

```

1   Course: Statistical_Inference
2   Lesson: Multiple_Testing
3
4
5   - Class: text
6   Output: "Multiple_Testing. (Slides for this and other Data Science courses may be
found at github https://github.com/DataScienceSpecialization/courses/. If you care to
use them, they must be downloaded as a zip file and viewed locally. This lesson
corresponds to 06_Statistical_Inference/12_MultipleTesting.)"
7
8   - Class: text
9   Output: In this lesson, we'll discuss multiple testing. You might ask, "What's that?"
10
11  - Class: text
12  Output: Given that data is valuable and we'd like to get the most out of it, we might
use it to test several hypotheses. If we have an alpha level of .05 and we test 20
hypotheses, then on average, we expect one error, just by chance.
13
14  - Class: text
15  Output: Another potential problem is that after running several tests, only the
lowest p-value might be reported OR all p-values under some threshold might be
considered significant. Undoubtedly, some of these would be false.
16
17  - Class: text
18  Output: Luckily, we have clever ways of minimizing errors in this situation. That's
what we'll address. We'll define specific error measures and then statistical ways
of correcting or limiting them.
19
20
21  - Class: text
22  Output: Multiple testing is particularly relevant now in this age of BIG data.
Statisticians are tasked with questions such as "Which variables matter among the
thousands measured?" and "How do you relate unrelated information?"
23
24  - Class: mult_question
25  Output: Since multiple testing addresses compensating for errors let's review what we
know about them. A Type I error is
26  AnswerChoices: rejecting a false hypothesis; failing to reject a false hypothesis;
rejecting a true hypothesis; failing to reject a true hypothesis
27  CorrectAnswer: rejecting a true hypothesis
28  AnswerTests: omnitest(correctVal='rejecting a true hypothesis')
29  Hint: Eliminate the two choices that are not errors. A Type I error involves rejection.
30
31  - Class: mult_question
32  Output: In an American court, an example of a Type I error is
33  AnswerChoices: convicting an innocent person; acquitting a guilty person; letting the
indicted off on a technicality
34  CorrectAnswer: convicting an innocent person
35  AnswerTests: omnitest(correctVal='convicting an innocent person')
36  Hint: In an American court, the null hypothesis is that the accused is innocent. If
he is convicted when he really is innocent then the null hypothesis is rejected
incorrectly.
37
38
39  - Class: mult_question
40  Output: A Type II error is
41  AnswerChoices: rejecting a false hypothesis; failing to reject a false hypothesis;
rejecting a true hypothesis; failing to reject a true hypothesis
42  CorrectAnswer: failing to reject a false hypothesis
43  AnswerTests: omnitest(correctVal='failing to reject a false hypothesis')
44  Hint: Eliminate the two choices that are not errors. A Type II error involves failing
to reject.
45
46  - Class: mult_question
47  Output: In an American court, an example of a Type II error is
48  AnswerChoices: convicting an innocent person; acquitting a guilty person; letting the
indicted off on a technicality
49  CorrectAnswer: acquitting a guilty person
50  AnswerTests: omnitest(correctVal='acquitting a guilty person')

```

51 **Hint:** In an American court, the null hypothesis is that the accused is innocent. If this is accepted (not rejected) by the jury and the defendant is in fact guilty a Type II error has been made.

52

53 - **Class:** mult_question

54 **Output:** Good. Let's continue reviewing. The null hypothesis

55 **AnswerChoices:** represents the status_quo and is assumed true; tells us the origins of the number 0; is never true; is a big nothing that statisticians like to gossip about

56 **CorrectAnswer:** represents the status_quo and is assumed true

57 **AnswerTests:** omnitest(correctVal='represents the status_quo and is assumed true')

58 **Hint:** Really? Only one choice seems reasonable.

59

60 - **Class:** text

61 **Output:** The p-value is "the probability under the null hypothesis of obtaining evidence as or more extreme than your test statistic (obtained from your observed data) in the direction of the alternative hypothesis." Of course p-values are related to significance or alpha levels, which are set before the test is conducted (often at 0.05).

62

63 - **Class:** mult_question

64 **Output:** If a p-value is found to be less than alpha (say 0.05), then the test result is considered statistically significant, i.e., surprising and unusual, and the null hypothesis (the status quo) is ?

