# 2021.7.16

using GEO data to explore DEGs between autistic spectrum disorder patients (neural diversity people) and normal people (neural typical people)

## GSE42133

| data source | [GSE42133](GSE42133) |
|---|---|
| title | Disrupted functional neworks in autism underlie early brain maldevelopment and provide accurate classification |
| Organism | Homo sapiens |
| Experiment type | Expression profiling by array |
| Status | Public on Mar 24, 2015 |

use GEO2R to analyze

```
# Version info: R 3.2.3, Biobase 2.30.0, GEOquery 2.40.0, limma 3.26.8
################################################################
#   Differential expression analysis with limma
library(GEOquery)
library(limma)
library(umap)

# load series and platform data from GEO

gset <- getGEO("GSE42133", GSEMatrix =TRUE, AnnotGPL=TRUE)
if (length(gset) > 1) idx <- grep("GPL10558", attr(gset, "names")) else idx <-
1
gset <- gset[[idx]]

# make proper column names to match toptable
fvarLabels(gset) <- make.names(fvarLabels(gset))

# group membership for all samples
gsms <- paste0("111100101011110001001101110000001010000000001100000",
        "11010011100101010110001011111110111111011111111001",
        "11001011111011110111111111110111111101110011111110")
sml <- strsplit(gsms, split="")[[1]]
```

```r
# log2 transformation
ex <- exprs(gset)
qx <- as.numeric(quantile(ex, c(0., 0.25, 0.5, 0.75, 0.99, 1.0), na.rm=T))
LogC <- (qx[5] > 100) ||
          (qx[6]-qx[1] > 50 && qx[2] > 0)
if (LogC) { ex[which(ex <= 0)] <- NaN
  exprs(gset) <- log2(ex) }

# assign samples to groups and set up design matrix
gs <- factor(sml)
groups <- make.names(c("control","ASD"))
levels(gs) <- groups
gset$group <- gs
design <- model.matrix(~group + 0, gset)
colnames(design) <- levels(gs)

fit <- lmFit(gset, design)  # fit linear model

# set up contrasts of interest and recalculate model coefficients
cts <- paste(groups[1], groups[2], sep="-")
cont.matrix <- makeContrasts(contrasts=cts, levels=design)
fit2 <- contrasts.fit(fit, cont.matrix)

# compute statistics and table of top significant genes
fit2 <- eBayes(fit2, 0.01)
tT <- topTable(fit2, adjust="fdr", sort.by="B", number=250)

tT <- subset(tT,
select=c("ID","adj.P.Val","P.Value","t","B","logFC","Gene.symbol","Gene.title")
)
write.table(tT, file=stdout(), row.names=F, sep="\t")
