

Course: Regression_Models
Lesson: MultiVar_Examples2

- **Class:** text

Output: "MultiVar_Examples2. (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 Regression_Models/02_02_multivariateExamples.)"

- **Class:** text

Output: This is the second lesson in which we'll look at some regression models with more than one independent variable. We'll begin with the InsectSprays data which we've taken the liberty to load for you. This data is part of R's datasets package. It shows the effectiveness of different insect sprays. We've used the code from the slides to show you a boxplot of the data.

- **Class:** exact_question

Output: How many Insect Sprays are in this dataset?

CorrectAnswer: 6

AnswerTests: omnitest(correctVal=6)

Hint: How many letters are in the sequence of letters A, B, C, D, E, F ?

- **Class:** text_question

Output: From the boxplot, which spray has the largest median?

CorrectAnswer: B

AnswerTests: omnitest(correctVal='B')

Hint: The median of each spray is indicated by the thick horizontal line in each box. Which spray has its median furthest from 0?

- **Class:** cmd_question

Output: Let's first try to get a better understanding of the dataset InsectSprays. Use the R function dim to find the dimensions of the data.

CorrectAnswer: dim(InsectSprays)

AnswerTests: omnitest(correctExpr='dim(InsectSprays)')

Hint: Type "dim(InsectSprays)" at the R prompt.

- **Class:** cmd_question

Output: The R function dim says that InsectSprays is a 72 by 2 array. Use the R function head to look at the first 15 elements of InsectSprays.

CorrectAnswer: head(InsectSprays,15)

AnswerTests: omnitest(correctExpr='head(InsectSprays,15)')

Hint: Type "head(InsectSprays,15)" at the R prompt.

- **Class:** cmd_question

Output: So this dataset contains 72 counts, each associated with a particular different spray. The counts are in the first column and a letter identifying the spray in the second. To save you some typing we've created 6 arrays with just the count data for each spray. The arrays have the names sx, where x is A,B,C,D,E or F. Type one of the names (your choice) of these arrays to see what we're talking about.

CorrectAnswer: sA

AnswerTests: ANY_of_exprs('sA','sB','sC','sD','sE','sF')

Hint: Type sB at the R prompt.

- **Class:** cmd_question

Output: As a check, run the R command summary on the second column of the dataset to see how many entries we have for each spray. (Recall that the expression M[,2] yields the second column of the array M.)

CorrectAnswer: summary(InsectSprays[,2])

AnswerTests: omnitest(correctExpr='summary(InsectSprays[,2])')

Hint: Type "summary(InsectSprays[,2])" at the R prompt.

