7 Two-factor Experiments, Continued

In the last lecture and in lab we dealt with the parameters **Stata** (and most software packages) use to fit the additive and the interaction models for two-way ANOVA. The lincom command was one way we learned to deal with the parameters. That probably is not the easiest approach in many problems, however. If we learn a little more about the parameters, some information is fairly immediate.

Let's continue with the insecticide problem, where we have 4 poisons and 3 doses. Remember the pattern of population cell means as follows (ignoring marginal means for now):

		Dose	
Insecticide	1	2	3
1	μ_{11}	μ_{12}	μ_{13}
2	μ_{21}	μ_{22}	μ_{23}
3	μ_{31}	μ_{32}	μ_{33}
4	μ_{41}	μ_{42}	μ_{43}

Consider the full interaction model first. The parameterization for this is $\mu_{ij} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij}$, i = 1, ..., 4; j = 1, ..., 3. We dodged the issue of constraints last time, but recall the problem: There are 12 real parameters (the μ_{ij}), but 20 new parameters (1 + 4 + 3 + 12). We need to put 8 constraints (restrictions) on these new parameters to bring us back down to 12. An old standard textbook solution to this, and one that makes the math look a lot simpler (for marginal means at least) is

$$0 = \sum_{i} \alpha_{i} = \sum_{j} \beta_{j} = \sum_{i} (\alpha \beta)_{ij} = \sum_{j} (\alpha \beta)_{ij}$$

(that looks like 9 constraints but one is redundant so it is 8). Software packages like **Stata** and **SAS** use an algorithm called the sweep algorithm that makes a completely different and much more useful set of constraints more natural, though. Effectively, they start adding parameters in the model and as soon as they hit a redundant one, they set the new parameter to 0. The new constraints become

$$0 = \alpha_4 = \beta_3 = (\alpha\beta)_{41} = (\alpha\beta)_{42} = (\alpha\beta)_{43} = (\alpha\beta)_{13} = (\alpha\beta)_{23} = (\alpha\beta)_{33}$$

i.e. if there are I levels of i and J levels of j then any time i reaches level I or j reaches level J then the parameter becomes 0. It is a little easier to see 8 constraints here.

If we plug in these constraints and rewrite all 12 cell means, we see the following pattern:

		Dose	
Insecticide	1	2	3
1	$\mu + \alpha_1 + \beta_1 + (\alpha\beta)_{11}$	$\mu + \alpha_1 + \beta_2 + (\alpha\beta)_{12}$	$\mu + \alpha_1$
2	$\mu + \alpha_2 + \beta_1 + (\alpha\beta)_{21}$	$\mu + \alpha_2 + \beta_2 + (\alpha\beta)_{22}$	$\mu + \alpha_2$
3	$\mu + \alpha_3 + \beta_1 + (\alpha\beta)_{31}$	$\mu + \alpha_3 + \beta_2 + (\alpha\beta)_{32}$	$\mu + \alpha_3$
4	$\mu + \beta_1$	$\mu + \beta_2$	μ

At first this may not seem like much simplification, but let's examine it a bit more carefully. μ is often referred to as the grand mean (this comes from the old textbook parameterization) but here we see $\mu = \mu_{43}$. The last cell of the table has become the *reference group* with all other parameters being deviations from that reference group. β_2 is the difference between doses 2 and 3 for the 4th poison, β_1 is the difference between doses 1 and 3 for the 4th poison. α_3 is the difference between poisons 3 and 4 for the 3^{rd} dose, α_2 is the difference between poisons 2 and 4 for the 3^{rd} dose, and α_1 is the difference between poisons 1 and 4 for the 3^{rd} dose. The difference between poisons 3 and 4 for the 2^{nd} dose is $\alpha_3 + (\alpha\beta)_{32}$ rather than just α_3 (if the difference does not depend on poison then there is no interaction).