65 **AnswerChoices:** accepted; rejected; revised; renamed the aleph null hypothesis

66 **CorrectAnswer:** rejected

67 **AnswerTests:** omnitest(correctVal='rejected')

68 **Hint:** Accepted (failed to be rejected) or rejected are the only real choices here. A low p-value is a low probability. This means your data is unusual and is closer to the alternative hypothesis than the null.

69

70 - **Class:** figure

71 **Output:** Now consider this chart copied from http://en.wikipedia.org/wiki/Familywise_error_rate. Suppose we've tested m null hypotheses, m_0 of which are actually true, and m-m_0 are actually false. Out of the m tests R have been declared significant, that is, the associated p-values were less than alpha, and m-R were nonsignificant, or boring results.

72 **Figure:** errorMat.R

73 **FigureType:** new

74

75 - **Class:** mult_question

76 **Output:** Looking at the chart, which variables are known?

77 **AnswerChoices:** m and R; m_0, and m; S,T,U,V; A,B,C

78 **CorrectAnswer:** m and R

79 **AnswerTests:** omnitest(correctVal='m and R')

80 **Hint:** The number of hypotheses tested (m) and the number declared significant (R) are known. The variable m_0 represents the unknowable, the number of true hypotheses. S, T, U, and V are unobservable random variables.

81

82 - **Class:** mult_question

83 **Output:** In testing the m_0 true null hypotheses, V results were declared significant, that is, these tests favored the alternative hypothesis. What type of error does this represent?

84 **AnswerChoices:** Type I; Type II; Type III; a serious one

85 **CorrectAnswer:** Type I

86 **AnswerTests:** omnitest(correctVal='Type I')

87 **Hint:** By declaring the test result significant the true null hypothesis was rejected, like convicting an innocent person.

88

89 - **Class:** text

90 **Output:** Another name for a Type I error is False Positive, since it is falsely claiming a significant (positive) result.

91

92 - **Class:** mult_question

93 **Output:** Of the m-m_0 false null hypotheses, T were declared nonsignificant. This means that these T null hypotheses were accepted (failed to be rejected). What type of error does this represent?

94 **AnswerChoices:** Type I; Type II; Type III; a serious one

95 **CorrectAnswer:** Type II

96 **AnswerTests:** omnitest(correctVal='Type II')

97 **Hint:** By declaring the test result nonsignificant the false null hypothesis was
 98 accepted (failed to be rejected), like letting a guilty person go free.

99 - **Class:** text
 100 **Output:** Another name for a Type II error is False Negative, since it is falsely
 claiming a nonsignificant (negative) result.

101 - **Class:** text
 102 **Output:** A rose by any other name, right? Consider the fraction V/R .

103 - **Class:** text
 104 **Output:** The observed R represents the number of test results declared significant.
 105 These are 'discoveries', something different from the status quo. V is the number of
 106 those falsely declared significant, so V/R is the ratio of FALSE discoveries. Since V
 is a random variable (i.e., unknown until we do an experiment) we call the expected
 value of the ratio, $E(V/R)$, the False Discovery Rate (FDR).

107 - **Class:** text
 108 **Output:** A rose by any other name, right? How about the fraction V/m_0 ? From the
 109 chart, m_0 represents the number of true H_0 's and m_0 is unknown. V is the number of
 those falsely declared significant, so V/m_0 is the ratio of FALSE positives. Since V
 is a random variable (i.e., unknown until we do an experiment) we call the expected
 value of the ratio, $E(V/m_0)$, the FALSE POSITIVE rate.

110 - **Class:** mult_question
 111 **Output:** Another good name for the false positive rate would be
 112 **AnswerChoices:** false alarm rate; the Type II rate; a rose; a thorn
 113 **CorrectAnswer:** false alarm rate
 114 **AnswerTests:** omnitest(correctVal='false alarm rate')
 115 **Hint:** False positives are Type I errors so one of the only two sensible answers is
 116 incorrect.

117 - **Class:** mult_question
 118 **Output:** The false positive rate would be closely related to
 119 **AnswerChoices:** the Type I error rate; the Type II error rate; a thorny rose;
 120 **CorrectAnswer:** the Type I error rate
 121 **AnswerTests:** omnitest(correctVal='the Type I error rate')
 122 **Hint:** False positives are Type I errors so one of the only two sensible answers is
 123 incorrect.

