

# Multiple Testing

Chapter 5, Lab 4

*OpenIntro Biostatistics*

## Topics

- Experiment-wise error
- Controlling experiment-wise error for correlated data

The previous lab in this chapter introduced the Bonferroni correction as a method for controlling Type I error when conducting pairwise  $t$ -tests following an ANOVA. This lab more formally introduces the multiple testing problem and discusses one specific approach for controlling Type I error: controlling the **experiment-wise error rate**, the probability of making at least one Type I error in a set of hypothesis tests.

The material in this lab is an extension to Section 5.5 of *OpenIntro Biostatistics*.

## Introduction

Recall that making a Type I error, rejecting the null hypothesis when the null hypothesis is true, occurs with probability  $\alpha$ . The Type I error rate is controlled by rejecting the null hypothesis only when the  $p$ -value of a test is smaller than  $\alpha$ . For a single hypothesis test conducted at significance level  $\alpha = 0.05$ , there is a 5% chance of incorrectly rejecting the null hypothesis.

1. With each iteration of the following code, 100 values are randomly sampled from a standard normal distribution and a one-sample  $t$ -test of  $H_0 : \mu = 0$  is conducted, with the  $p$ -value stored. The vector `reject` records whether the  $p$ -value for a test is smaller than  $\alpha$ .

Run the simulation to obtain an empirical estimate of the Type I error rate when one hypothesis test is conducted. Confirm that the Type I error rate is approximately 5%.

```
#set parameters
num.iterations = 1000
alpha = 0.05

#set seed
set.seed(2018)

#create empty list
p.values.A = vector("numeric", num.iterations)

#run simulation
for(k in 1:num.iterations){

  A = rnorm(100)
```

```

p.values.A[k] = t.test(A, mu = 0)$p.val
}

#view results
reject = (p.values.A <= alpha)
table(reject)

```

When conducting more than one  $t$ -test, the significance level  $\alpha$  used in each test controls the error rate for that test. The **experiment-wise error rate** is the chance that at least one test will incorrectly reject  $H_0$  when all tested null hypotheses are true.

2. When two hypothesis tests are conducted at  $\alpha = 0.05$ , is the probability of making at least one Type I error equal to 0.05?
  - a) Modify the simulation code to estimate the experiment-wise error rate (i.e., the ‘overall’ Type I error rate) when two hypothesis tests are conducted; let the vector B contain a set of 100 values randomly drawn from a standard normal distribution.
  - b) Using an algebraic approach, calculate the probability of making at least one Type I error when conducting two hypothesis tests, if the null hypotheses in both cases are true. Assume independence between the tests.
3. Suppose that 100 independent two-sample  $t$ -tests are conducted. What is the probability of at least one incorrect rejection of  $H_0$  at  $\alpha = 0.05$ , given that in all cases there is no difference between the population group means?

### Multiple testing in the Golub leukemia dataset

The Golub leukemia dataset was introduced in Chapter 1 (Lab 3). To investigate whether gene expression profiling could be a tool for classifying acute leukemia type, Golub and co-authors used DNA microarrays to measure the expression level of 7,129 genes from children known to have either acute myeloblastic leukemia (AML) or acute lymphoblastic leukemia (ALL). The goal of the experiment was to identify genes that are differentially expressed between individuals with AML versus ALL.

The analysis from Chapter 1 used a “data-driven” approach, searching for genes that appeared substantially differentially expressed relative to the distribution of differences in mean expression levels between AML and ALL patients. No claims were made regarding whether observed differences were more extreme than expected by chance alone.

In this lab, a hypothesis testing approach will be used to assess whether, for a particular gene  $i$ , there is significant evidence that the mean expression level among ALL patients is different from the mean expression level among AML patients. For simplicity, the analysis will be conducted on a subset of 100 genes rather than the full set of 7,129 genes.

4. Run the following code to load the Golub data and prepare it for analysis; the code was introduced in Chapter 1, Lab 3.

```

#load the data
library(oibistat)

```

```

data(golub)

#remove phenotype information (in the first 6 columns)
gene.matrix = as.matrix(golub[, -(1:6)])

#set parameters
num.genes = ncol(gene.matrix)
num.genes.subset = 100

#create matrix with expression data from subset of 100 genes
set.seed(2401)
gene.index.set = sample(1:num.genes, size = num.genes.subset, replace = FALSE)
gene.matrix.sample = gene.matrix[, gene.index.set]

#create logical variable for cancer type
leuk.type = (golub$cancer == "aml")