With these constraints then $(\mu_{32} - \mu_{42}) - (\mu_{33} - \mu_{43}) = (\alpha\beta)_{32}$. Recall last lecture we said that **no** interaction (parallel profiles in an interaction plot) was this: If i, i', j, j' are legal indexes, then $\mu_{ij} - \mu_{ij'} = \mu_{i'j} - \mu_{i'j'}$, which is to say the difference between doses j and j' is the same for insecticide i as for insecticide i'; and $\mu_{ij} - \mu_{i'j} = \mu_{ij'} - \mu_{i'j'}$, which is to say the difference between insecticides i and i' is the same for dose j as it is for dose j'. Last week we saw that $[\mu_{ij} - \mu_{ij'}] - [\mu_{i'j} - \mu_{i'j'}] = [(\alpha\beta)_{ij} - (\alpha\beta)_{ij'}] - [(\alpha\beta)_{i'j'}] - (\alpha\beta)_{i'j'}]$. The constraints that $0 = (\alpha\beta)_{Ij} = (\alpha\beta)_{iJ}$ allow us to simplify greatly the lincom command for many such terms. Since $(\mu_{32} - \mu_{42}) - (\mu_{33} - \mu_{43}) = (\alpha\beta)_{32}$ (only for these constraints, though!) then a simple lincom(_b[poison[3]*dose[2]]) gets that term. Better yet, we can have it automatically printed.

Stata Implementation

Anova problems actually are specialized regression problems (we will grapple with this idea later). What we want are regression estimates of all the effects $(\mu, \alpha_i, \beta_j, (\alpha\beta)_{ij})$. The regress option with anova gets that for us in a form that matches the lincom syntax. This can be treated as a post-estimation command, i.e. after issuing the anova time poison dose poison*dose command (and examining the the ANOVA table, interaction plots, etc.) just type another command anova, regress to get the following results

. anova S	, r our	egr ce	ess 	SS	df		MS		Number of obs	=	48
Res	Mod idu	el al 	+	2.20435628 .800724989	11 36	. 20	0396025 2242361		Prob > F R-squared Adi R-squared	- - -	9.01 0.0000 0.7335 0.6521
	Tot	al	İ	3.00508126	47	.06	3937899		Root MSE	=	.14914
	ti	me		Coef.	Std.	Err.	t	P> t	[95% Conf.	Int	terval]
_cons				. 325	.0749	5694	4.36	0.000	.1737663	.4	1762337
poison		1 2 3 4		115 .01 09 (dropped)	.105 .105 .105	5457 5457 5457	-1.09 0.09 -0.85	0.283 0.925 0.399	3288767 2038767 3038767	.()988767 2238767 1238767
dose		1 2 3		.285 .3425 (dropped)	.105 .105	5457 5457	2.70 3.25	$\substack{0.010\\0.003}$.0711233 .1286233	.4	1988767 5563767
poison*	dos 1	e 1		- 0825	1491	1387	-0 55	0 584	- 3849674		0199674
	1	23		2325	:1493	1387	-1.56	0.128	5349673	.(5699674
	222	1 2 3		.26 .1375 (dropped)	.1492 .1492	1387 1387	$\substack{1.74\\0.92}$	0.090 0.363	0424673 1649673	. [5624674 1399674
	3 3 3 4	1 2 3 1		(dropped) (dropped)	.1492 .1492	1387 1387	0.32 -1.36	0.752 0.183	2549674 5049673	.:	3499674)999674
	4 4	2 3		(dropped) (dropped)							

Let's try to make sense of these coefficients by relating them both to the table of sample cell means and to an interaction plot from the last lecture. First, try to reconstruct the sample cell means from the coefficients. Recall that the full interaction model imposes no restrictions. Various questions arise. Why is poison highly significant yet none of the poison coefficients is significant? – Get your answer from the interaction plot and what these coefficients are estimating. One of the large differences appears to be between poisons 1 and 2 at dose level 2. How would you estimate that difference? How would you test for specific interactions (i.e. different slopes in the above plot?). We will spend some time examining all this.

