

# GTEX Hippocampus DEG

Kennedi Todd

3/12/2021

Set working directory.

Load packages

```
##  
## Attaching package: 'dplyr'  
## The following objects are masked from 'package:stats':  
##  
##     filter, lag  
## The following objects are masked from 'package:base':  
##  
##     intersect, setdiff, setequal, union  
## Loading required package: limma  
##  
## Attaching package: 'gplots'  
## The following object is masked from 'package:stats':  
##  
##     lowess  
##  
## Attaching package: 'gridExtra'  
## The following object is masked from 'package:dplyr':  
##  
##     combine  
## Loading required package: foreach  
## Loading required package: scales  
## Loading required package: Biobase  
## Loading required package: BiocGenerics  
## Warning: package 'BiocGenerics' was built under R version 4.0.5  
## Loading required package: parallel  
##  
## Attaching package: 'BiocGenerics'  
## The following objects are masked from 'package:parallel':  
##  
##     clusterApply, clusterApplyLB, clusterCall, clusterEvalQ,  
##     clusterExport, clusterMap, parApply, parCapply, parLapply,  
##     parLapplyLB, parRapply, parSapply, parSapplyLB
```

```

## The following object is masked from 'package:gridExtra':
##
##      combine

## The following object is masked from 'package:limma':
##
##      plotMA

## The following objects are masked from 'package:dplyr':
##
##      combine, intersect, setdiff, union

## The following objects are masked from 'package:stats':
##
##      IQR, mad, sd, var, xtabs

## The following objects are masked from 'package:base':
##
##      anyDuplicated, append, as.data.frame, basename, cbind, colnames,
##      dirname, do.call, duplicated, eval, evalq, Filter, Find, get, grep,
##      grepl, intersect, is.unsorted, lapply, Map, mapply, match, mget,
##      order, paste, pmax, pmax.int, pmin, pmin.int, Position, rank,
##      rbind, Reduce, rownames, sapply, setdiff, sort, table, tapply,
##      union, unique, unsplit, which.max, which.min

## Welcome to Bioconductor
##
##      Vignettes contain introductory material; view with
##      'browseVignettes()'. To cite Bioconductor, see
##      'citation("Biobase")', and for packages 'citation("pkgname")'.

##
## Attaching package: 'variancePartition'

## The following object is masked from 'package:limma':
##
##      classifyTestsF

```

Read in files.

```

file1 <- read.delim("Brain_Hippocampus", header = TRUE, sep = "\t")
file2 <- read.delim("GTEx_Analysis_v8_Annotations_SubjectPhenotypesDS.txt",
                     header = TRUE, sep = "\t")
gtf.file <- "gencode.v26.GRCh38.genes.gtf"
gtf.gr <- rtracklayer::import(gtf.file)
gtf.df <- as.data.frame(gtf.gr)
gtf.df.condensed <-
  unique(gtf.df[,c("type","gene_type","gene_name","gene_id","seqnames")])

```

Reformat and rename files.

```

gtf.genes <- gtf.df.condensed[gtf.df.condensed$type == "gene",]

metadata <- file2
colnames(metadata) <- c("subject_id","sex","age","hardy_scale")

file1$X <- NULL
counts <- file1
colnames(counts)[1:2] <- c("gene_id","gene_name")

```

Merge to view info in one file.

```
counts.info <- left_join(counts, gtf.genes, by = "gene_id")
counts.info <- counts.info[,c(1,202:205,3:199)]
colnames(counts.info)[4] <- "gene_name"
```

Look to see if any NA values or replicated gene names.

```
table(is.na(counts.info$gene_name))

##
## FALSE
## 56200

table(duplicated(counts.info$gene_name))
```

```
##
## FALSE TRUE
## 54592 1608
```

There are no NA values. However, 1,608 gene names are repetitive.

View the repetitive gene names.

```
table(counts.info$gene_name[duplicated(counts.info$gene_name)])
```

	5_8S_rRNA	5S_rRNA	7SK	ACA59	ACA64
##	4	15	5	1	5
##	ACTR3BP2	AKAP17A	ALG1L9P	AMD1P2	ASMT
##	1	1	1	1	1
##	ASMTL	ASMTL-AS1	C2orf61	CD99	CD99P1
##	1	1	1	1	1
## Clostridiales-1		CRLF2	CSF2RA	CTSLP2	CYB561D2
##	2	1	1	1	1
##	DDX11L16	DHRSX	DHRSX-IT1	DNAJC9-AS1	DPH3P2
##	1	1	1	1	1
##	ELFN2	ELOCP24	FABP5P13	FAS-AS1	GOLGA8M
##	1	1	1	1	1
##	GTPBP6	IL3RA	IL9R	KRT18P53	LINC00102
##	1	1	1	1	1
##	LINC00106	LINC00484	LINC00685	LINC00901	LINC01115
##	1	1	1	1	1
##	LINC01238	LINC01297	LINC01347	LINC01422	LINC01481
##	1	1	1	1	1
##	LINC01598	LLOCYNC03-29C1.1	LYNX1	MAL2	Metazoa_SRP
##	1	1	1	1	166
##	MIR3179-3	MIR3180-4	MIR3690	MIR6089	NBPF13P
##	1	1	1	1	1
##	OR7E47P	P2RY8	PLCXD1	PMS2P6	PPP2R3B
##	1	1	1	1	1
##	pRNA	PROX1-AS1	RAET1E-AS1	RGS5	RNA5-8S5
##	8	1	1	1	3
##	RNA5SP498	RP11-309M23.1	RP13-297E16.4	RP13-297E16.5	RP13-465B17.4
##	1	1	1	1	1
##	RP13-465B17.5	RPL14P5	RPS23P5	SCARNA11	SCARNA15
##	1	1	1	3	1
##	SCARNA16	SCARNA17	SCARNA18	SCARNA20	SCARNA21