124 - **Class:** text
 125 **Output:** We call the probability of at least one false positive, $\Pr(V \geq 1)$ the Family
 126 Wise Error Rate (FWER).

127 - **Class:** text
 128 **Output:** So how do we control the False Positive Rate?

129 - **Class:** text
 130 **Output:** Suppose we're really smart, calculate our p-values correctly, and declare
 131 all tests with $p < \alpha$ as significant. This means that our false positive rate is
 132 at most α , on average.

133 - **Class:** mult_question
 134 **Output:** Suppose we perform 10,000 tests and $\alpha = .05$. How many false positives do
 135 we expect on average?
 136 **AnswerChoices:** 500; 5000; 50; 50000
 137 **CorrectAnswer:** 500
 138 **AnswerTests:** omnitest(correctVal='500')
 139 **Hint:** Multiply 10000 by .05 to get the correct answer.

140 - **Class:** text
 141 **Output:** You got it! 500 false positives seems like a lot. How do we avoid so many?

142 - **Class:** text
 143 **Output:** We can try to control the family-wise error rate (FWER), the probability of
 144 at least one false positive, with the Bonferroni correction, the oldest multiple
 145 testing correction.

146 - **Class:** text
 147

```

148 Output: It's very straightforward. We do  $m$  tests and want to control the FWER at
    level  $\alpha$  so that  $\Pr(V \geq 1) < \alpha$ . We simply reduce  $\alpha$  dramatically. Set
     $\alpha_{\text{fwer}}$  to be  $\alpha/m$ . We'll only call a test result significant if its p-value  $<$ 
     $\alpha_{\text{fwer}}$ .
149
150 - Class: mult_question
151 Output: Sounds good, right? Easy to calculate. What would be a drawback with this
    method?
152 AnswerChoices: too many results will pass; too many results will fail; requires too
    much math
153 CorrectAnswer: too many results will fail
154 AnswerTests: omnitest(correctVal='too many results will fail')
155 Hint: Dividing  $\alpha$  by  $m$  makes your cutoff value very small so you might not get any
    significant results, much less false ones.
156
157 - Class: text
158 Output: Another way to limit the false positive rate is to control the false
    discovery rate (FDR). Recall this is  $E(V/R)$ . This is the most popular correction when
    performing lots of tests. It's used in lots of areas such as genomics, imaging,
    astronomy, and other signal-processing disciplines.
159
160 - Class: text
161 Output: Again, we'll do  $m$  tests but now we'll set the FDR, or  $E(V/R)$  at level  $\alpha$ .
    We'll calculate the p-values as usual and order them from smallest to largest,  $p_1,$ 
 $p_2, \dots, p_m$ . We'll call significant any result with  $p_i \leq (\alpha \cdot i)/m$ . This is the
    Benjamini-Hochberg method (BH). A p-value is compared to a value that depends on its
    ranking.
162
163 - Class: text
164 Output: This is equivalent to finding the largest  $k$  such that  $p_k \leq (k \cdot \alpha)/m$ ,
    (for a given  $\alpha$ ) and then rejecting all the null hypotheses for  $i=1, \dots, k$ .
165
166
167
168 - Class: text
169 Output: Like the Bonferroni correction, this is easy to calculate and it's much less
    conservative. It might let more false positives through and it may behave strangely
    if the tests aren't independent.
170
171 - Class: figure
172 Output: Now consider this chart copied from the slides. It shows the p-values for 10
    tests performed at the  $\alpha=.2$  level and three cutoff lines. The p-values are shown
    in order from left to right along the x-axis. The red line is the threshold for No
    Corrections (p-values are compared to  $\alpha=.2$ ), the blue line is the Bonferroni
    threshold,  $\alpha=.2/10 = .02$ , and the gray line shows the BH correction. Note that it
    is not horizontal but has a positive slope as we expect.
173 Figure: corrMat1.R
174 FigureType: new
175
176 - Class: mult_question
177 Output: With no correction, how many results are declared significant?
178 AnswerChoices: 2; 4; 6 ;8
179 CorrectAnswer: 4
180 AnswerTests: omnitest(correctVal='4')
181 Hint: How many points fall below the red line?
182
183 - Class: mult_question
184 Output: With the Bonferroni correction, how many tests are declared significant?
185 AnswerChoices: 2; 4; 6; 8
186 CorrectAnswer: 2
187 AnswerTests: omnitest(correctVal='2')
188 Hint: How many points fall below the blue line?
189
190 - Class: mult_question
191 Output: So the Bonferroni passed only half the results that the No Correction
    (comparing p-values to  $\alpha$ ) method passed. Now look at the BH correction. How many
    tests are significant with this scale?
192 AnswerChoices: 1; 3; 5; 7
193 CorrectAnswer: 3

