

Sign up **x**

I have a question about which is the best way to specify an interaction in a regression model. Consider the following data:

r	regression	interaction
---	------------	-------------

 Manuel Ramón
810 1 8 13

which indicates (not definitively) that the underlying models are the same, just expressed (to the lm internals) differently.

If you define your interaction as `paste(ds, dr)` instead of `paste(dr, ds)` your parameter estimates will change again, in interesting ways.

Note how in your model summary for `lm1` the coefficient estimate for `ss2` is 4.94 lower than in the summary for `lm2`, with the coefficient for `rr2:ss2` being 4.95 (if you print to 3 decimal places, the difference goes away). This is another indication that an internal rearrangement of terms has occurred.

I can't think of any advantage to doing it yourself, but there may be one with more complex models where you don't want a full interaction term but instead only some of the terms in the "cross" between two or more factors.

answered Dec 2 '11 at 22:40



jbowman

12.8k 1 26 53

The only advantage I see to define the interactions as in `lm2` is that it is easy to perform multiple comparisons for the interaction term. What I do not quite understand is why different results are obtained if, in principle, it seems that the 2 approaches are the same. — Manuel Ramón Dec 3 '11 at 20:27

4 The approaches are the same, but the exact parameterizations of the model being estimated are different, so the results appear different. Consider a model with two binary regressors x_1, x_2 and an interaction. You have four categories, but you can write the model several different ways, e.g., let 1 be a constant term, with variables $(1, x_1, x_2, x_1 * x_2)$ or $(x_1, x_2, x_1 * x_2, (1 - x_1) * (1 - x_2))$, or others. The variables are just linear combinations of each other. The coefficient estimates will be different, but the model is really the same. — jbowman Dec 3 '11 at 20:43

Therefore, although different, both approaches are correct, aren't it? — Manuel Ramón Dec 3 '11 at 20:56

Right. Mathematically the matrices of independent variables in the various formulations are just linear transforms of each other, so the parameter estimates of one model can be calculated from the parameter estimates of another if one knows how the two models were actually set up. — jbowman Dec 4 '11 at 22:13

You might understand this behavior better if you look at the model matrices.

```
model.matrix(lm1 <- lm(y ~ r*s, data=d))
model.matrix(lm2 <- lm(y ~ r + s + rs, data=d))
```

When you look at these matrices, you can compare the constellations of $s_2=1$ with the other variables (i.e. when $s_2=1$, which values do the other variables take?). You will see that these constellations differ slightly, which just means that the base category is different. Everything else is essentially the same. In particular, note that in your `lm1`, the coefficient on `ss2` equals the coefficients `ss2+rsr1s2` of `lm2`, i.e. $3.82=8.76-4.95$, short of rounding errors.

For instance, executing the following code gives you exactly the same output as using the automatic setting of R:

```
d$rs <- relevel(d$rs, "r1s1")
summary(lm1 <- lm(y~ factor(r) + factor(s) + factor(rs), data=d))
```

This also provides a quick answer to your question: the really only reason to change the way factors are set up is to provide expositional clarity. Consider the following example: Suppose you regress wage on a dummy for high school completion interacted with a factor indicating if you belong to a minority.

That is: $wage = \alpha + \beta edu + \gamma edu * minority + \epsilon$

If said minority factor takes value 1 if you do belong to a minority, the coefficient β can be interpreted as a wage difference for non-minority individuals who have completed high school. If this is your coefficient of interest, then you should code it as such. Otherwise, suppose the minority factor takes the value of 1 if you do not belong to a minority. Then, in order to see how much more non-minority individuals earn when they complete high school, you would have to "manually" compute $\beta + \gamma$. Note though that all information is contained in the estimates though, and substantial results do not change by setting up the factors differently!

edited Apr 12 '14 at 19:45

answered Apr 12 '14 at 19:32



coffeinjunky

516 4 12