##		4		3		1		6		3
##	SCARNA24		SCARNA6		SHOX		SLC25A6		snoMBII-202	
##	1		1		1		1		1	
##	snoMe28S-Am2634		SNORA1		SNORA11		SNORA12		SNORA15	
##	2		4		6		2		3	
##	SNORA17		SNORA18		SNORA19		SNORA2		SNORA20	
##	1		7		3		2		4	
##	SNORA21		SNORA22		SNORA24		SNORA25		SNORA26	
##	1		3		1		18		10	
##	SNORA27		SNORA3		SNORA30		SNORA31		SNORA32	
##	4		4		2		25		2	
##	SNORA33		SNORA35		SNORA36		SNORA38		SNORA4	
##	3		2		1		1		4	
##	SNORA40		SNORA41		SNORA42		SNORA43		SNORA44	
##	20		1		3		4		1	
##	SNORA46		SNORA48		SNORA50		SNORA51		SNORA57	
##	1		12		1		12		2	
##	SNORA58		SNORA62		SNORA63		SNORA64		SNORA67	
##	2		6		10		5		6	
##	SNORA68		SNORA69		SNORA7		SNORA70		SNORA71	
##	5		1		4		29		2	
##	SNORA72		SNORA73		SNORA74		SNORA75		SNORA76	
##	6		3		6		7		1	
##	SNORA77		SNORA79		SNORA8		SNORA81		SNORA9	
##	2		1		4		3		3	
##	SNORD11		SNORD111		SNORD112		SNORD113		SNORD115	
##	1		1		2		4		1	
##	SNORD116		SNORD18		SNORD19		SNORD19B		SNORD2	
##	2		1		1		1		1	
##	SNORD22		SNORD23		SNORD27		SNORD28		SNORD30	
##	1		1		1		1		1	
##	SNORD33		SNORD37		SNORD38		SNORD39		SNORD41	
##	1		2		3		3		1	
##	SNORD42		SNORD45		SNORD46		SNORD5		SNORD51	
##	1		3		2		2		1	
##	SNORD56		SNORD59		SNORD60		SNORD63		SNORD65	
##	6		1		2		3		3	
##	SNORD66		SNORD67		SNORD70		SNORD74		SNORD75	
##	1		1		1		6		1	
##	SNORD77		SNORD78		SNORD81		snoU109		snoU13	
##	4		1		2		5		31	
##	snoU2_19		snoZ6		SPATA13		SPRY3		TRPC6P	
##	2		3		1		1		1	
##	U1		U2		U3		U4		U6	
##	12		18		49		6		32	
##	U7		U8		uc_338		VAMP7		Vault	
##	5		21		31		1		1	
##	WASH6P		WASIR1		Y_RNA		ZBED1			
##	1		1		737		1			

View the gene\_type for the repetitive genes. This will help determine if we should remove repetitive gene names.

```



```

##                                antisense          lincRNA
##                               12                  15
##                                miRNA          misc_RNA
##                               3                  949
##      processed_pseudogene processed_transcript
##                               9                  1
##      protein_coding           rRNA
##                               21                 23
##                                scaRNA       sense_intronic
##                               22                  2
##                                snoRNA          snRNA
##                               468                 73
##      sRNA transcribed_unprocessed_pseudogene
##                               2                  3
##      unprocessed_pseudogene
##                               5

```


```

---

—CREATE DGE OBJECT—

Remove duplicated names

```

keep <- !duplicated(counts.info$gene_name) # true when no duplicate
counts.info.unique <- counts.info[keep,]
dim(counts.info)
```

```
## [1] 56200   202
```

```
dim(counts.info.unique)
```

```
## [1] 54592   202
```

Went from 56,200 genes to 54,592.

RNAs were mainly lost.

Remove MT genes.

```

keep <- counts.info.unique$seqnames != "chrM"
counts.info.noMT <- counts.info.unique[keep,]
```

Make gene\_name the row name.

```
rownames(counts.info.unique.noMT) <- counts.info.unique.noMT$gene_name
```

Change sex column from numbers to letters.

```

# male = 1, female = 2
sex_nums <- metadata$sex
sex_letters <- vector()
for (i in 1:length(sex_nums)) {
  if (sex_nums[i] == 1) {
    sex_letters <- c(sex_letters, "M")
  }
  else if (sex_nums[i] == 2) {
    sex_letters <- c(sex_letters, "F")
  }
}
```

```
}
```

```
head(cbind(sex_nums, sex_letters))
```

```
##      sex_nums sex_letters
## [1,] "2"       "F"
## [2,] "1"       "M"
## [3,] "1"       "M"
## [4,] "1"       "M"
## [5,] "1"       "M"
## [6,] "2"       "F"
```

```
metadata$sex <- sex_letters
metadata$subject_id <- make.names(metadata$subject_id)
```

Subset counts data.

```
counts <- counts.info.unique.noMT[, 6:202]
```

Make counts columns same order as metadata rows.

```
metadata_all_samples <- metadata$subject_id
counts_hip_samples <- colnames(counts)

new_names_counts <- vector()
new_names_metadata <- vector()

for (i in 1:length(counts_hip_samples)) {
  for (j in 1:length(metadata_all_samples)) {
    checkpoint <- str_detect(counts_hip_samples[i], metadata_all_samples[j])
    if (checkpoint) {
      new_names_counts <- c(new_names_counts, counts_hip_samples[i])
      new_names_metadata <- c(new_names_metadata, metadata_all_samples[j])
    }
  }
}

head(cbind(new_names_counts, new_names_metadata)) #check
```

```
##      new_names_counts      new_names_metadata
## [1,] "GTEX.11DXW.0011.R1a.SM.DNZZD" "GTEX.11DXW"
## [2,] "GTEX.11EI6.0011.R1a.SM.D093L" "GTEX.11EI6"
## [3,] "GTEX.11GS4.0011.R1a.SM.D0129" "GTEX.11GS4"
## [4,] "GTEX.11GSO.0011.R1b.SM.57WD3" "GTEX.11GSO"
## [5,] "GTEX.11GSP.0011.R1a.SM.9QEJ3" "GTEX.11GSP"
## [6,] "GTEX.11ONC.0011.R1a.SM.57WD4" "GTEX.11ONC"
```

Re-order dataframes.

```
counts <- counts[, new_names_counts]
tempdf <- as.data.frame(new_names_metadata)
colnames(tempdf) <- "subject_id"
metadata <- left_join(x = tempdf, y = metadata, by = "subject_id")

all.equal(colnames(counts), new_names_counts)

## [1] TRUE
```

```
metadata$counts_id <- colnames(counts)
```

Create dge object

```
dge <- DGEList(counts,
                 samples = metadata,
                 genes = counts.info.unique.noMT[,1:5])
```

Reformat type of sample columns.

```
dge$samples[,c(5:7)] <-
  lapply(dge$samples[,c(5:7)], factor)
```