		Dose		
Insecticide	1	2	3	Insect marg
1	.413	.320	.210	.314
2	.880	.815	.335	.677
3	.568	.375	.235	.393
4	.610	.668	.325	.534
Dose marg	.618	.544	.277	.480



The Additive Model

Since interaction was not significant (there is not much data so this *might* just be poor power) we should see how all this looks when we fit the no-interaction (additive) model. This is a highly restricted model and we will not reproduce all the sample cell means from this model. Now we fit $\mu_{ij} = \mu + \alpha_i + \beta_j$ using the command, with ensuing results given below:

anova time poison dose					
-	Number of obs Root MSE	= .	48 1 . 158179	R-squared Adj R-squared	= 0.6503 = 0.6087
Source	Partial SS	df	f MS	F	Prob > F
Model	1.95421877	5	5.39084375	5 15.62	0.0000
poison	.921206282	З	.30706876	1 12.27	0.0000
dose	1.03301249	2	2 .51650624	6 20.64	0.0000
Residual	1.05086249	42	2 .02502053	6	
Total	3.00508126	47	.06393789	9	

All constraints are as previously described, and easily seen from the following:

•	anova, regres	S				
	Source	SS	df	MS	Number of obs =	48
_	+				F(5, 42) =	15.62
	Model	1.95421877	5	.390843755	Prob > F = 0	0.0000
	Residual	1.05086249	42	.025020536	R-squared = (0.6503
	+				Adj [®] R-squared = (0.6087

	Total	3.00508126	47 .063	937899		Root MSE	= .15818
	time	Coef.	Std. Err.	t	P> t	[95% Conf.	Interval]
_cons poison		.3310417	.0559247	5.92	0.000	.2181811	.4439022
1	1 2 3 4	22 .1425 1416667 (dropped)	.0645762 .0645762 .0645762	-3.41 2.21 -2.19	$\begin{array}{c} 0.001 \\ 0.033 \\ 0.034 \end{array}$	3503201 .0121799 2719868	0896799 .2728201 0113466
aose	1 2 3	.34125 .268125 (dropped)	.0559247 .0559247	6.10 4.79	0.000 0.000	.2283895 .1552645	.4541105 .3809855

What we have fit now is the much simpler structure for population cell means (all the parameters have very easy interpretations – what are they?):

		Dose	
Insecticide	1	2	3
1	$\mu + \alpha_1 + \beta_1$	$\mu + \alpha_1 + \beta_2$	$\mu + \alpha_1$
2	$\mu + \alpha_2 + \beta_1$	$\mu + \alpha_2 + \beta_2$	$\mu + \alpha_2$
3	$\mu + \alpha_3 + \beta_1$	$\mu + \alpha_3 + \beta_2$	$\mu + \alpha_3$
4	$\mu + \beta_1$	$\mu + \beta_2$	μ

You should confirm that you no longer reproduce the sample cell means \bar{y}_{ij} from the estimated regression coefficients, but you do reproduce the sample marginal means \bar{y}_{i} . and \bar{y}_{j} . We can look at an interaction plot of predicted cell means, but note that we have *forced* it to look this way.



Estimable Functions

What we have been covering is expanded upon at considerable length in the **SAS** manual and many textbooks under the topic of *estimable functions*. This is a fairly advanced topic, but the gist of it is that only linear combinations of *population cell means* can legally be estimated (things of the form $\sum_i \sum_j c_{ij} \mu_{ij}$ for some constants c_{ij} – we have been using 0, 1 and -1 as constants). Anything we estimate or test has to be a linear combination of population cell means. In particular, μ and α_i , for instance, are not estimable since there is ambiguity about what they are until constraints are put on them. Different constraints give different interpretations. What we have been doing is relating everything back to the μ_{ij} in order to keep it all "legal". I can run an analysis on two different packages (or even the same package with different options) and get considerably different estimates of α_1 reported. As long as I stick to estimable functions, though, I always get the same estimate.

Parameters in the One-Way Problem

A look back at the one-way ANOVA problem shows the constraints can make some things simpler there too. In that case we have a model $y_{ij} = \mu_i + \epsilon_{ij}$; $i = 1, \ldots, I$; $j = 1, \ldots, n_i$. Let's do the same type of decomposition of μ_i , i.e. $\mu_i = \mu + \alpha_i$, with the same problem that we have I "real" group means and now I + 1 new parameters. We need a constraint, and the one **Stata** and **SAS** impose is $\alpha_I = 0$ (set the first redundant parameter to 0).