```

```

194 AnswerTests: omnitest(correctVal='3')
195 Hint: How many points fall below the gray line?.
196
197 - Class: text
198 Output: So the BH correction which limits the FWER is between the No Correction and
the Bonferroni. It's more conservative (fewer significant results) than the No
Correction but less conservative (more significant results) than the Bonferroni. Note
that with this method the threshold is proportional to the ranking of the values so
it slopes positively while the other two thresholds are flat.
199
200 - Class: text
201 Output: Notice how both the Bonferroni and BH methods adjusted the threshold (alpha)
level of rejecting the null hypotheses. Another equivalent corrective approach is to
adjust the p-values, so they're not classical p-values anymore, but they can be
compared directly to the original alpha.
202
203 - Class: text
204 Output: Suppose the p-values are  $p_1, \dots, p_m$ . With the Bonferroni method you
would adjust these by setting  $p'_i = \max(m * p_i, 1)$  for each p-value. Then if you
call all  $p'_i < \alpha$  significant you will control the FWER.
205
206
207 - Class: figure
208 Output: To demonstrate some of these concepts, we've created an array of p-values for
you. It is 1000-long and the result of a linear regression performed on random normal
x,y pairs so there is no true significant relationship between the x's and y's.
209 Figure: genNoTrue.R
210 FigureType: new
211
212 - Class: cmd_question
213 Output: Use the R command head to see the first few entries of the array pValues.
214 CorrectAnswer: head(pValues)
215 AnswerTests: omnitest(correctExpr='head(pValues)')
216 Hint: Type head(pValues) at the command prompt.
217
218 - Class: cmd_question
219 Output: Now count the number of entries in the array that are less than the value
.05. Use the R command sum, and the appropriate Boolean expression.
220 CorrectAnswer: sum(pValues < 0.05)
221 AnswerTests: omnitest(correctExpr='sum(pValues < 0.05)')
222 Hint: Type sum(pValues < 0.05) at the command prompt.
223
224 - Class: cmd_question
225 Output: "So we got around 50 false positives, just as we expected (.05*1000=50). The
beauty of R is that it provides a lot of built-in statistical functionality. The
function p.adjust is one example. The first argument is the array of pValues. Another
argument is the method of adjustment. Once again, use the R function sum and a
boolean expression using p.adjust with method="bonferroni" to control the FWER."
226 CorrectAnswer: sum(p.adjust(pValues,method="bonferroni") < 0.05)
227 AnswerTests: omnitest(correctExpr='sum(p.adjust(pValues,method="bonferroni") <
0.05)')
228 Hint: Type sum(p.adjust(pValues,method="bonferroni") < 0.05) at the command prompt.
229
230 - Class: cmd_question
231 Output: "So the correction eliminated all the false positives that had passed the
uncorrected alpha test. Repeat the same experiment, this time using the method "BH"
to control the FDR."
232 CorrectAnswer: sum(p.adjust(pValues,method="BH") < 0.05)
233 AnswerTests: omnitest(correctExpr='sum(p.adjust(pValues,method="BH") < 0.05)')
234 Hint: Type sum(p.adjust(pValues,method="BH") < 0.05) at the command prompt.
235
236 - Class: figure
237 Output: So the BH method also eliminated all the false positives. Now we've generated
another 1000-long array of p-values, this one called pValues2. In this data, the
first half ( 500 x/y pairs) contains x and y values that are random and the second
half contain x and y pairs that are related, so running a linear regression model on
the 1000 pairs should find some significant (not random) relationship.
238 Figure: gen50True.R
239 FigureType: new