---

#### TTR TPM HISTOGRAM

---

Data is already in TPM (transcripts per million). Create histogram of TPM TTR values.

```
ttr.counts <- as.data.frame(dge$counts[["TTR",]])
colnames(ttr.counts) <- "ttr.counts"

log2.ttr.counts <- log2(ttr.counts + 0.01)
colnames(log2.ttr.counts) <- "log2.ttr.counts"

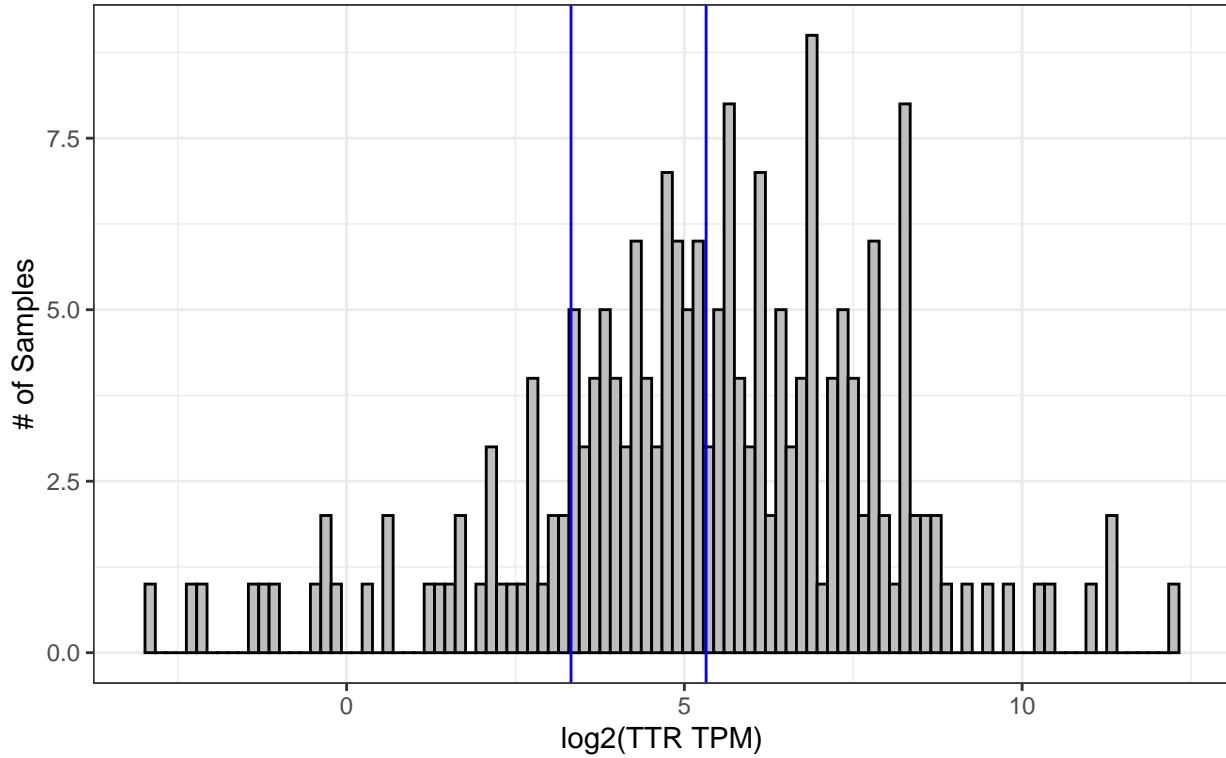
h1 <- ggplot(ttr.counts, aes(x = ttr.counts)) +
  geom_histogram(bins = 100, fill = "gray", color = "black") +
  labs(title = "A.", x=NULL, y=NULL) +
  xlab("TTR TPM") + ylab("# of Samples") +
  geom_vline(xintercept = 20, col = "red") +
  theme_bw()

h2 <- ggplot(ttr.counts, aes(x = ttr.counts)) +
  geom_histogram(bins = 100, fill = "gray", color = "black") +
  labs(title = "B.", x=NULL, y=NULL) +
  xlab("TTR TPM") + ylab("# of Samples") +
  xlim(0, 500) + geom_vline(xintercept = 20, col = "red") +
  theme_bw()

h3 <- ggplot(log2.ttr.counts, aes(x = log2.ttr.counts)) +
  geom_histogram(bins = 100, fill = "gray", color = "black") +
  labs(title = "GTEx\nHippocampus", x=NULL, y=NULL) +
  xlab("log2(TTR TPM)") + ylab("# of Samples") +
  theme_bw() +
  geom_vline(xintercept = log2(20) + 1, col = "blue") +
  geom_vline(xintercept = log2(20) - 1, col = "blue")
```

```
h3
```

GTEX  
Hippocampus



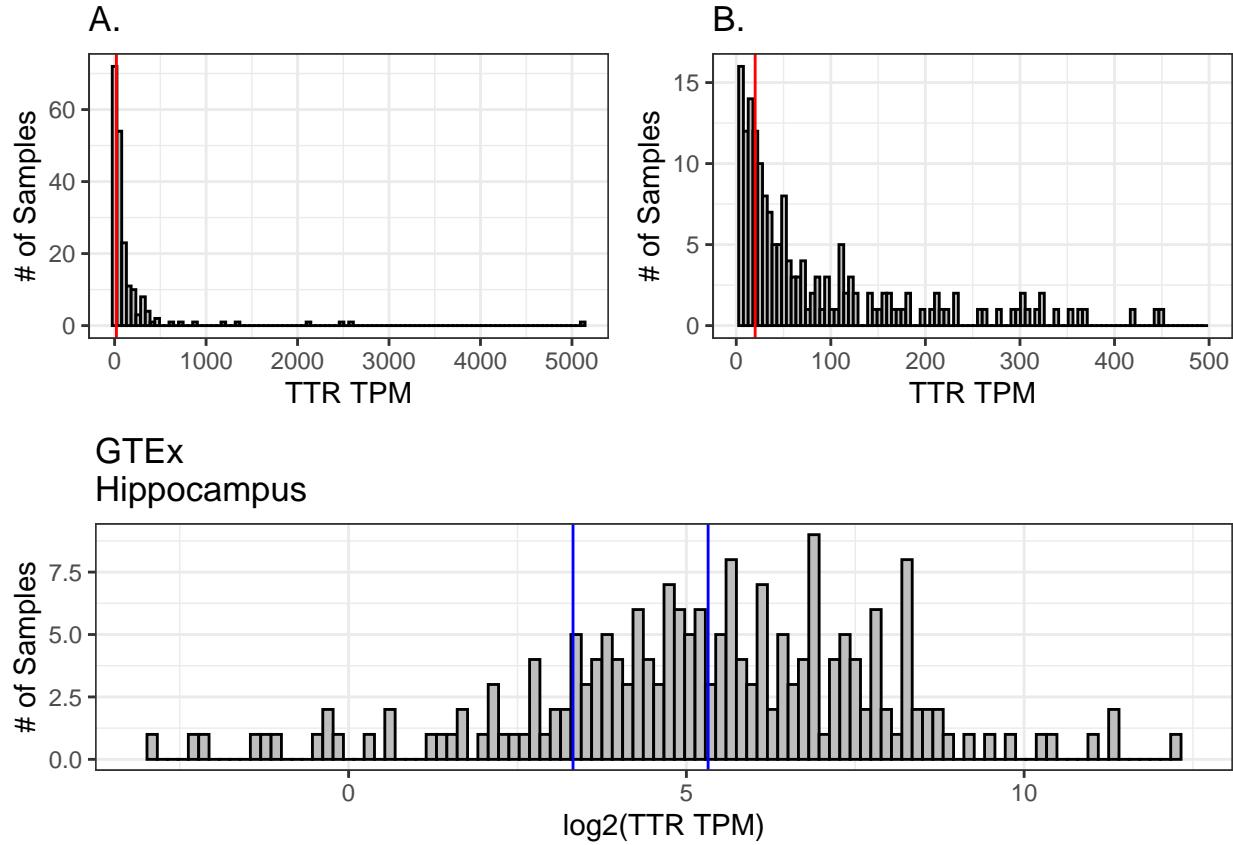
Arrange graphs in grid.

```
plots1 <- list(h1,h2,h3)

layout1 <- rbind(c(1,2),c(3))

grid1 <- grid.arrange(grobs = plots1, layout_matrix = layout1)

## Warning: Removed 9 rows containing non-finite values (stat_bin).
## Warning: Removed 2 rows containing missing values (geom_bar).
```