```
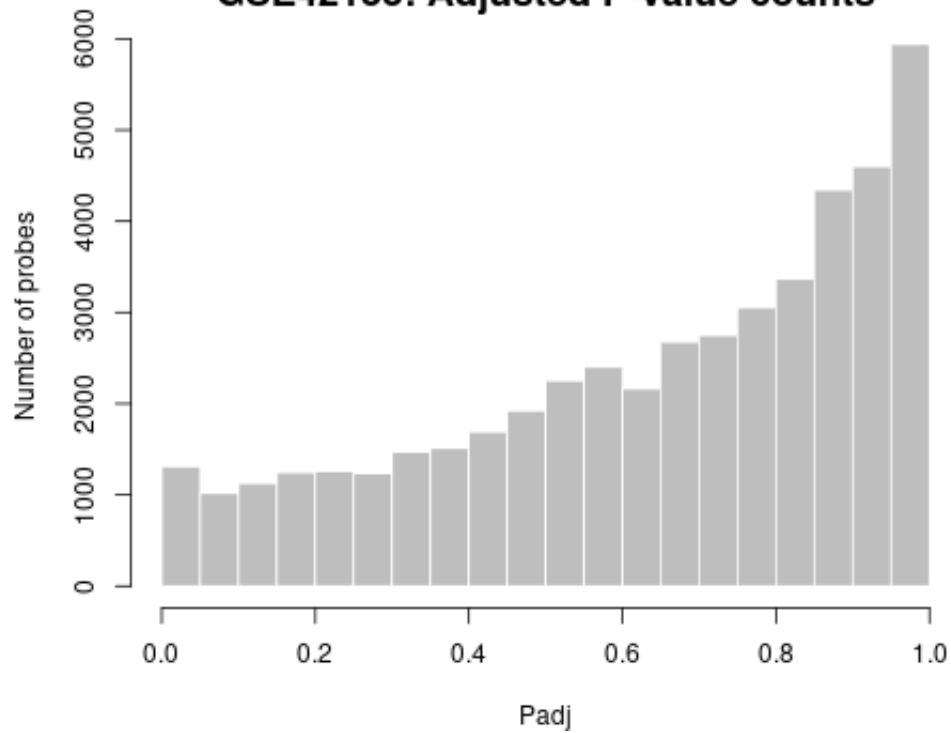
```r
# Visualize and quality control test results.
# Build histogram of P-values for all genes. Normal test
# assumption is that most genes are not differentially expressed.
tT2 <- topTable(fit2, adjust="fdr", sort.by="B", number=Inf)
hist(tT2$adj.P.Val, col = "grey", border = "white", xlab = "P-adj",
  ylab = "Number of genes", main = "P-adj value distribution")
```
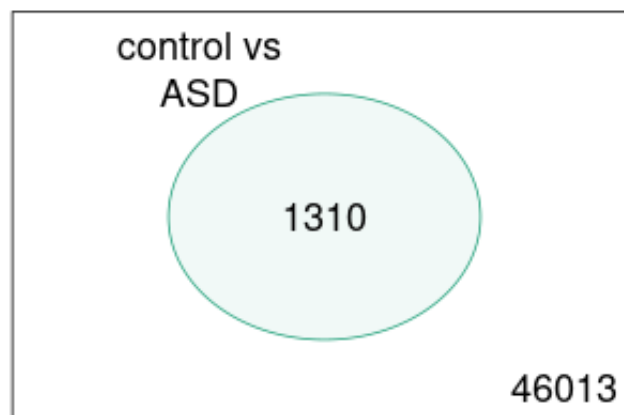
GSE42133: Adjusted P-value counts

```
# summarize test results as "up", "down" or "not expressed"
dT <- decideTests(fit2, adjust.method="fdr", p.value=0.05)

# Venn diagram of results
vennDiagram(dT, circle.col=palette())
# download significant genes in genes.tsv
```
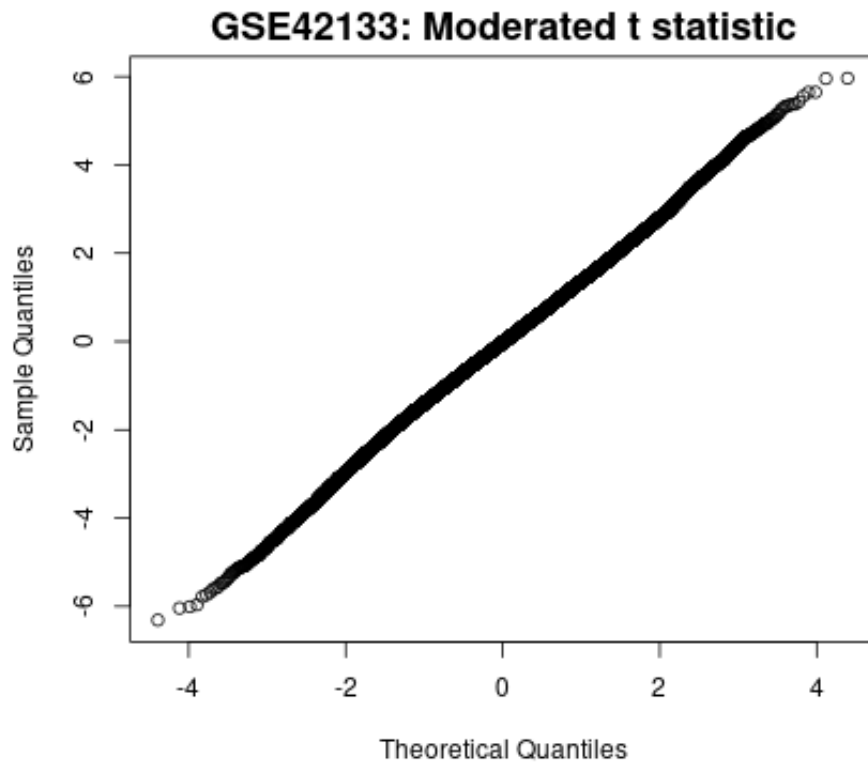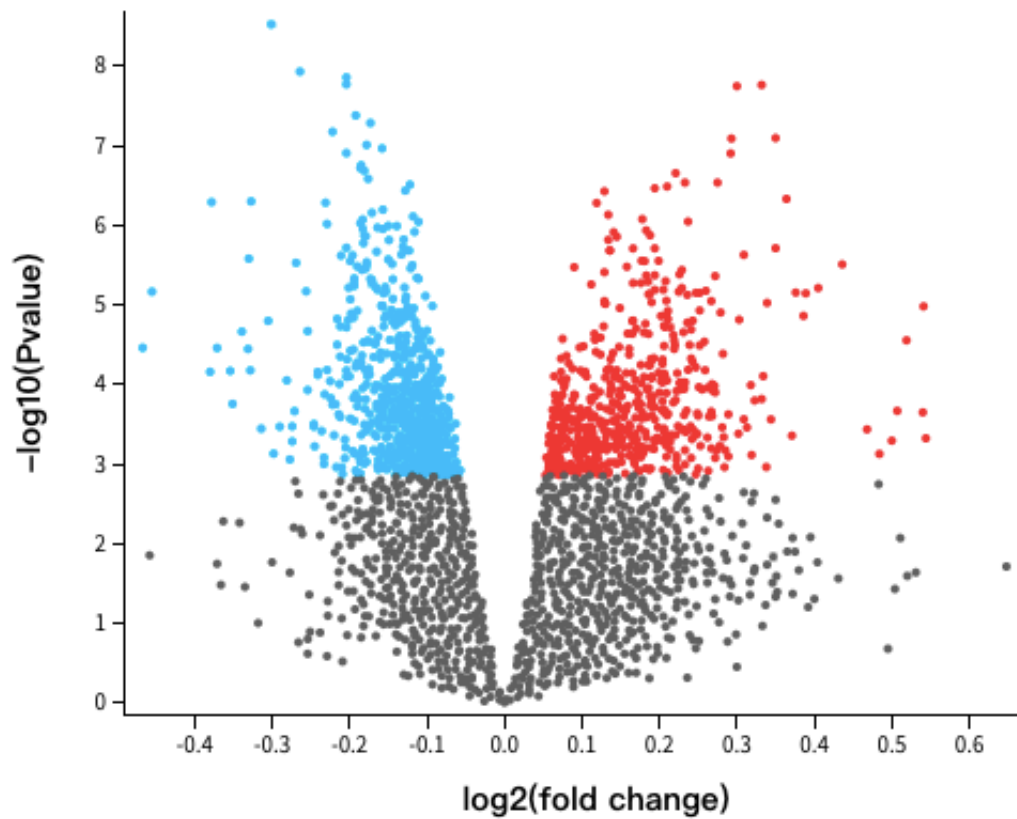
GSE42133: limma, Padj<0.05

```