```

```

240
241 - Class: cmd_question
242 Output: We also created a 1000-long array of character strings, trueStatus. The
first 500 entries are "zero" and the last are "not zero". Use the R function tail to
look at the end of trueStatus.
243 CorrectAnswer: tail(trueStatus)
244 AnswerTests: omnitest(correctExpr='tail(trueStatus)')
245 Hint: Type tail(trueStatus) at the command prompt.
246
247 - Class: cmd_question
248 Output: Once again we can use R's greatness to count and tabulate for us. We can call
the R function table with two arguments, a boolean such as pValues2<.05, and the
array trueStatus. The boolean obviously has two outcomes and each entry of trueStatus
has one of two possible values. The function table aligns the two arguments and
counts how many of each combination (TRUE,"zero"), (TRUE,"not zero"), (FALSE,"zero"),
and (FALSE,"not zero") appear. Try it now.
249 CorrectAnswer: table(pValues2 < 0.05, trueStatus)
250 AnswerTests: omnitest(correctExpr='table(pValues2 < 0.05, trueStatus)')
251 Hint: Type table(pValues2 < 0.05, trueStatus) at the command prompt.
252
253 - Class: text
254 Output: "We see that without any correction all 500 of the truly significant
(nonrandom) tests were correctly identified in the \"not zero\" column. In the zero
column (the truly random tests), however, 24 results were flagged as significant."
255
256 - Class: cmd_question
257 Output: What is the percentage of false positives in this test?
258 CorrectAnswer: 24/500
259 AnswerTests: equiv_val(.048)
260 Hint: "Divide 24 by 500 to get the percentage."
261
262 - Class: text
263 Output: "Just as we expected - around 5% or .05*100."
264
265 - Class: cmd_question
266 Output: "Now run the same table function, however, this time use the call to p.adjust
with the \"bonferroni\" method in the boolean expression. This will control the FWER."
267 CorrectAnswer: "table(p.adjust(pValues2,method=\"bonferroni\") < 0.05, trueStatus)"
268 AnswerTests: omnitest(correctExpr='table(p.adjust(pValues2,method=\"bonferroni\") <
0.05, trueStatus)')
269 Hint: "Type table(p.adjust(pValues2,method=\"bonferroni\") < 0.05, trueStatus) at the
command prompt."
270
271 - Class: text
272 Output: Since the Bonferroni correction method is more conservative than just
comparing p-values to alpha all the truly random tests are correctly identified in
the zero column. In other words, we have no false positives. However, the threshold
has been adjusted so much that 23 of the truly significant results have been
misidentified in the not zero column.
273
274 - Class: cmd_question
275 Output: "Now run the same table function one final time. Use the call to p.adjust
with \"BH\" method in the boolean expression. This will control the false discovery
rate."
276 CorrectAnswer: "table(p.adjust(pValues2,method=\"BH\") < 0.05, trueStatus)"
277 AnswerTests: omnitest(correctExpr='table(p.adjust(pValues2,method=\"BH\") < 0.05,
trueStatus)')
278 Hint: "Type table(p.adjust(pValues2,method=\"BH\") < 0.05, trueStatus) at the command
prompt."
279
280 - Class: text
281 Output: "Again, the results are a compromise between the No Corrections and the
Bonferroni. All the significant results were correctly identified in the \"not zero\"
column but in the random (\"zero\") column 13 results were incorrectly identified.
These are the false positives. This is roughly half the number of errors in the other
two runs."
282
283 - Class: figure
284 Output: Here's a plot of the two sets of adjusted p-values, Bonferroni on the left

```

and BH on the right. The x-axis indicates the original p-values. For the Bonferroni, (adjusting by multiplying by 1000, the number of tests), only a few of the adjusted values are below 1. For the BH, the adjusted values are slightly larger than the original values.

Figure: plot2.R

FigureType: new

- **Class:** text

Output: We'll conclude by saying that multiple testing is an entire subfield of statistical inference. Usually a basic Bonferroni/BH correction is good enough to eliminate false positives, but if there is strong dependence between tests there may be problems. Another correction method to consider is "BY".

- **Class:** text

Output: Congrats! We hope you liked the multiple concepts and questions you saw in this lesson.

- **Class:** mult_question

Output: "Would you like to receive credit for completing this course on Coursera.org?"

CorrectAnswer: NULL

AnswerChoices: Yes;No

AnswerTests: coursera_on_demand()

Hint: ""