What is the implication? Now $\mu = \mu_I$ and $\alpha_i = \mu_i - \mu_I$, i.e. the *last* group has become a reference group and all the parameter estimates (the α_i) are deviations from this reference group. In the CHDS example this means we could get mostly for free the t-tests comparing non-smokers to heavy smokers and light smokers to heavy smokers. It might be more convenient to reorder so that nonsmokers were last, so that the easy tests would compare the two smoking groups to nonsmokers. The lincom command can fill in the last comparison in either case ($\mu_1 - \mu_2 = \alpha_1 - \alpha_2$). The idea of estimable functions still applies; we need to make sure we are looking at linear combinations of the means.

Unbalanced Data

Returning to the two-way problem, we wrote a model

$$y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}; \ i = 1, 2, \dots, I; \ j = 1, 2, \dots, J; \ k = 1, 2, \dots, K$$

where the sample size in each cell (i, j) was the same number K. Unbalanced data allow the sample size to vary with cell, so now

$$y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}; \ i = 1, 2, \dots, I; \ j = 1, 2, \dots, J; \ k = 1, 2, \dots, K_{ij}$$

What difference does it make? — In practice, not much. The reason the topic gets mentioned is that there are disagreeable aspects to unbalanced designs. There are multiple ways to formulate SSs and F-tests. **SAS** provides 4, Types I-IV, and **Stata** provides the same as **SAS** Types I and III. With balanced data all the types agree, but for unbalanced data they do not all agree. The only reason this is not much of a practical problem is that most analysts use Type III (**Stata**'s default) and don't anguish much over it. The t-tests on coefficients are not obviously affected (i.e. **lincom** results), although comparing main effects in the presence of interaction is a subtle business in unbalanced designs (the preferred approach being least squares means, and Stata makes those awkward to get).

Let's return to the unbalanced example of rat insulin levels from the last lecture. The ANOVA table indicates that the vein and time effects are significant, with p-values of .0001 and .0014, respectively, but that the interaction is not significant (p-value=.575). Recall that the jugular and portal profiles are reasonably parallel, which is consistent with a lack of interaction. Looking at the estimates below, let us figure out as an in-class exercise how to interpret the various coefficient estimates, and how to test for significance of the important effects. We really ought to see if we can transform to a scale where the residuals look better, too.

. anova,regres Source	s SS	df	MS		Number of obs	= 48
Model Residual	99478.4833 102914.183	5 19 42 2	0895.6967 2450.3377		F(5, 42) Prob > F R-squared Adi B-squared	= 8.12 = 0.0000 = 0.4915 = 0.4310
Total	202392.667	47 43	806.22695		Root MSE	= 49.501
insulin	Coef.	Std. Err	t. t	P> t	[95% Conf.	Interval]
_cons	128.5	14.28967	8.99	0.000	99.66227	157.3377
1 2	-67.16667 (dropped)	31.95268	-2.10	0.042	-131.6498	-2.68354
ucinttino	-46.58333 44.4 (dropped)	20.20865 21.19501	-2.31 2.09	$0.026 \\ 0.042$	-87.36604 1.626732	-5.800623 87.17327
vein*time 1 1 1 2 1 3 2 1 2 2 2 3	11.85 -26.23333 (dropped) (dropped) (dropped) (dropped)	41.41541 40.9194	0.29 -0.64	0.776 0.525	-71.72969 -108.812	95.42969 56.34536
. anova insuli	n vein time.					
		Number of Root MSE	obs = = 49	48 .0037	R-squared Adj R-squared	= 0.4779 = 0.4424
	Source	Partial	SS df	MS	F	Prob > F
	Model	96732.56	i94 3	32244.18	398 13.43	0.0000
	vein time	51594.46 50332.28	85 1 93 2	51594.40 25166.14	58521.4914710.48	0.0000 0.0002
	Residual	105660.0	97 44	2401.365	585	
anous rogrog	Total	202392.6	67 47	4306.226	395	
Source	SS	df	MS		Number of obs	= 48
Model Residual	96732.5694 105660.097	3 32 44 24	244.1898 01.36585		F(3, 44) Prob > F R-squared	= 13.43 = 0.0000 = 0.4779 = 0.4424
Total	202392.667	47 43	806.22695		Root MSE	= 0.4424 = 49.004
insulin	Coef.	Std. Err	t.	P> t	[95% Conf.	Interval]
_cons	129.6685	13.03897	9.94	0.000	103.3902	155.9468
time	-73.00912 (dropped)	15.75087	-4.64	0.000	-104.7529	-41.26531
1 2 3	-42.54816 35.58493 (dropped)	17.42256 17.82622	-2.44 2.00	0.019 0.052	-77.66102 3414601	-7.435306 71.51132

