   `chip name or GEO platform id` = col_character()
## )
GA_affy %<>% dplyr::filter(
  `characteristics: treatment` ==
    ifelse(params$treatment == "stim", yes = "H1N1", "Mock")
)
GA_affy %<>% dplyr::rename(ptidno = `characteristics: ptidno`)
targets %<>% left_join(GA_affy, by = "ptidno")
tb_case = targets[targets$ptb_case == "case", ]
tb_control = targets[targets$ptb_case == "control", ]
tb = left_join(tb_control, tb_case, by = "pair",
  suffix = c(".control", ".case"))

# take diff within pairs
tb_exrs = exprs(eset)
X = tb_exrs[, tb$id.control]
Y = tb_exrs[, tb$id.case]
D = Y - X

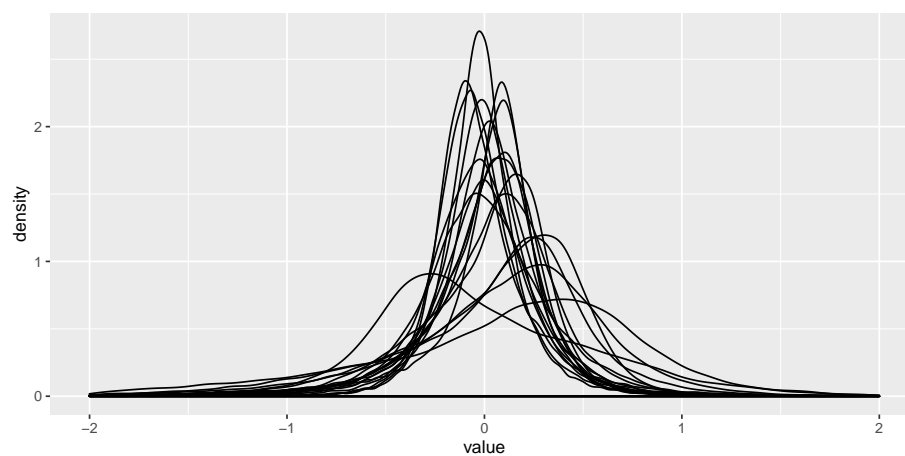
# implement my own paired t-test
```

Case-Control with Gestational Age Matching

```
t_stat = apply(D, MARGIN = 1,  
              function(x) t.test(x)$statistic)  
hist(t_stat, breaks = 100)
```

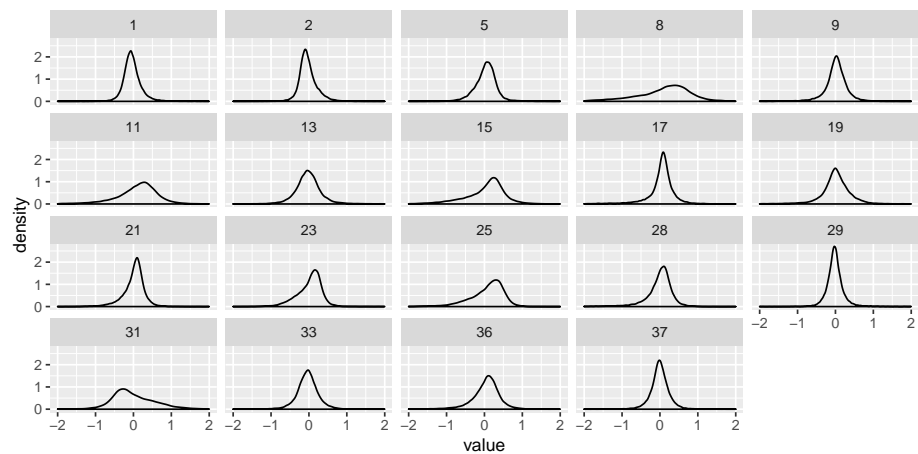


```
# plot  
tb_D = D %>% as_tibble %>% mutate(gene_id = 1:nrow(D))  
tb_D %<>% gather(pair, value, -gene_id)  
tb_D$pair %<>% factor(levels = sort(unique(as.integer(tb_D$pair))))  
ggplot(tb_D, aes(value, group = pair)) +  
  geom_density() +  
  xlim(c(-2, 2))  
## Warning: Removed 2997 rows containing non-finite values (stat_density).
```



```
ggplot(tb_D, aes(value)) +  
  geom_density() +  
  facet_wrap(~pair) +  
  xlim(c(-2, 2))  
## Warning: Removed 2997 rows containing non-finite values (stat_density).
```


Case-Control with Gestational Age Matching



The variability does not seem to be related to `run_day` and `run_batch`.

