# ASSIGNMENT 11

## Exercise 2.2.8

### Question 1

1. In this question you will write an R loop to perform 1000 t-tests on randomly generated normal data with mean 0. Therefore the null hypothesis is always true but you will see that one still obtains small $p$ values (therefore indicating falsely signifcant results or false positives).

```
1  pvalues <- sapply(1:1000, function(i) {
2    t.test(rnorm(20,0,1))$p.value
3  })
```
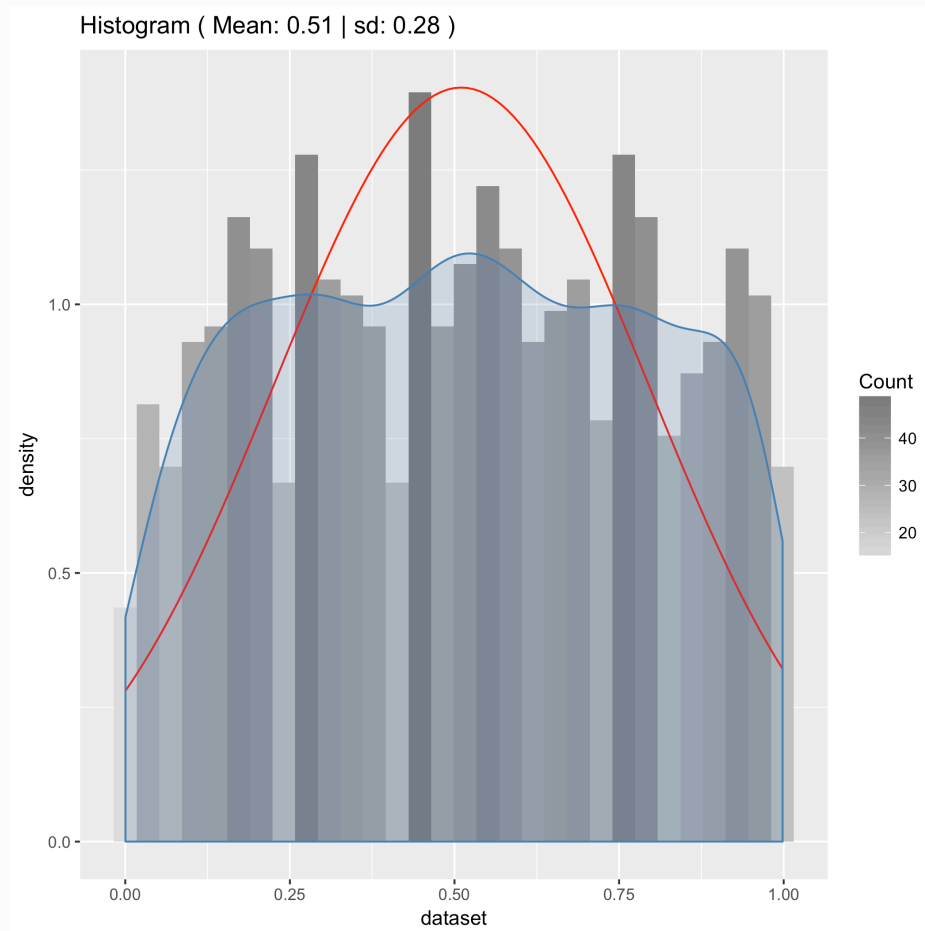
### Question 2

2. Using the which function, count the fraction of times you observe $p$ values less than 0.1, 0.05, 0.01.

```
1  length(which(pvalues<0.1))
2  # 102
3  length(which(pvalues<0.05))
4  # 50
5  length(which(pvalues<0.01))
6  # 10
```

### Question 3

3. Make a histogram of the $p$ values obtained. What distribution is this?

```
1   library('ggplot2')
2   draw_histogram <- function(dataset) {
3     dist_mean <- mean(dataset)
4     dist_sd <- sd(dataset)
5     gg <- ggplot(as.data.frame(dataset), aes(dataset))
6     gg <- gg + geom_histogram(aes(y=..density.., fill=..count..))
7     gg <- gg + scale_fill_gradient("Count", low="#DCDCDC", high="#7C7C7C")
8     gg <- gg + stat_function(fun=dnorm, color="red",
9                             args=list(mean=dist_mean,sd=dist_sd))
10    # Adds a density plot on top
11    gg <- gg + geom_density(alpha = 0.2, fill="steelblue", colour="steelblue")
12    gg <- gg + ggtitle(paste("Histogram", "( Mean:", round(dist_mean,2), '|',
13                            "sd:", signif(dist_sd,2), ")"))
14    return(gg)
15  }
16
17  draw_histogram(pvalues)
18  ggsave('Assignment/histogram.png')
```

Histogram ( Mean: 0.51 | sd: 0.28 )