To see one difference with the unbalanced design, consider the following two ANOVA tables; the first is the usual one, the second is an optional one. Note the differences for SS of main effects. If the data were balanced, the default (SAS Type III SS) and the sequential (SAS Type I SS) would be the same. The second form even depends upon the order terms are entered into the model.

anova	insulin	n vein time	vein*time					
			Number of obs	=	48	R-squared	=	0.4915
			Root MSE	= 49	.5009	Adj ^R -squared	=	0.4310
	_	Source	Partial SS	df	MS	F	P:	rob > F
		Model	99478.4833	5	19895.696	.12 8.12		0.0000
		vein	48212.7037	1	48212.703	19.68		0.0001
		vein*time	2745.9139	2	1372.9569	95 0.56		0.0014 0.5752
		Residual	102914.183	42	2450.337	77		
	-	Total	202392.667	47	4306.2269	95		
anova	insulin	n vein time	vein*time,seq					
			Number of obs	=	48	R-squared	=	0.4915
			Root MSE	= 49	.5009	Adj R-squared	=	0.4310
		Source	Seq. SS	df	MS	F	P	rob > F
	-	Model	99478.4833	5	19895.696	67 8.12		0.0000
		vein	46400.2801	1	46400.280	18.94		0.0001
		vein*time	2745.9139	2	1372.9569	$\frac{10.27}{95}$ 0.56		0.0002
		Residual	102914.183	42	2450.337	(7		

Regression on Dummy Variables: Stata's xi

The way anova works is that it creates special variables called indicator or dummy variables for categorical variables and performs regression on them. One dummy variable is created for each level of a categorical variable and has a value of 1 if the observation has that level of the category else the value is 0. We do not need to worry much about this if we can use the **anova** command, but if we want to do all this in logistic regression we have to get explicit about it. All the big statistics packages do this, and if we were using SAS I could hide it all from you, but Stata requires you to learn about it. The xi facility in Stata is one we will need for many problems.

Following is the insecticide data analyzed as a regression problem using xi:

. xi: regress	time i.poison	i.dose	i.poi	son*i.d	ose		
i.poison	_Ipoison_	1-4	_ (n	aturall	y coded;	; _Ipoison_1 or	mitted)
1.dose	_laose_l-	3	(n	aturall	y coaea;	; _ldose_l oml	ttea)
i.poi~n*i.dose	e _IpoiXdos	_#_#	(c	oded as	above)		
Source	SS	df	M	S		Number of obs $F(11, 36)$	= 48 = 9.01
Model Residual	2.20435628 .800724989	11 36	.20039	6025 2361		Prob > F R-squared	= 0.0000 = 0.7335 = 0.6521
Total	3.00508126	47	.06393	7899		Root MSE	= .14914
time	Coef.	Std. E	rr	t	P> t	[95% Conf.	[Interval]
_Ipoison_2 _Ipoison_3 _Ipoison_4 _Idose_2 _Idose_3 _Ipoison_2 _Ipoison_3 _Ipoison_4	.4675 .155 .1975 0925 2025 (dropped) (dropped) (dropped)	.1054 .1054 .1054 .1054 .1054	57 57 57 57 57 57	4.43 1.47 1.87 -0.88 -1.92	0.000 0.150 0.069 0.386 0.063	.2536233 0588767 0163767 3063767 4163767	.6813767 .3688767 .4113767 .1213767 .0113767