- **Class:** text

Output: It's not surprising that with 72 counts we'd have 12 count for each of the 6 sprays. In this lesson we'll consider multilevel factor levels and how we interpret linear models of data with more than 2 factors.

```

52 - Class: cmd_question
53 Output: Use the R function apply to find out the classes of the columns of the data.
54 CorrectAnswer: apply(InsectSprays,class)
55 AnswerTests: omnittest(correctExpr='apply(InsectSprays,class)')
56 Hint: Type "apply(InsectSprays,class)" at the R prompt.
57
58 - Class: text
59 Output: The class of the second "spray" column is factor. Recall from the slides that
the equation representing the relationship between a particular outcome and several
factors contains binary variables, one for each factor. This data has 6 factors so
we need 6 dummy variables. Each will indicate if a particular outcome (a count) is
associated with a specific factor or category (insect spray).
60
61 - Class: cmd_question
62 Output: Using R's lm function, generate the linear model in which count is the
dependent variable and spray is the independent. Recall that in R formula has the
form y ~ x, where y depends on the predictor x. The data set is InsectSprays. Store
the model in the variable fit.
63 CorrectAnswer: fit <- lm(count ~ spray, InsectSprays)
64 AnswerTests: creates_lm_model('fit <- lm(count ~ spray, InsectSprays)')
65 Hint: Type "fit <- lm(count ~ spray, InsectSprays)" at the R prompt.
66
67 - Class: cmd_question
68 Output: Using R's summary function, look at the coefficients of the model. Recall
that these can be accessed with the R construct x$coef.
69 CorrectAnswer: summary(fit)$coef
70 AnswerTests: omnittest(correctExpr='summary(fit)$coef')
71 Hint: Type "summary(fit)$coef" at the R prompt.
72
73 - Class: cmd_question
74 Output: Notice that R returns a 6 by 4 array. For convenience, store off the first
column of this array, the Estimate column, in a variable called est. Remember the R
construct for accessing the first column is x[,1].
75 CorrectAnswer: est <- summary(fit)$coef[,1]
76 AnswerTests: omnittest(correctExpr='est <- summary(fit)$coef[,1]')
77 Hint: Type "est <- summary(fit)$coef[,1]" at the R prompt.
78
79
80 - Class: text
81 Output: Notice that sprayA does not appear explicitly in the list of Estimates. It is
there, however, as the first entry in the Estimate column. It is labeled as
"(Intercept)". That is because sprayA is the first in the alphabetical list of the
levels of the factor, and R by default uses the first level as the reference against
which the other levels or groups are compared when doing its t-tests (shown in the
third column).
82
83 - Class: cmd_question
84 Output: What do the Estimates of this model represent? Of course they are the
coefficients of the binary or dummy variables associated with sprays. More
importantly, the Intercept is the mean of the reference group, in this case sprayA,
and the other Estimates are the distances of the other groups' means from the
reference mean. Let's verify these claims now. First compute the mean of the sprayA
counts. Remember the counts are all stored in the vectors named sx. Now we're
interested in finding the mean of sA.
85 CorrectAnswer: mean(sA)
86 AnswerTests: omnittest(correctExpr='mean(sA)')
87 Hint: Type "mean(sA)" at the R prompt.
88
89 - Class: mult_question
90 Output: What do you think the mean of sprayB is?
91 AnswerChoices: 15.3333; 0.83333; -12.41667; I haven't a clue
92 CorrectAnswer: 15.3333
93 AnswerTests: omnittest(correctVal='15.3333')
94 Hint: Adding the value of the Intercept to the Estimate for sprayB yields the
empirical mean of sprayB.
95
96 - Class: cmd_question
97 Output: Verify this now by using R's mean function to compute the mean of sprayB.
98 CorrectAnswer: mean(sB)

```

```

99   AnswerTests: omnitest(correctExpr='mean(sB)')
100  Hint: Type "mean(sB)" at the R prompt.
101
102  - Class: cmd_question
103  Output: Let's generate another model of this data, this time omitting the intercept.
  We can easily use R's lm function to do this by appending " - 1" to the formula,
  e.g., count ~ spray - 1. This tells R to omit the first level. Do this now and store
  the new model in the variable nfit.
104  CorrectAnswer: nfit <- lm(count ~ spray - 1, InsectSprays)
105  AnswerTests: creates_lm_model('nfit <- lm(count ~ spray - 1, InsectSprays)')
106  Hint: Type "nfit <- lm(count ~ spray - 1, InsectSprays)" at the R prompt.
107
108  - Class: cmd_question
109  Output: Now, as before, look at the coefficient portion of the summary of nfit.
110  CorrectAnswer: summary(nfit)$coef
111  AnswerTests: omnitest(correctExpr='summary(nfit)$coef')
112  Hint: Type "summary(nfit)$coef" at the R prompt.
113
114  - Class: text
115  Output: Notice that sprayA now appears explicitly in the list of Estimates. Also
  notice how the values of the columns have changed. The means of all the groups are
  now explicitly shown in the Estimate column. Remember that previously, with an
  intercept, sprayA was excluded, its mean was the intercept, and the values for the
  other sprays (estimates, standard errors, and t-tests) were all computed relative to
  sprayA, the reference group. Omitting the intercept clearly affected the model.
116
117  - Class: mult_question
118  Output: What values does the Estimate column now show?
119  AnswerChoices: The means of all 6 levels; The variances of all 6 levels; I have no idea
120  CorrectAnswer: The means of all 6 levels
121  AnswerTests: omnitest(correctVal='The means of all 6 levels')
122  Hint: The numbers should look familiar, especially for sprayA and sprayB. What values
  have you computed for these two sprays?
123
124  - Class: mult_question
125  Output: Without an intercept (reference group) the tests are whether the expected
  counts (the groups means) are different from zero. Which spray has the least
  significant result?
126  AnswerChoices: sprayC; sprayF; sprayB; sprayA
127  CorrectAnswer: sprayC
128  AnswerTests: omnitest(correctVal='sprayC')
129  Hint: Which spray has the highest probability?
130
131  - Class: text
132  Output: Clearly, which level is first is important to the model. If you wanted a
  different reference group, for instance, to compare sprayB to sprayC, you could refit
  the model with a different reference group.
133
134  - Class: cmd_question
135  Output: The R function relevel does precisely this. It re-orders the levels of a
  factor. We'll do this now. We'll call relevel with two arguments. The first is the
  factor, in this case InsectSprays$spray, and the second is the level that we want to
  be first, in this case "C". Store the result in a new variable spray2.
136  CorrectAnswer: spray2 <- relevel(InsectSprays$spray, "C")
137  AnswerTests: omnitest(correctExpr='spray2 <- relevel(InsectSprays$spray, "C")')
138  Hint: Type "spray2 <- relevel(InsectSprays$spray, \"C\")" at the R prompt.
139
140  - Class: cmd_question
141  Output: Now generate a new linear model and put the result in the variable fit2.
142  CorrectAnswer: fit2 <- lm(count ~ spray2, InsectSprays)
143  AnswerTests: creates_lm_model('fit2 <- lm(count ~ spray2, InsectSprays)')
144  Hint: Type "fit2 <- lm(count ~ spray2, InsectSprays)" at the R prompt.
145
146
147  - Class: cmd_question
148  Output: As before, look at the coef portion of the summary of this new model fit2.
  See how sprayC is now the intercept (since it doesn't appear explicitly in the list).
149  CorrectAnswer: summary(fit2)$coef
150  AnswerTests: omnitest(correctExpr='summary(fit2)$coef')