Lets looks at groups +/- 1 of the cutoff

```
cutoff <- log2(20)
clean_threshold <- cutoff - 1 # values <= this are clean
clean_threshold

## [1] 3.321928
cutoff

## [1] 4.321928
dirty_threshold <- cutoff + 1 # values >= this are dirty
dirty_threshold

## [1] 5.321928

Assign clean or dirty

clean_or_dirty <- vector()
vector.log2.ttr.counts <- as.vector(log2.ttr.counts$log2.ttr.counts)

for (i in 1:length(vector.log2.ttr.counts)){
  if (vector.log2.ttr.counts[i] <= clean_threshold){
    clean_or_dirty <- c(clean_or_dirty, "clean")
  }
  else if (vector.log2.ttr.counts[i] >= dirty_threshold){
    clean_or_dirty <- c(clean_or_dirty, "dirty")
  }
}
```

```

    else {
      clean_or_dirty <- c(clean_or_dirty, "neither")
    }
  }
table(clean_or_dirty)

## clean_or_dirty
##   clean   dirty neither
##     34     100     63

dge$samples$group <- factor(clean_or_dirty)
head(dge$samples)

##                                     group lib.size norm.factors subject_id sex   age
## GTEX.11DXW.0011.R1a.SM.DNZZD neither 295069.0                 1 GTEX.11DXW   M 40-49
## GTEX.11EI6.0011.R1a.SM.D093L   dirty 255731.1                 1 GTEX.11EI6   M 60-69
## GTEX.11GS4.0011.R1a.SM.D0129   dirty 286639.2                 1 GTEX.11GS4   M 60-69
## GTEX.11GSO.0011.R1b.SM.57WD3 neither 290899.0                 1 GTEX.11GSO   M 60-69
## GTEX.11GSP.0011.R1a.SM.9QEJ3   clean 288836.8                 1 GTEX.11GSP   F 60-69
## GTEX.11ONC.0011.R1a.SM.57WD4   clean 214828.4                 1 GTEX.11ONC   M 60-69
##                                     hardy_scale counts_id
## GTEX.11DXW.0011.R1a.SM.DNZZD           2 GTEX.11DXW.0011.R1a.SM.DNZZD
## GTEX.11EI6.0011.R1a.SM.D093L           4 GTEX.11EI6.0011.R1a.SM.D093L
## GTEX.11GS4.0011.R1a.SM.D0129           2 GTEX.11GS4.0011.R1a.SM.D0129
## GTEX.11GSO.0011.R1b.SM.57WD3           2 GTEX.11GSO.0011.R1b.SM.57WD3
## GTEX.11GSP.0011.R1a.SM.9QEJ3           2 GTEX.11GSP.0011.R1a.SM.9QEJ3
## GTEX.11ONC.0011.R1a.SM.57WD4           2 GTEX.11ONC.0011.R1a.SM.57WD4

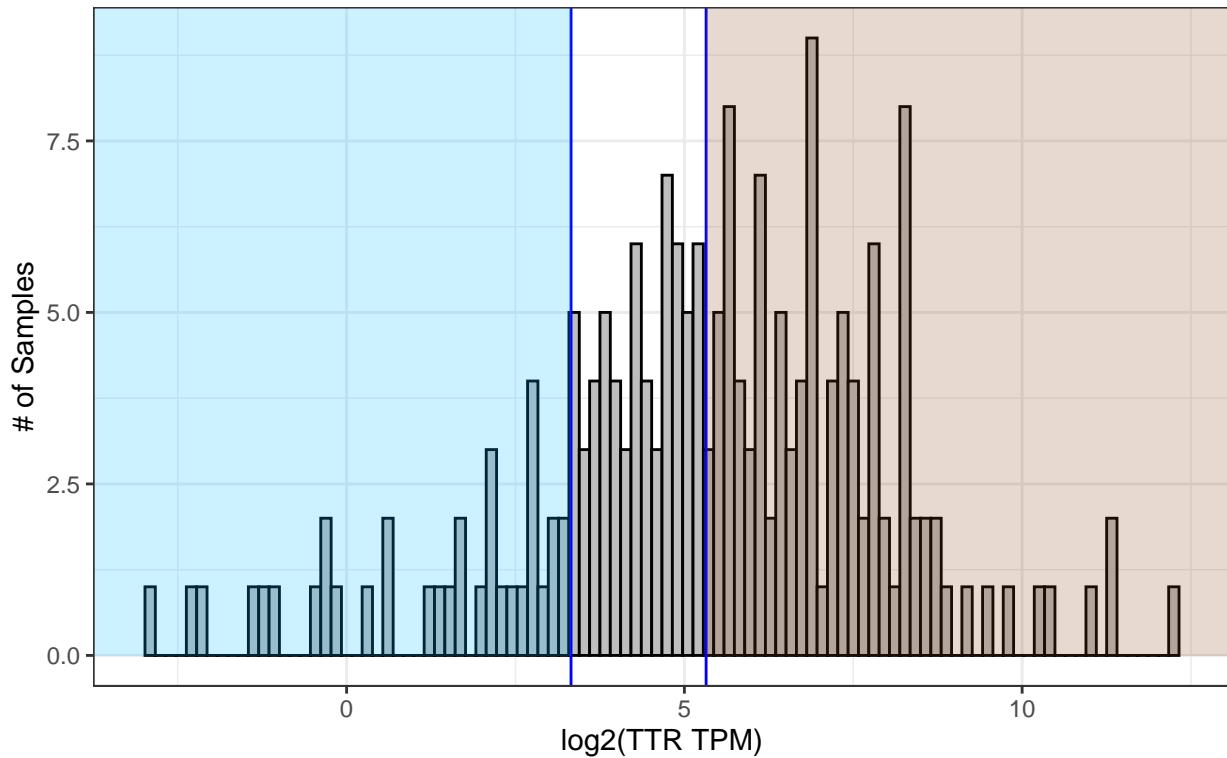
h4 <- h3 +
  annotate("rect",
    xmin = -Inf,
    xmax = clean_threshold,
    ymin = 0,
    ymax=Inf,
    alpha=0.2,
    fill="deepskyblue") +
  annotate("rect",
    xmin = dirty_threshold,
    xmax = Inf,
    ymin = 0,
    ymax=Inf,
    alpha=0.2,
    fill="chocolate4")

h5 <- h1 +
  annotate("rect",
    xmin = -Inf,
    xmax = clean_threshold,
    ymin = 0,
    ymax=Inf,
    alpha=0.2,
    fill="deepskyblue") +
  annotate("rect",
    xmin = dirty_threshold,
    xmax = Inf,
    ymin = 0,

