<u>STAT 525</u>

Chapter 7 General Linear Test and Multicollinearity

Dr. Qifan Song

General Linear Test

- Comparison of a <u>full</u> model and <u>reduced</u> model that involves a subset of full model predictors (i.e., hierarchical structure)
- Involves a comparison of unexplained SS
- Consider a full model with k predictors and reduced model with l predictors (l < k)
- Can show that under null hypothesis

$$F^{\star} = \frac{(\mathsf{SSE}(\mathsf{R}) - \mathsf{SSE}(\mathsf{F}))/(k-l)}{\mathsf{SSE}(\mathsf{F})/(n-k-1)} \sim F \text{ distribution}$$

• Degrees of freedom for F^* are the number of <u>extra</u> variables and the error degrees of freedom for the full model

- Testing the Null hypothesis that the regression coefficients for the <u>extra</u> variables are all zero.
- Examples:
 - X_1, X_2, X_3, X_4 vs $X_1, X_2 \longrightarrow H_0$: $\beta_3 = \beta_4 = 0$ if we are sure that β_5, \ldots are exactly 0's
 - X_1, X_2, X_4 vs $X_1 \longrightarrow H_0$: $\beta_2 = \beta_4 = 0$ if we are sure that β_3, β_5, \ldots are exactly 0's
 - X_1, X_2, X_3, X_4 vs $X_1 \longrightarrow H_0$: $\beta_2 = \beta_3 = \beta_4 = 0$ if we are sure that β_5, \ldots are exactly 0's
- Because SSM+SSE=SSTO, can also compare using explained SS (SSM)

Extra SS and Notation

- Consider $H_0 : X_1, X_3$ vs $H_a : X_1, X_2, X_3, X_4$
- Null can also be written H_0 : $\beta_2 = \beta_4 = 0$
- Write SSE(F) as $SSE(X_1, X_2, X_3, X_4)$
- Write SSE(R) as $SSE(X_1, X_3)$
- Difference in SSE's is the extra SS
- Write as

 $SSE(X_2, X_4 | X_1, X_3) = SSE(X_1, X_3) - SSE(X_1, X_2, X_3, X_4)$

Recall SSM can also be used

 $SSM(X_2, X_4 | X_1, X_3) = SSM(X_1, X_2, X_3, X_4) - SSM(X_1, X_3) \Longrightarrow$ $SSM(X_1, X_2, X_3, X_4) = SSM(X_1, X_3) + SSM(X_2, X_4 | X_1, X_3)$

General Linear Test in Terms of Extra SS

• Can rewrite F test as

$$F^{\star} = \frac{\mathsf{SSE}(X_2, X_4 | X_1, X_3) / (4 - 2)}{\mathsf{SSE}(X_1, X_2, X_3, X_4) / (n - 5)}$$

- Under H_0 , $F^* \sim F(2, n 5)$; under H_1 , $F^* \sim$ noncentral F(2, n 5)
- If reject, conclude either X_2 or X_4 or both contain additional useful information to predict Y in a linear model with X_1 and X_3
- Example: Consider predicting GPA with HS grades, do SAT scores add any useful information?
- If neither H₀ nor H₁ is correct, p-value has no rigorous statistical meaning, but still serves as a comparison tool.

Special Cases

• Consider testing individual predictor X_i based on

 $SSE(X_i|X_1, ..., X_{i-1}, X_{i+1}, ..., X_{p-1})$

- These are related to SAS's indiv parameter *t*-tests

 $F(1, n-p) = t^2(n-p)$

- Can decompose SSM variety of ways
 - Decomposition of $SSM(X_1, X_2, X_3)$

= SSM $(X_1) +$ SSM $(X_2|X_1) +$ SSM $(X_3|X_2,X_1)$

- = SSM (X_2) + SSM $(X_1|X_2)$ + SSM $(X_3|X_2,X_1)$
- = SSM(X₃) + SSM(X₂|X₃) + SSM(X₁|X₂, X₃)
- Stepwise sum of squares called Type I SS

Type I SS and Type II SS

- Type I and Type II are very different
 - Type I is sequential, so it depends on model statement

 Type II is conditional on all others, so it does not depend on model statement

• For example, model y = x1 x2 x3 yields

Туре І	Туре II
$SSM(X_1)$	$SSM(X_1 X_2,X_3)$
$SSM(X_2 X_1)$	$SSM(X_2 X_1,X_3)$
$SSM(X_3 X_1,X_2)$	$SSM(X_3 X_1,X_2)$

 Could variables be explaining same SS and "canceling" each other out, such that we need to cautions about testing results?