# create Q-Q plot for t-statistic
t.good <- which(!is.na(fit2$F)) # filter out bad probes
qqt(fit2$t[t.good], fit2$df.total[t.good], main="Moderated t statistic")
```



```
# volcano plot (log P-value vs log fold change)
colnames(fit2) # list contrast names
ct <- 1         # choose contrast of interest
volcanoplot(fit2, coef=ct, main=colnames(fit2)[ct], pch=20,
  highlight=length(which(dT[,ct]!=0)), names=rep('+', nrow(fit2)))
```

# Volcano plot
## GSE42133: Disrupted functional neworks in autism underlie early brain...
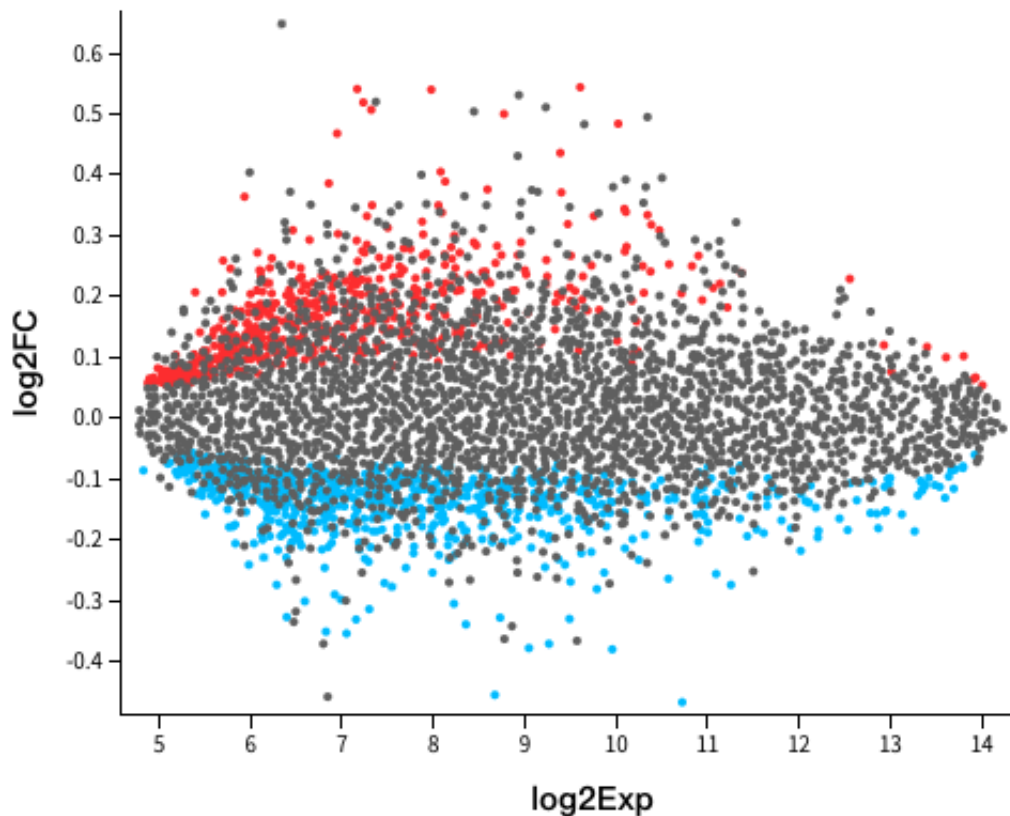### control vs ASD, Padj<0.05



```
# MD plot (log fold change vs mean log expression)
# highlight statistically significant (p-adj < 0.05) probes
plotMD(fit2, column=ct, status=dT[,ct], legend=F, pch=20, cex=1)
abline(h=0)
```
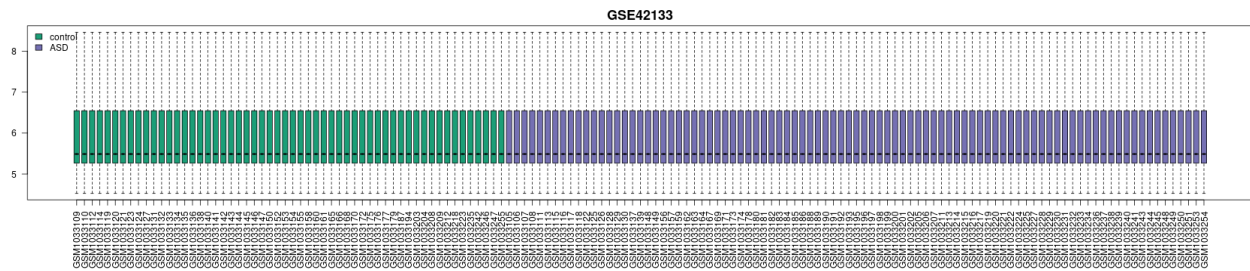