```

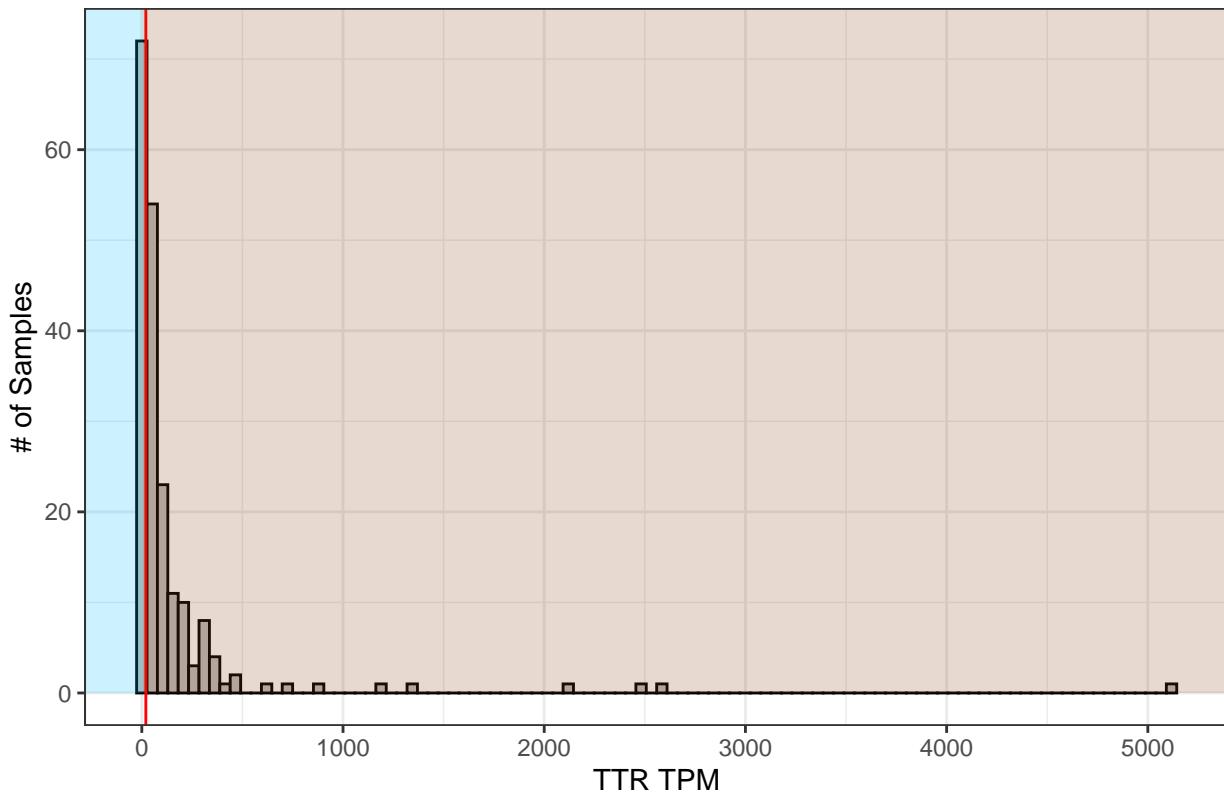
```
ymax=Inf,  
alpha=0.2,  
fill="chocolate4")  
h4
```

GTEEx  
Hippocampus



```
h5
```

A.



---

FILTER LOWLY EXPRESSED GENES

---

TPM to CPM

```
cpm <- cpm(dge)
```

Filter

```
# produces a logical matrix
thresh <- cpm > 2 # cpm greater than 2

# keep genes that have at least 3 TRUES in each row of thresh
keep <- rowSums(thresh) >= 3
summary(keep)
```

```
##      Mode    FALSE     TRUE
## logical   30225   24330

# redefine object, filter by logical
dge.filtered <- dge[keep, keep.lib.sizes = FALSE]

dim(dge)
```

```
## [1] 54555    197
```

```
dim(dge.filtered)
```

```
## [1] 24330    197
```

Went from 54,592 genes to 18,883 genes.

---

### TRIMMED MEAN OF M-VALUES

---

Now, we want to normalize gene expression distributions. We do this by normalizing library size differences in each sample. The method Trimmed Mean of M-values (TMM) is used with the calcNormFactors in the edgeR package. The normalization factors calculated here are used as a scaling factor for the library size.

Trimmed mean of M values (TMM) normalization estimates sequencing depth after excluding genes for which the ratio of counts between a pair of experiments is too extreme or for which the average expression is too extreme. The edgeR software implements a TMM normalization (Biostar Handbook 2020).

```
# creates another large DGEList
dge.filtered.norm <- calcNormFactors(dge.filtered, method = "TMM")

# view norm factors
dge.filtered.norm$samples$norm.factors

## [1] 0.9208014 0.9397341 1.0664196 1.1657822 1.1034621 1.0757837 0.9692125
## [8] 1.0345089 1.1286779 0.9074237 1.0777276 1.1876111 1.0649142 0.9585218
## [15] 0.9699986 0.9005879 0.9866129 0.9987741 0.9415839 0.9032043 1.0152109
## [22] 1.0121925 1.0371544 0.8992394 1.0453845 0.9016306 1.0659189 1.1260111
## [29] 0.9450652 0.9096229 0.8623565 0.8448163 1.1555415 0.9492569 0.9825655
## [36] 1.0605614 1.0104486 0.8361482 0.8753509 0.9654817 1.0096480 0.8608656
## [43] 0.8320688 0.9922814 1.1346967 0.9275243 1.0780612 1.0901378 1.0437759
## [50] 0.8624394 1.0029232 1.0540529 1.1874402 0.8926142 1.0670313 1.0547220
## [57] 1.0384668 0.7318076 0.9500912 1.1119956 0.8796994 0.9476418 0.9584379
## [64] 1.0863173 1.0988247 1.0302707 1.0276728 1.1872606 0.9843724 0.7001994
## [71] 1.0412301 1.1075697 1.1262024 0.8636979 1.0865561 1.1177207 1.0470575
## [78] 1.0812113 1.0651687 1.0803114 0.8892324 1.0766973 0.9631251 1.0482810
## [85] 1.1339810 1.0978243 0.9991410 0.9488282 1.0089296 0.9492502 1.0222190
## [92] 0.9941390 1.0646432 0.9794406 0.9119222 0.9667743 0.9246916 0.9639162
## [99] 1.0380671 0.9249168 0.9808264 1.0000709 1.1134358 0.9578064 0.9663822
## [106] 0.9866986 0.6440762 0.8832843 1.0252993 0.9320145 1.1172923 0.9903614
## [113] 1.0588432 0.9264131 1.0554269 0.9254882 0.9490689 0.8831684 1.1087182
## [120] 0.8840364 1.0230063 1.0438604 0.9622401 0.8130978 1.0115068 1.0412311
## [127] 0.9575043 0.9875283 0.9423856 1.0974353 0.9971180 1.0841743 1.0317499
## [134] 0.9764619 1.0516356 1.0403124 0.8842430 0.8214545 1.1038727 0.9825239
## [141] 0.8676495 1.0425724 0.8893092 0.8093808 1.1624619 1.0098526 1.0763740
## [148] 1.1643346 0.9752810 1.0302402 1.0366424 0.8129844 1.0329518 1.0349895
## [155] 0.9732349 1.0351201 0.9280301 0.9094309 1.0713761 0.9071833 1.1432864
## [162] 1.1622661 1.0131992 0.9996608 1.0791143 1.0572388 1.0152764 0.9515026
## [169] 1.1424870 0.9718733 1.1850443 1.2132065 1.0591215 0.9334512 1.1576479
## [176] 1.0268450 1.0141231 1.0059438 1.1929826 0.9018986 0.9574237 1.0124440
## [183] 1.0337595 1.0564354 1.0587401 1.1084535 1.0184384 1.0201374 1.0178207
## [190] 0.9576280 1.0526098 1.0964601 1.0651772 0.9379478 0.9691619 1.1172171
## [197] 1.0581417

summary(dge.filtered.norm$samples$norm.factors)

##      Min.    1st Qu.     Median      Mean    3rd Qu.      Max.
## 0.6441  0.9488  1.0124  1.0047  1.0652  1.2132
```