```

```

151     Hint: Type "summary(fit2)$coef" at the R prompt.
152
153 - Class: mult_question
154 Output: According to this new model what is the mean of spray2C?
155 AnswerChoices: 2.083333; 12.416667; 14.583333; The model doesn't tell me.
156 CorrectAnswer: 2.083333
157 AnswerTests: omnitest(correctVal='2.083333')
158 Hint: Recall that the intercept is the mean of the reference group, in this case
    sprayC, so look at the value in the (Intercept) row of the Estimate column.
159
160 - Class: cmd_question
161 Output: Verify your answer with R's mean function using the array sC as the argument.
162 CorrectAnswer: mean(sC)
163 AnswerTests: omnitest(correctExpr='mean(sC)')
164 Hint: Type "mean(sC)" at the R prompt.
165
166
167 - Class: mult_question
168 Output: According to this new model what is the mean of spray2A?
169 AnswerChoices: 14.50000; 12.416667; 14.583333; I don't have a clue
170 CorrectAnswer: 14.50000
171 AnswerTests: omnitest(correctVal='14.50000')
172 Hint: Recall that when there is an intercept, the mean of a level that's not the
    reference, is the intercept + the coefficient (or estimate) of that level, in this
    case spray2a, so you'll have to add together two numbers. Alternatively, just look
    back and see what the mean was for the original model.
173
174 - Class: text
175 Output: Remember that with this model sprayC is the reference group, so the t-test
    statistics (shown in column 3 of the summary coefficients) compare the other sprays
    to sprayC. These can be computed by hand using the Estimates and standard error from
    the original model (fit) which used sprayA as the references.
176
177 - Class: cmd_question
178 Output: The slides show the details of this but here we'll demonstrate by calculating
    the spray2B t value. Subtract fit's sprayC coefficient (fit$coef[3]) from sprayB's
    (fit$coef[2]) and divide by the standard error which we saw was 1.6011. The result is
    spray2B's t value. Do this now.
179 CorrectAnswer: (fit$coef[2]-fit$coef[3])/1.6011
180 AnswerTests: omnitest(correctExpr='(fit$coef[2]-fit$coef[3])/1.6011')
181 Hint: Type "(fit$coef[2]-fit$coef[3])/1.6011" at the R prompt.
182
183 - Class: text
184 Output: We glossed over some details in this lesson. For instance, counts can never
    be 0 so the assumption of normality is violated. We'll explore this issue more when
    we discuss Poisson GLMs. For now be glad that you've concluded this second lesson on
    multivariable linear models.
185
186 - Class: mult_question
187 Output: "Would you like to receive credit for completing this course on
    Coursera.org?"
188 CorrectAnswer: NULL
189 AnswerChoices: Yes;No
190 AnswerTests: coursera_on_demand()
191 Hint: ""
192
193

```