## Meandiff plot
## GSE42133: Disrupted functional neworks in autism underlie early brain...
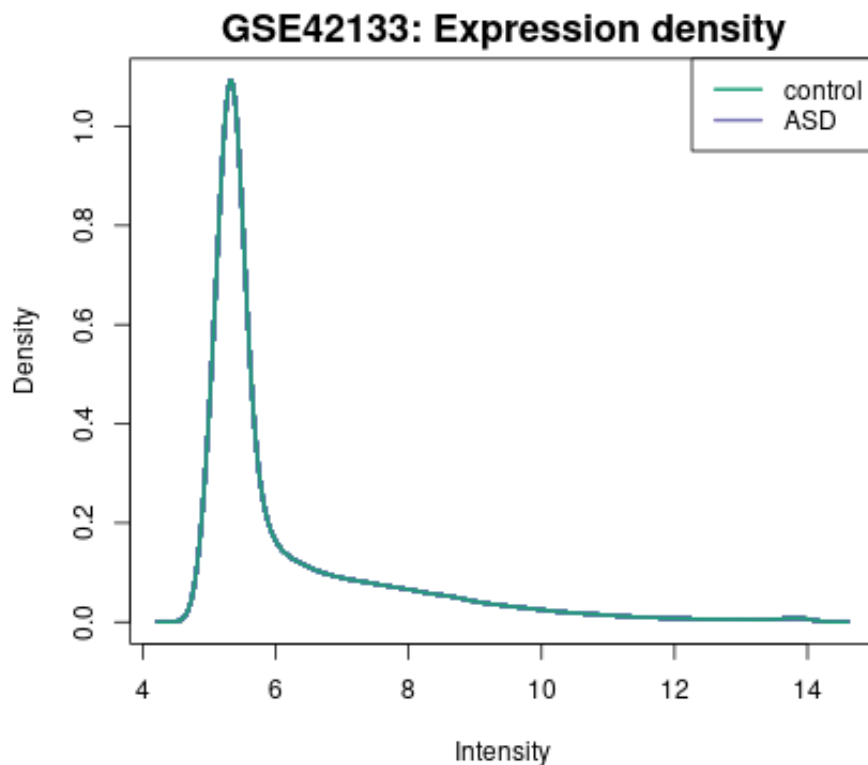### control vs ASD, Padj<0.05



```
###############################################################
# General expression data analysis
ex <- exprs(gset)

# box-and-whisker plot
dev.new(width=3+ncol(gset)/6, height=5)
ord <- order(gs)  # order samples by group
palette(c("#1B9E77", "#7570B3", "#E7298A", "#E6AB02", "#D95F02",
         "#66A61E", "#A6761D", "#B32424", "#B324B3", "#666666"))
par(mar=c(7,4,2,1))
title <- paste ("GSE42133", "/", annotation(gset), sep ="")
boxplot(ex[,ord], boxwex=0.6, notch=T, main=title, outline=FALSE, las=2,
col=gs[ord])
legend("topleft", groups, fill=palette(), bty="n")
dev.off()
```
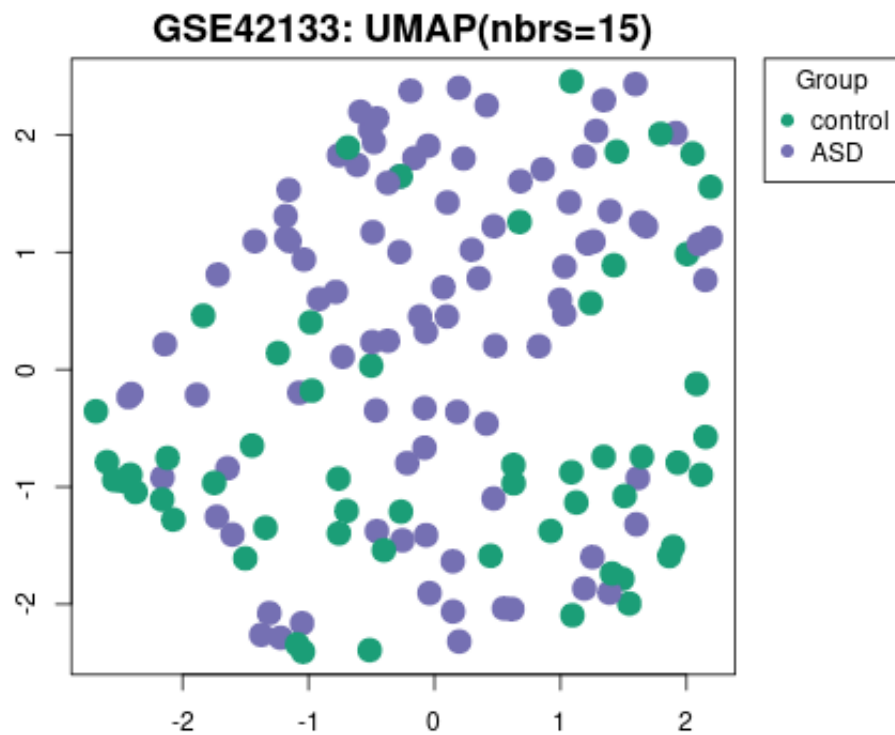
GSE42133

```
# expression value distribution
par(mar=c(4,4,2,1))
title <- paste ("GSE42133", "/", annotation(gset), " value distribution", sep
="")
plotDensities(ex, group=gs, main=title, legend ="topright")
```



GSE42133: Expression density

```
# UMAP plot (dimensionality reduction)
ex <- na.omit(ex) # eliminate rows with NAs
ex <- ex[!duplicated(ex), ]  # remove duplicates
ump <- umap(t(ex), n_neighbors = 15, random_state = 123)
par(mar=c(3,3,2,6), xpd=TRUE)
plot(ump$layout, main="UMAP plot, nbrs=15", xlab="", ylab="", col=gs, pch=20,
cex=1.5)
legend("topright", inset=c(-0.15,0), legend=levels(gs), pch=20,
col=1:nlevels(gs), title="Group", pt.cex=1.5)
library("maptools")  # point labels without overlaps
pointLabel(ump$layout, labels = rownames(ump$layout), method="SANN", cex=0.6)
```

GSE42133: UMAP(nbrs=15)

```
# mean-variance trend, helps to see if precision weights are needed
plotSA(fit2, main="Mean variance trend, GSE42133")
```



GSE42133 : Mean-variance trend