```
table(targets$pair, targets$`characteristics: run_day`)
##
##      2019-02-11 2019-02-19 2019-02-24 2019-02-26 2019-03-01 2019-03-07
##  1           2           0           0           0           0
##  2           2           0           0           0           0
##  3           2           0           0           0           0
##  4           2           0           0           0           0
##  5           1           1           0           0           0
##  6           1           1           0           0           0
##  7           0           2           0           0           0
##  8           0           2           0           0           0
##  9           0           1           1           0           0
## 10           0           1           1           0           0
## 11           0           0           1           1           0
## 12           0           0           1           1           0
## 13           0           0           0           2           0
## 14           0           0           0           0           2
## 15           0           0           0           1           1
## 16           0           0           0           0           0
## 17           0           0           0           0           2
## 18           0           0           0           0           2
## 19           0           0           0           0           2

table(targets$pair, targets$`characteristics: run_batch`)
##
##      1 2 3
##  1 2 0 0
##  2 2 0 0
##  3 0 2 0
##  4 0 0 2
##  5 1 0 1
##  6 1 1 0
##  7 2 0 0
##  8 0 2 0
##  9 1 1 0
## 10 1 1 0
```

Case-Control with Gestational Age Matching

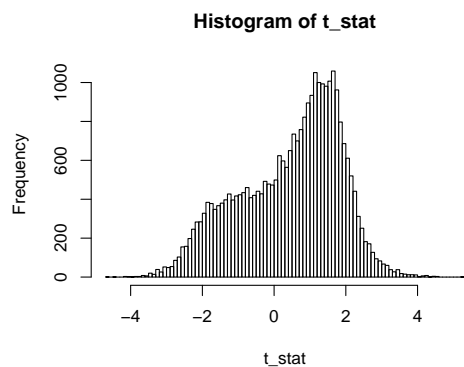
```
## 11 2 0 0
## 12 2 0 0
## 13 2 0 0
## 14 2 0 0
## 15 2 0 0
## 16 2 0 0
## 17 2 0 0
## 18 1 1 0
## 19 0 2 0
```

9 Power Analysis

Power analysis using local FDR methodology.

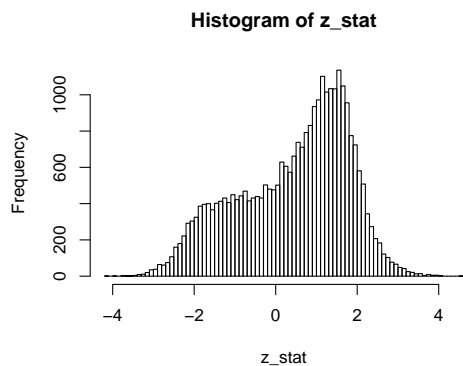
```
# fit model
targets = eset@phenoData@data
pair = factor(targets$pair)
treat = factor(targets$ptb_case, levels = c("control", "case"))
design = model.matrix(~ treat + pair)
fit = lmFit(eset, design)

# ordinary t-statistic
t_stat = fit$coef[, "treatcase"] / fit$stdev.unscaled[, "treatcase"] / fit$sigma
z_stat = qnorm(pt(t_stat, df = ncol(exprs(eset))-2))
hist(t_stat, breaks = 100)
```

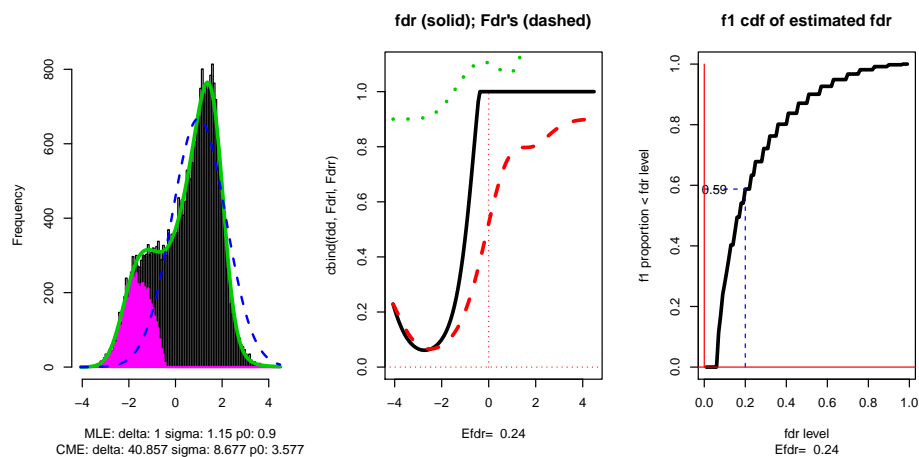