Looking to be around ~1.

---

### Variance Partition

---

See what options we have to choose from.

```

colnames(dge.filtered.norm$samples)

## [1] "group"         "lib.size"        "norm.factors"   "subject_id"    "sex"
## [6] "age"           "hardy_scale"     "counts_id"

Specify variables to consider

#register(SnowParam(4)) # work in parallel

# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
#geneExpr <- as.matrix(dge.filtered.norm$counts)
#info <- as.data.frame(dge.filtered.norm$samples)
#info$hardy_scale <- as.factor(info$hardy_scale)

# Age is USUALLY continuous so model it as a fixed effect "age"
# However, this metadata categorizes it into bins so it's categorical
# Sex is categorical, so model them as random effects "(1/sex)"
# Note the syntax
#form <- ~ (1/group) +
# (1/age) +
# (1/sex) +
# (1/hardy_scale)

#varPart <- fitExtractVarPartModel(geneExpr, form, info)
#vp <- sortCols(varPart)

```

Violin plots

```
#plotVarPart(vp)
```

Percent bars

```
#plotPercentBars(vp[1:10,])
```

```
# sort genes based on variance explained by group
#varPart.df <- as.data.frame(varPart)
#order.varPart.df <- varPart.df[order(varPart.df$group, #decreasing = TRUE),]
#order.varPart.df["group"]
```

---

### DESIGN MATRIX

---

```
sex <- dge.filtered.norm$samples$sex
age <- dge.filtered.norm$samples$age
```

Create a design matrix and specify the samples we want to compare to each other.

```
# interaction computes a factor which represents the interaction of the given factors
group <- interaction(dge.filtered$samples$group)

design <- model.matrix(~ 0 + group + sex + age)
colnames(design) <- make.names(colnames(design))
colnames(design)[1:2] <- c("clean","dirty")
```

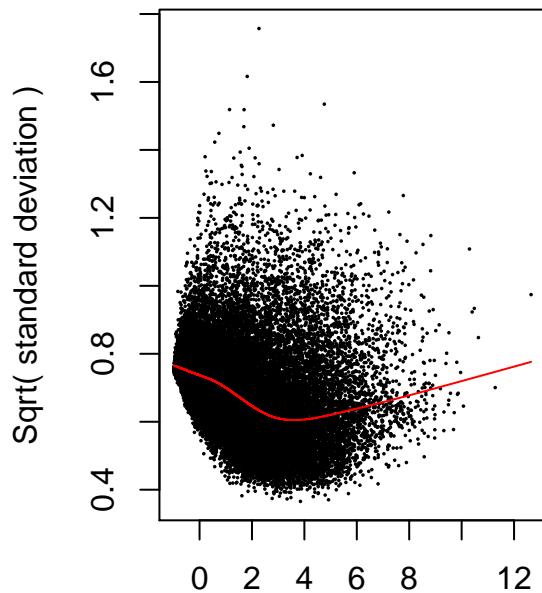
Run voom with quality weights. This combines observational-level weights with sample-specific quality weights in an experiment. Normalizes expression intensities so that the log-ratios have similar distributions across a set of samples. To quantile normalize, add normalize.method = “quantile”.

Quality weights in RNA-seq data improves data quality when samples with high variance are present. This

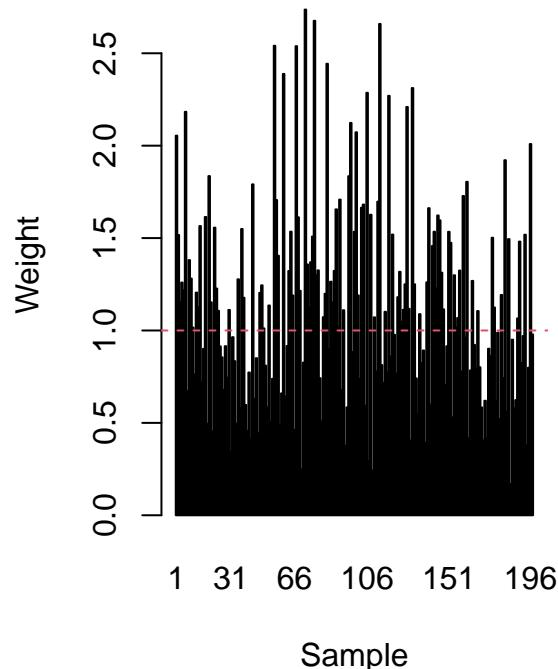
variation is modeled by taking into account global intensity-dependent trends using voom. Sample-specific variability is taken into account with a log-linear model that shares parameters between genes. <https://doi.org/10.1093/nar/gkv412>

```
v <- voomWithQualityWeights(dge.filtered.norm, design, plot = TRUE)
```