As shown by the above histogram, the data seems to follow a continuous uniform distribution.
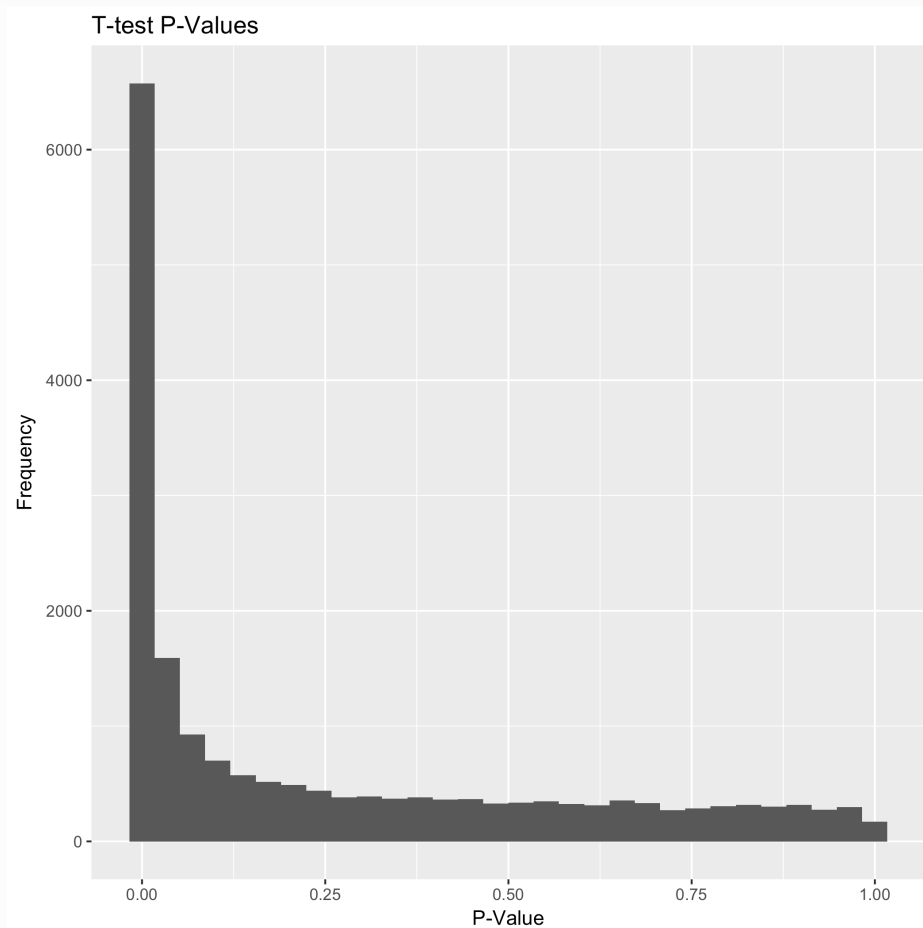
# Exercise 2.2.12

## Question 1

1. Repeat the analysis above to examine the genes that are differentially expressed after 24 hour exposure to cadmium. Provide a histogram of the t-test p values for this dataset.

```r
1    library("GEOquery")
2    library('ggplot2')
3
4    gds       <- getGEO(filename='GDS3420_full.soft')
5    geData    <- Table(gds)
6    geData2   <- geData[which( geData['Platform_ORF'] != "" ),]
7
8    # Identify Samples for each factor
9    hrs4          <- unlist( strsplit( Meta(gds)$sample_id[1], ',' ) )
10   hrs24         <- unlist( strsplit( Meta(gds)$sample_id[2], ',' ) )
11
12   # separate the data for the two factors
13   geData_hrs4  <- data.matrix(geData2[, names(geData2) %in% hrs4])
14   geData_hrs24 <- data.matrix(geData2[, names(geData2) %in% hrs24])
15
16   top_hits <- function (dataMatrix, genes, identifiers, hits_type='upregulated',
17                         fc=2, pvalue=2.57e-6) {
18     numNA       <- apply(dataMatrix, MARGIN=1, function(x) sum(is.na(x)) )
19     indices     <- which(numNA <= 5)
20
21     dataMatrix <- dataMatrix[indices, ]
22     genenames <- genes[indices, ]
23     symbols    <- identifiers[indices]
24
25     foldchange <- apply(dataMatrix, MARGIN=1, function(x){
26       median( exp(as.numeric(x)), na.rm=T)
27     })
28
29     tresults <- apply(dataMatrix, MARGIN=1, function(x){
30       t.test(as.numeric(x))$p.value
31     })
32
33     results <- data.frame(pvalue=tresults, FC=foldchange,
34                           id=genenames, symbol=symbols)
35
36     gg <- ggplot(results, aes(tresults)) + geom_histogram()
37     gg <- gg + scale_fill_gradient("Count", low="#DCDCDC", high="#7C7C7C")
38     gg <- gg + ggtitle("T-test P-Values") + xlab('P-Value') + ylab('Frequency')
39
40     if (hits_type == 'upregulated') {
41       hits <- results[ which( results$pvalue < pvalue & results$FC > fc ), ]
42     } else if (hits_type == 'downregulated') { #
43       hits <- results[ which( results$pvalue < pvalue & results$FC < fc ), ]
44     } else {
45       cat("ERROR: The hit_type variable in the top_hits() function must be
     'upregulated' or 'downregulated'.", file=stderr())
46       quit(save = "no", status = 1, runLast = FALSE)
47     }
48     hits <- hits[order(hits$FC, decreasing=T),]
49
50     return (list(hits=hits, histogram=gg))
51   }
52
53   results_hrs4  <- top_hits(geData_hrs4, geData2["Gene ID"], geData2[,
     'IDENTIFIER'])
54   results_hrs24 <- top_hits(geData_hrs24, geData2["Gene ID"], geData2[,
     'IDENTIFIER'])
55
56   results_hrs24$histogram
57   ggsave('Assignment/results_hrs24.png')
```

T-test P-Values

## Question 2

2. How many genes are significantly up-regulated?

```
1   results_hrs24$hits
2   dim(results_hrs24$hits)
3   # [1] 99   4
```

As shown above, there 99 hits that are significantly up-regulated.

## Question 3

3. How many genes are significantly down regulated ($FC < 0.5$)

```
1   down_regulated_hrs24 <- top_hits(geData_hrs24, geData2["Gene ID"],
2                                     geData2[, 'IDENTIFIER'], 'downregulated', 0.5)
3   down_regulated_hrs24$hits
4   dim(down_regulated_hrs24$hits)
5   # [1] 6 4
```

As show above, there are 6 genes that are significantly down regulated.