Example: Body Fat (p.256)

- Twenty healthy female subjects
- Y is body fat via underwater weighing
- Underwater weighing is expensive/difficult
- X_1 is triceps skinfold thickness
- X₂ is thigh circumference
- X₃ is midarm circumference

• Investigate the model with all three predictors:

```
data a1;
    infile 'U:\Ch07ta01.txt';
    input skinfold thigh midarm fat;
proc reg data=a1;
    model fat=skinfold thigh midarm /ss1 ss2;
run;
```

			Analysis	of	Varianc	е					
			Sum o	of	M	ear	ı				
Source		DF	Square	es	Squ	are	e F	Valu	le	Pr	> F
Model		3	396.9846	51	132.32	820		21.5	52	<.	0001
Error		16	98.4048	39	6.15	031	L				
Corrected	Total	19	495.3898	50							
Root MSE			2.47998		R-Square		0.80	014			
Dependent	Mean	2	0.19500		Adj R-Sq		0.76	541			
Coeff Var		1	2.28017								
		P	arameter	Es	timates						
		Para	meter	S	tandard						
Variable	DF	Est	imate		Error	t	Value		Pr 2	> t	
Intercept	1	117.	08469	9	9.78240		1.17		0	.257	8
skinfold	1	4.	33409		3.01551		1.44		0	.169	9
thigh	1	-2.	85685		2.58202		-1.11		0	.284	9
midarm	1	-2.	18606		1.59550		-1.37		0	.189	6

Conclusions

- Set of three variables helpful in predicting body fat (P < 0.0001)
- None of the individual parameters is significant
 - Addition of each predictor to a model containing the other two is not helpful
 - Example of multicollinearity
 - Will discuss more in next topic
- Will now focus on extra SS

• Output Using SS1 & SS2 gives an additonal two columns

Parameter Estimates

Parameter

Variable	DF	Estimate	Type I SS	Type II SS
Intercept	1	117.08469	8156.76050	8.46816
skinfold	1	4.33409	352.26980	12.70489
thigh	1	-2.85685	33.16891	7.52928
midarm	1	-2.18606	11.54590	11.54590

More than 90% of Type I SS of skinfold can also be explained by thigh and midarm

• Investigate the model via general linear tests: fat=skinfold

```
proc reg data=a1;
    model fat=skinfold;
run;
```

		Ana	lysis (of Va	arian	ce				
			Sur	n of		Mea	n			
Source		DF	Squa	ares		Squar	e F	7 Val	ue	Pr > F
Model		1	352.20	5980	352	.2698	0	44.	30	<.0001
Error		18	143.1	1970	7	.9510	9			
Corrected	Total	19	495.38	8950						
Root MSE			2.8197	7	R-Sq	uare		0.71	.11	
Dependent	Mean	2	0.1950)	Adj 1	R-Sq		0.69	50	
Coeff Var		1	3.9627	1						
		Pa	ramete	r Est	timat	es				
		Param	eter	Star	ndard					
Variable	DF	Estim	ate	E	ror	t Va	lue	Р	'r >	t
Intercept	1	-1.49	610	3.3	1923	-0	.45		0.6	6576
skinfold	1	0.85	719	0.12	2878	6	.66		<.(0001

• Skinfold now helpful. Note the change in coefficient estimate and standard error compared to the full model.

- Does this variable alone do the job?
- Perform general linear test

```
proc reg data=a1;
    model fat=skinfold thigh midarm;
    thimid: test thigh, midarm;
run; quit;
```

Test thimid Results for Dependent Variable fat

		Mean		
Source	DF	Square	F Value	Pr > F
Numerator	2	22.35741	3.64	0.0500
Denominator	16	6.15031		

• Appears there is additional information in the variables. Perhaps the addition of one more variable would be helpful.