_Idose_2 _Idose_3 _IpoiXdo~2_2 _IpoiXdo~2_3 _IpoiXdo~3_2 _IpoiXdo~3_3 _IpoiXdo~4_2 _IpoiXdo~4_2	(dropped) (dropped) .0275 3425 1 13 .15	.1491387 .1491387 .1491387 .1491387 .1491387 .1491387	0.18 -2.30 -0.67 -0.87 1.01	0.855 0.028 0.507 0.389 0.321	2749674 6449674 4024674 4324674 1524674 2840674	.3299674 0400326 .2024674 .1724674 .4524674 2100674
_IpoiXdo~4_2	.15	.1491387		0.321	1524674	.4524674
cons	.4125	.0745694	5.53	0.000	.2612663	.5637337

You should do this and then examine the variables Stata has placed in the data set. This does not give us the tests from the ANOVA table. The test on interaction is easy:

```
. testparm _IpoiXdo*
 (1)
       _{IpoiXdos_2_2} = 0
 (2)
       _{IpoiXdos_2_3 = 0}
 (3)
       _{IpoiXdos_3_2} = 0
 (4)
       _{IpoiXdos_3_3 = 0}
 (5)
       _{IpoiXdos_4_2} = 0
 (6)
       _{IpoiXdos_4_3 = 0}
       F(
                  36) =
           6
                            1.87
                            0.1123
             Prob > F =
```

but the tests on main effects (equality of marginal means) are a lot less obvious (don't worry, you won't have to do it this way for anova problems). Because this is such a mess, we probably would not test for main effects if interaction were present if we had to go through these steps. With no interaction the procedure is just like the test above.

```
. test _Ipoison_2 + (_IpoiXdos_2_2+_IpoiXdos_2_3)/3 = _Ipoison_3 + (_IpoiXdos_3_2+_IpoiXdo
> s_3_3)/3 = _Ipoison_4 + (_IpoiXdos_4_2+_IpoiXdos_4_3)/3 =0
(1) _Ipoison_2 - _Ipoison_3 + .3333333 _IpoiXdos_2_2 + .3333333 _IpoiXdos_2_3 - .333333
> 3 _IpoiXdos_3_2 - .3333333 _IpoiXdos_3_3 = 0
(2) _Ipoison_2 - _Ipoison_4 + .3333333 _IpoiXdos_2_2 + .3333333 _IpoiXdos_2_3 - .333333
> 3 _IpoiXdos_4_2 - .3333333 _IpoiXdos_4_3 = 0
(3) _Ipoison_2 + .3333333 _IpoiXdos_2_2 + .3333333 _IpoiXdos_2_3 = 0
F( 3, 36) = 13.81
Prob > F = 0.0000
. test _Idose_2+(_IpoiXdos_2_2+_IpoiXdos_3_2+_IpoiXdos_4_2)/4 =_Idose_3+(_IpoiXdos_2_3+_I
> poiXdos_3_3+_IpoiXdos_4_3)/4=0
(1) _Idose_2 - _Idose_3 + .25 _IpoiXdos_2_2 - .25 _IpoiXdos_2_3 + .25 _IpoiXdos_3_2 - .
> 25 _IpoiXdos_3_3 + .25 _IpoiXdos_4_2 - .25 _IpoiXdos_4_3 = 0
(2) _Idose_2 + .25 _IpoiXdos_4_2 - .25 _IpoiXdos_4_3 = 0
F( 2, 36) = 23.22
Prob > F = 0.0000
```

Annoyingly, Stata has changed constraints on us (now the first level of a categorical variable gets zeroed out). We can set what level gets zeroed out (and thus becomes the reference level) with the **char** command

```
. char poison[omit] 4
. char dose[omit] 3
. xi: regress time i.poison i.dose i.poison*i.dose
```

Execute these commands and confirm the original parameters from **anova,regress** are reproduced. We rarely use **xi** and explicit regression on dummy variables in anova problems. We will need it with logistic regression, though.