```
hist(z_stat, breaks = 100)
```

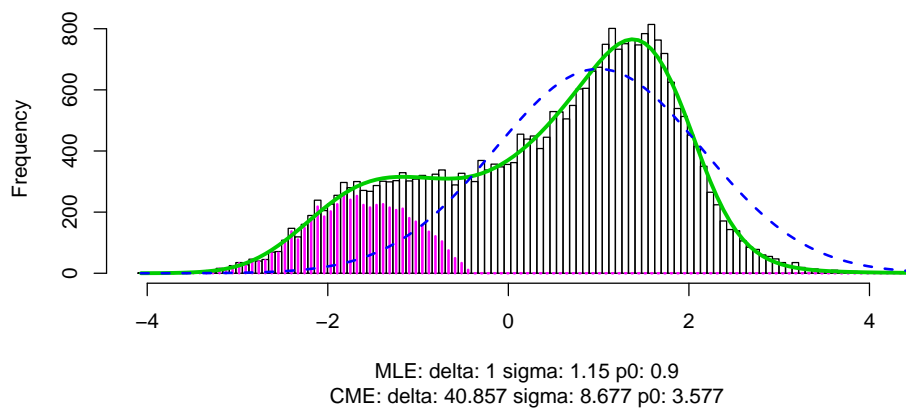
Case-Control with Gestational Age Matching



```
locfdr_res = locfdr(z_stat, df = 7, plot = 4)
## Warning in locfdr(z_stat, df = 7, plot = 4): f(z) misfit = 1.7. Rerun with
## increased df
```



```
locfdr_res$Efdr
##      Efdr      Eleft      Eright      Efdrttheo      Eleft0      Eright0
## 0.2397774 0.2397774 1.0000000 0.3355690 0.4467380 0.1891478
locfdr_res = locfdr(z_stat, df = 7, plot = 1)
## Warning in locfdr(z_stat, df = 7, plot = 1): f(z) misfit = 1.7. Rerun with
## increased df
```



Case-Control with Gestational Age Matching

Large values $E_{\text{fdr}} > 0.4$ indicate low power (according to Section 3 of “Size, power and false discovery rates”, Efron 2007).

10 Gene Set Enrichment Analysis

Standard KEGG analysis.

```
fit %<>% eBayes
pathway_results = keggga(fit, species.KEGG = "hsa")
## No DE genes
topKEGG(pathway_results)
## data frame with 0 columns and 0 rows
```

Session Info

```
sessionInfo()
## R version 3.5.1 (2018-07-02)
## Platform: x86_64-apple-darwin15.6.0 (64-bit)
## Running under: macOS 10.15.3
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/3.5/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/3.5/Resources/lib/libRlapack.dylib
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
## [1] stats4      parallel    stats       graphics    grDevices   utils       datasets
## [8] methods     base
##
## other attached packages:
## [1] pd.hta.2.0_3.12.2          DBI_1.0.0
## [3] RSQLite_2.1.1              locfdr_1.1-8
## [5] tidyr_1.0.0                ggrepel_0.8.1
## [7] tibble_2.1.3               stringr_1.4.0
## [9] readxl_1.3.1               dplyr_0.8.3
## [11] readr_1.3.1                statmod_1.4.32
## [13] magrittr_1.5               ggfortify_0.4.7
## [15] ggplot2_3.2.1              MatchIt_3.0.2
## [17] genefilter_1.64.0          affycoretools_1.54.0
## [19] hta20transcriptcluster.db_8.7.0 org.Hs.eg.db_3.7.0
## [21] AnnotationDbi_1.44.0       limma_3.38.3
## [23] oligo_1.46.0               Biostrings_2.50.2
## [25] XVector_0.22.0             IRanges_2.16.0
## [27] S4Vectors_0.20.1          Biobase_2.42.0
## [29] oligoClasses_1.44.0       BiocGenerics_0.28.0
## [31] BiocStyle_2.10.0
##
```

Case-Control with Gestational Age Matching