### voom: Mean–variance trend



### Sample-specific weights



```
group_colors <- c("deepskyblue", "chocolate4", "gray")[v$targets$group] # blue is clean, brown is dirty
v$targets$sex <- as.factor(v$targets$sex)
point_shapes <- c("F", "M")[v$targets$sex] # will say F or M for points
head(data.frame(v$targets$group, group_colors, v$targets$sex, point_shapes)) # Check everything matches

##   v.targets.group group_colors v.targets.sex point_shapes
## 1      neither      gray          M           M
## 2       dirty    chocolate4          M           M
## 3       dirty    chocolate4          M           M
## 4      neither      gray          M           M
## 5      clean    deepskyblue         F           F
## 6      clean    deepskyblue         M           M

#create a png file# png(filename = "mds_1and2.png", width = 400, height = 400)
#default height and width are 480, default units are "px" (pixels)

mds <- plotMDS(
  v, # our data object
  top = 10, # only looking at top 10 genes
  pch = point_shapes,
  cex = 1, # point size
  dim.plot = c(5,6), # specifying principal components to be plotted
  plot = TRUE, # if TRUE then plot is created on current graphics device
```

```

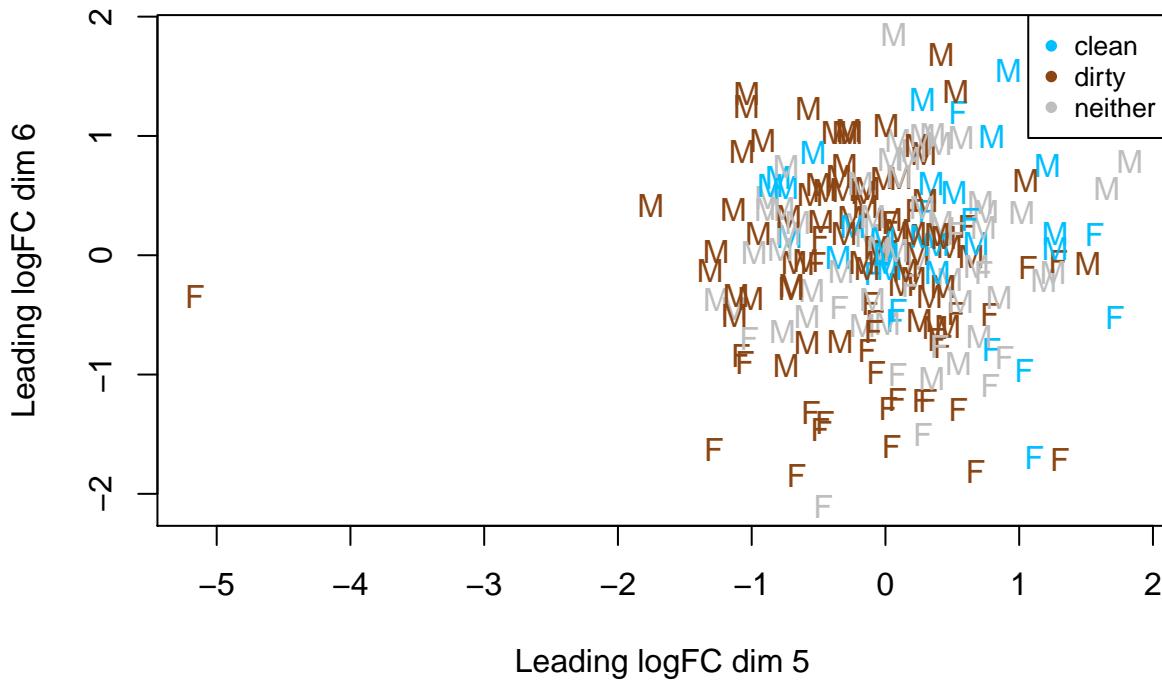
    col = group_colors # assigning our colors for each genotype
)

legend(
  "topright",
  pch = 16,
  legend = c("clean","dirty","neither"),
  col = c("deepskyblue","chocolate4","gray"),
  cex = 0.8
)

title("Top 10 Genes")

```

## Top 10 Genes




---

### -HIERARCHIAL HEAT CLUSTERING-

```

ds[order(factor(Month, levels = c("Apr", "May", "Jan", "Feb", "Mar"))), .SD, ID]
current_order <- dge.filtered.norm$samples$counts_id

# reorder levels
group_order <- dge.filtered.norm$samples
group_order <- group_order[order(factor(group, levels = c("dirty", "neither", "clean"))),]
sample_group_order <- group_order$counts_id
dge.group <- dge.filtered.norm[, sample_group_order]

# sample group colors
sample_group_color <- c("deepskyblue", "chocolate4", "gray")[dge.group$samples$group]
as.data.frame(dge.group$samples$group, sample_group_color)

##          dge.group$samples$group
## chocolate4           dirty

```









```

#      Colv = sample_group_order,
#      colCol = sample_group_color,
#      srtCol = 45,
#      xlab = "Samples"
#      )

```

---

## LINEAR MODELS

---

Fit linear models for comparisons of interest. Linear modeling in limma is carried out using the lmFit and contrasts.fit functions.

What do lm.fit and contrasts.fit do? 1) Fit a separate model to the expression values for each gene. 2) An empirical Bayes moderation is carried out by borrowing info across all the gene to obtain more precise estimates of gene-wise variability. 3) The model's residual variances are plotted against average expression values in the next figure.

It can be seen from this plot that the variance is no longer dependent on the mean expression level.

```

# Fits linear model for each gene given a series of arrays
fit <- lmFit(
  v, # object containing log-expression values for a series of arrays
  design # design matrix
)

```

Contrast design for differential expression.

```

contrasts <-
  makeContrasts(
    dirty_vs_clean = dirty - clean,
    levels = colnames(design)
  )

head(contrasts)

##           Contrasts
## Levels       dirty_vs_clean
##   clean          -1
##   dirty           1
##   groupneither   0
##   sexM            0
##   age30.39        0
##   age40.49        0

```

Assign allComparisons to a vector for later use.

```

allComparisons <- colnames(contrasts)
allComparisons # check

```

```

## [1] "dirty_vs_clean"

```

Run contrast analysis. Given a linear model fit, compute estimated coefficients and standard errors for a given set of contrasts.

```

vfit <- contrasts.fit(fit, contrasts = contrasts)

```

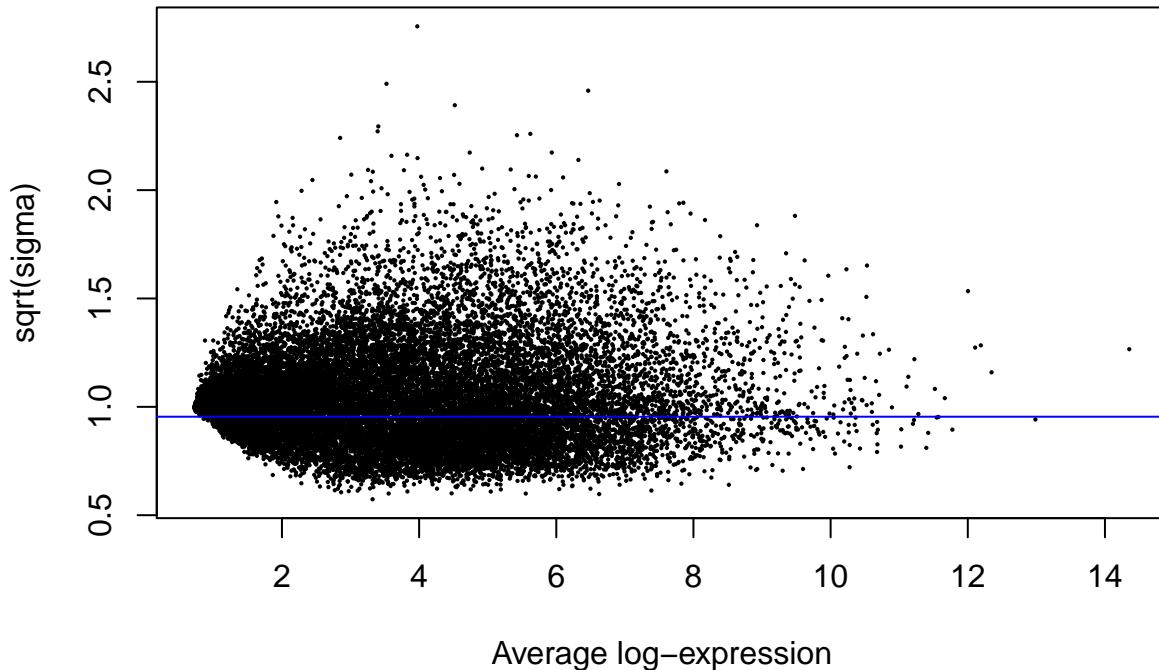
Compute differential expression based on the empirical Bayes moderation of the standard errors towards a common value.

```

veBayesFit <- eBayes(vfit)
plotSA(veBayesFit, main = "Final Model: Mean-variance Trend")