```

- a) Produce a plot to show the association between cancer type and expression levels for the gene in the first column of `gene.matrix.sample`. Describe what you see. What do the values TRUE and FALSE mean?
  - b) Conduct a *t*-test comparing expression levels between AML and ALL patients for the gene in the first column of `gene.matrix.sample`. Summarize the results.
5. Run the following code to conduct 100 *t*-tests, comparing expression levels between AML and ALL patients for each of the 100 genes. The *p*-values are stored in the vector `p.vals`.

```

#set parameters
alpha = 0.05

#create empty vector to store results
p.vals = vector("numeric", num.genes.subset)

#conduct a t-test for each gene
for(k in 1:num.genes.subset){

  p.vals[k] = t.test(gene.matrix.sample[, k] ~ leuk.type,
                    alternative = "two.sided", mu = 0,
                    var.equal = FALSE, conf.level = 0.95)$p.val

}

#create table of results
gene.names = colnames(gene.matrix.sample)
results = cbind(gene.names, p.vals)

#is the p-value smaller than alpha?
reject = (p.vals <= alpha)

```

```
#view results
```

- a) Plot a histogram of the  $p$ -values. Describe what you see.
  - b) The logical vector `reject` has value TRUE if the  $p$ -value of a test is smaller than  $\alpha = 0.05$ .
    - i. Of the 100  $t$ -tests conducted, how many result in rejecting the null hypothesis at  $\alpha = 0.05$ ? In other words, how many of the 100 genes are identified as being significantly differently expressed between AML and ALL patients?
    - ii. Do you think it is reasonable to expect that for all the genes identified as significantly differently expressed between AML and ALL patients, there is actually an association between expression level and leukemia type? In other words, is it reasonable to expect that for each rejection, the null hypothesis was correctly rejected?
6. One approach to controlling the experiment-wise error rate is the Bonferroni correction. Using the Bonferroni correction, how many genes out of the 100 are identified as being significantly differently expressed between AML and ALL patients?

From theory, it can be shown that the Bonferroni correction is too strict when comparisons are not independent, inflating the rate of false negatives (i.e., failing to reject when there is a difference between groups). In this setting, independence is not a realistic assumption; from a biological perspective, the expression level of each gene is unlikely to be completely independent of the expression level of another gene.

7. The following simulation estimates the experiment-wise error rate for 100 tests conducted on the Golub gene expression data under the null hypothesis that there is no difference between the population mean expression of AML versus ALL patients.

The simulation is based on randomly sampling values from a distribution related to the normal distribution (specifically, the multivariate normal) in which the variance is described not by a single value, but a matrix of values known as a covariance matrix. This matrix can be thought of as measuring the pairwise correlation between expression levels of each of the 100 genes being tested. Further details of the multivariate normal distribution are beyond the scope of this course.

```
#load package for using multivariate normal distribution
library("MASS")

#set seed
set.seed(2401)

#set parameters
num.patients = nrow(golub)
num.replicates = 1000
alpha = 0.05

#set parameters for multivariate normal distribution
theoretical.means = rep(0, len = num.genes.subset)
estimated.correlation = cov(gene.matrix.sample)
```

```

#create empty vectors to store results
gene.p.vals = vector("numeric", num.genes.subset)
min.p.vals = vector("numeric", num.replicates)

#run the simulation
for(k in 1:num.replicates){

  gene.expression.matrix = mvrnorm(num.patients, mu = theoretical.means,
                                   Sigma = estimated.correlation)

  for(j in 1:num.genes.subset){

gene.p.vals[j] = t.test(gene.expression.matrix[, j] ~ leuk.type)$p.val

  }

  min.p.vals[k] = min(gene.p.vals)

}

#view results
reject = (min.p.vals <= alpha)
table(reject)

```

- a) Run the simulation. What does the vector `gene.p.vals` contain, versus `min.p.vals`?
- b) Create a histogram of `min.p.vals`. Describe the section of the histogram that denotes the estimated experiment-wise error.
- c) From the simulation, what is the estimated experiment-wise error rate? Does this differ from the answer to Question 3?
- d) How does the estimated experiment-wise error change as  $\alpha$  decreases from 0.05? Test values of 0.005 and 0.001 for  $\alpha$ .
- e) Run the following code to find the value for  $\alpha$  that will yield an experiment-wise error of approximately 0.05. Compare this value to the  $\alpha^*$  value recommended by the Bonferroni correction.

```

sim.based.correction = quantile(min.p.vals, 0.05)
sim.based.correction

```

- f) Why is it not advisable to use an overly strict correction for  $\alpha$ ?
- g) Using the simulation-based adjustment to  $\alpha$ , how many genes out of the 100 are identified as being significantly differentially expressed between AML and ALL patients?