```
## loaded via a namespace (and not attached):
## [1] utf8_1.1.4 R.utils_2.8.0
## [3] tidyselect_0.2.5 htmlwidgets_1.3
## [5] grid_3.5.1 BiocParallel_1.16.6
## [7] munsell_0.5.0 codetools_0.2-16
## [9] preprocessCore_1.44.0 withr_2.1.2
## [11] colorspace_1.4-1 Category_2.48.1
## [13] OrganismDbi_1.24.0 knitr_1.22
## [15] rstudioapi_0.10 labeling_0.3
## [17] GenomeInfoDbData_1.2.0 hwriter_1.3.2
## [19] bit64_0.9-7 farver_2.0.1
## [21] vctrs_0.2.1 xfun_0.6
## [23] biovizBase_1.30.1 affxparser_1.54.0
## [25] R6_2.4.1 GenomeInfoDb_1.18.2
## [27] optmatch_0.9-11 locfit_1.5-9.1
## [29] AnnotationFilter_1.6.0 bitops_1.0-6
## [31] reshape_0.8.8 DelayedArray_0.8.0
## [33] assertthat_0.2.1 scales_1.1.0
## [35] nnet_7.3-12 gtable_0.3.0
## [37] affy_1.60.0 ggbio_1.30.0
## [39] svd_0.5 ensemblDb_2.6.8
## [41] rlang_0.4.2 zeallot_0.1.0
## [43] splines_3.5.1 rtracklayer_1.42.2
## [45] lazyeval_0.2.2 acepack_1.4.1
## [47] dichromat_2.0-0 checkmate_1.9.4
## [49] BiocManager_1.30.4 yaml_2.2.0
## [51] reshape2_1.4.3 abind_1.4-5
## [53] GenomicFeatures_1.34.8 backports_1.1.5
## [55] Hmisc_4.2-0 RBGL_1.58.2
## [57] tools_3.5.1 bookdown_0.9
## [59] ellipsis_0.3.0 affyio_1.52.0
## [61] gplots_3.0.1.1 ff_2.2-14
## [63] RColorBrewer_1.1-2 Rcpp_1.0.3
## [65] plyr_1.8.5 base64enc_0.1-3
## [67] progress_1.2.2 zlibbioc_1.28.0
## [69] purrr_0.3.3 RCurl_1.95-4.12
## [71] prettyunits_1.1.0 rpart_4.1-15
## [73] SummarizedExperiment_1.12.0 cluster_2.0.9
## [75] data.table_1.12.6 SparseM_1.77
## [77] ProtGenerics_1.14.0 matrixStats_0.55.0
## [79] hms_0.5.2 evaluate_0.13
## [81] xtable_1.8-4 XML_3.98-1.19
## [83] gcrma_2.54.0 gridExtra_2.3
## [85] RITools_0.1-17 compiler_3.5.1
## [87] biomaRt_2.38.0 KernSmooth_2.23-15
## [89] crayon_1.3.4 ReportingTools_2.22.1
## [91] R.oo_1.22.0 htmltools_0.3.6
## [93] GOstats_2.48.0 Formula_1.2-3
## [95] geneplotter_1.60.0 MASS_7.3-51.4
## [97] Matrix_1.2-17 cli_2.0.1
## [99] R.methodsS3_1.7.1 gdata_2.18.0
```

Case-Control with Gestational Age Matching

```
## [101] GenomicRanges_1.34.0      pkgconfig_2.0.3
## [103] GenomicAlignments_1.18.1    foreign_0.8-71
## [105] foreach_1.4.4               annotate_1.60.1
## [107] AnnotationForge_1.24.0      VariantAnnotation_1.28.13
## [109] digest_0.6.23               graph_1.60.0
## [111] rmarkdown_1.12              cellranger_1.1.0
## [113] htmlTable_1.13.2            edgeR_3.24.3
## [115] GSEABase_1.44.0             curl_4.3
## [117] Rsamtools_1.34.1            gtools_3.8.1
## [119] lifecycle_0.1.0             PFAM.db_3.7.0
## [121] fansi_0.4.1                 BSgenome_1.50.0
## [123] pillar_1.4.3                lattice_0.20-38
## [125] GGally_1.4.0                httr_1.4.0
## [127] survival_2.44-1.1           GO.db_3.7.0
## [129] glue_1.3.1                  iterators_1.0.10
## [131] bit_1.1-14                  Rgraphviz_2.26.0
## [133] stringi_1.4.5               blob_1.1.1
## [135] DESeq2_1.22.2               latticeExtra_0.6-28
## [137] caTools_1.17.1.2            memoise_1.1.0
```