```

## Final Model: Mean–variance Trend



Looking at N of differentially expressed genes (DEGs) with adjusted  $p < 0.01$  and  $\log_{2}FC > 2$ .

`decideTests()` identifies which genes are significantly differentially expressed for each contrast from a fit object containing p-values and test statistics.

```
sumTable <-  
  summary(decideTests(  
    vfit, # object  
    # by default the method = "separate"  
    adjust.method = "BH",  
    p.value = 0.05,  
    lfc = 1 # numeric, minimum absolute log2-fold change required  
  ))  
  
head(sumTable)  
  
##          dirty_vs_clean  
## Down              25  
## NotSig            24236  
## Up                69
```

### DEG TABLES

For each comparison... extract table of the top-ranked genes (DEGs) from a linear model fit AND output a table.

NOTE: Log2FC of 1 is equivalent to linear fold change of 2.

```
coef = 1  
for (i in allComparisons) {  
  
  # p < 1, log2fc > 0 -----  
  vTopTableAll <-
```

```

topTable(
  veBayesFit,
  coef = coef,
  n = Inf,
  p.value = 1,
  lfc = 0
)
#output as txt file
path <- paste("DEG_", i, "_FDRq1_Log2FC0.txt", sep = "")
write.table(
  vTopTableAll,
  path,
  sep = "\t",
  row.names = TRUE,
  quote = FALSE
)

# p < 0.05, log2fc > 1 -----
vTopTable1 <-
  topTable(
    veBayesFit,
    coef = coef,
    n = Inf,
    p.value = 0.05,
    lfc = 1
  )
path <- paste("DEG_", i, "_FDRq0.05_Log2FC1.txt", sep = "")
write.table(
  vTopTable1,
  path,
  sep = "\t",
  row.names = TRUE,
  quote = FALSE
)

# increment the coefficient -----
coef <- coef + 1
}

```

### —VOLCANO PLOTS—

Using the DEG tables we made, graph a volcano plot. Read in file.

```

dirty_vs_clean <-
  read.table(
    "DEG_dirty_vs_clean_FDRq1_Log2FC0.txt",
    header = TRUE,
    sep = "\t",
    stringsAsFactors = F
  )

dirty_vs_clean$gene_symbol <- rownames(dirty_vs_clean)

summary(decideTests(
  vfit, # object

```

```

    adjust.method = "BH", # by default the method = "separate"
    p.value = 0.05,
    lfc = 1
  ))}

##          dirty_vs_clean
## Down           25
## NotSig        24236
## Up            69

color_values <- vector()
max <- nrow(dirty_vs_clean)

for(i in 1:max){
  if (dirty_vs_clean$adj.P.Val[i] < 0.05){
    if (dirty_vs_clean$logFC[i] > 1){
      color_values <- c(color_values, 1) # 1 when logFC > 1 and pval < 0.05
    }
    else if (dirty_vs_clean$logFC[i] < -1){
      color_values <- c(color_values, 2) # 2 when logFC < -1 and pval < 0.05
    }
    else{
      color_values <- c(color_values, 3) # 3 when -1 <= logFC <= 1 pval < 0.05
    }
  }
  else{
    color_values <- c(color_values, 3) # 3 when pval >= 0.05
  }
}

dirty_vs_clean$color_p0.05_lfc1 <- factor(color_values)

hadjpval <- (-log10(max(dirty_vs_clean$adj.P.Val[dirty_vs_clean$adj.P.Val < 0.05], na.rm=TRUE)))
hadjpval

## [1] 1.302149

Subset 10 lowest FDRq for up and down regulated genes.

keep <- dirty_vs_clean$color_p0.05_lfc1 == 1
up10 <- dirty_vs_clean[keep,]
up10 <- up10[order(up10$adj.P.Val),]
up10 <- up10[c(1:10),c(4,6,10)]

keep <- dirty_vs_clean$color_p0.05_lfc1 == 2
down10 <- dirty_vs_clean[keep,]
down10 <- down10[order(down10$adj.P.Val),]
down10 <- down10[1:10, c(4,6,10)]

p <-
  ggplot(data = dirty_vs_clean,
         aes(x = logFC, # x-axis is logFC
              y = -log10(adj.P.Val), # y-axis will be log10 of adj.P.Val
              color = color_p0.05_lfc1)) + # color is based on factored Color column
  geom_point(alpha = 0.8, size = 3) + # create scatterplot, alpha makes points transparent
  theme_bw() +

```

```

theme(legend.position = "none") + # no legend
xlim(c(-3, 6)) + #ylim(c(0, 25)) + # x and y axis limits
scale_color_manual(values = c("red", "blue", "grey")) + # set factor colors
labs(
  title = "", # no main title
  x = expression(log[2](FC)), # x-axis title
  y = expression(-log[10] ~ "(FDR adjusted " ~ italic("p") ~ "-value)") # y-axis title
) +
theme(axis.title.x = element_text(size = 10),
      axis.text.x = element_text(size = 10)) +
theme(axis.title.y = element_text(size = 10),
      axis.text.y = element_text(size = 10)) +
geom_hline(yintercept = hadjpval,
            colour = "#000000",
            linetype = "dashed") +
geom_vline(xintercept = 1,
            colour = "#000000",
            linetype = "dashed") +
geom_vline(xintercept = -1,
            colour = "#000000",
            linetype = "dashed") +
ggtitle("GTEEx\nFDRq < 0.05, -1 > LogFC > 1") +
geom_text_repel(data = up10,
                aes(x = logFC, y = -log10(adj.P.Val), label = gene_name),
                color = "black",
                fontface="italic",
                max.overlaps =getOption("ggrepel.max.overlaps", default = 80)
                ) +
geom_text_repel(data = down10,
                aes(x = logFC, y = -log10(adj.P.Val), label = gene_name),
                color = "black",
                fontface = "italic",
                max.overlaps =getOption("ggrepel.max.overlaps", default = 80)
)

```

p

GTEx  
FDRq < 0.05, -1 > LogFC > 1