Partial Correlations

- Measures the strength of a linear relation between two variables taking into account other variables or after adjusting for other variables, while marginal t-test measure the existence of such a linear realtionship
- Procedure for X_i vs Y
 - Predict Y using other X's
 - Predict X_i using other X's
- Each residual represents what is not explained by the other variables
 - The correlation between residuals is the partial correlation
 - The regression coefficient residuals is the regression coefficient of $X_i\,$ in the full model
 - The test of association is the marginal t test in the full model
- Looking for <u>additional</u> information in X_i that better explains Y

Example: Body Fat

```
proc reg data=a1;
```

```
model fat=skinfold thigh midarm / pcorr2;
```

run;

Parameter Estimates

					Squared
	Parameter	Standard			Partial
DF	Estimate	Error	t Value	Pr > t	Corr Type II
1	117.08469	99.78240	1.17	0.2578	•
1	4.33409	3.01551	1.44	0.1699	0.11435
1	-2.85685	2.58202	-1.11	0.2849	0.07108
1	-2.18606	1.59550	-1.37	0.1896	0.10501
	DF 1 1 1 1	Parameter DF Estimate 1 117.08469 1 4.33409 1 -2.85685 1 -2.18606	ParameterStandardDFEstimateError1117.0846999.7824014.334093.015511-2.856852.582021-2.186061.59550	ParameterStandardDFEstimateErrort Value1117.0846999.782401.1714.334093.015511.441-2.856852.58202-1.111-2.186061.59550-1.37	ParameterStandardDFEstimateErrort ValuePr > t 1117.0846999.782401.170.257814.334093.015511.440.16991-2.856852.58202-1.110.28491-2.186061.59550-1.370.1896

- Squared partial correlation is also called coefficient of partial determination. Has similar interpretation to coefficient of multiple determination.
- Squared partial correlation = Type II SS/ (Type II SS + SSE).
- In this case, variables only explain approximately 10% of the remaining variation after the other two variables are fit.

Standardized Regression Model

- Can reduce round-off errors in calculations
- Standardization

$$\tilde{Y}_i = \frac{1}{\sqrt{n-1}} \left(\frac{Y_i - \overline{Y}}{s_Y} \right) \quad \text{and} \quad \tilde{X}_{ik} = \frac{1}{\sqrt{n-1}} \left(\frac{X_{ik} - \overline{X}_i}{s_{X_i}} \right)$$

- Puts regression coefficients in common units, such that they can be compared fairly
- A one SD change in X_i corresponds to $\tilde{\beta}_i$ SD increase in Y
- Can show

$$\beta_i = \left(\frac{s_Y}{s_{X_i}}\right) \tilde{\beta}_i$$

Example: Body Fat

```
proc reg data=a1;
```

```
model fat=skinfold thigh midarm / stb;
```

run;

Parameter Estimates

		Parameter	Standard			Standardized
Variable	DF	Estimate	Error	t Value	Pr > t	Estimate
Intercept	1	117.08469	99.78240	1.17	0.2578	0
skinfold	1	4.33409	3.01551	1.44	0.1699	4.26370
thigh	1	-2.85685	2.58202	-1.11	0.2849	-2.92870
midarm	1	-2.18606	1.59550	-1.37	0.1896	-1.56142

**Skinfold has highest standardized coefficient. Midarm does not appear to be as important a predictor. Perhaps best model includes skinfold and thigh.

Multicollinearity

- Numerical analysis problem is that the matrix $\mathbf{X}'\mathbf{X}$ is almost singular (linear dependent columns)
 - Makes it difficult to take the inverse
 - Generally handled with current algorithms
- Statistical problem: too much correlation among predictors
 - The coefficient estimation lacks interpretability.
 - Difficult to determine regression coefficients \longrightarrow Increased standard error
 - May not affect prediction accuracy if the testing samples follow similar multicollinear correlation.
- Want to refine model to remove redundancy in the predictors

Example

• Consider a two-predictor model

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i$$

- What is the estimate of β_1 ?
- Can show

$$b_1 = \frac{\tilde{b}_1 - \sqrt{\frac{s_Y^2}{s_{X_1}^2}} r_{12} r_{Y2}}{1 - r_{12}^2}$$

where \tilde{b}_1 is the estimate fitting Y vs X_1

Extreme Cases

- Consider X_1 and X_2 are uncorrelated
 - $r_{12} = 0$
 - $b_1 = \tilde{b}_1$ (fitting Y vs X_1)
 - Estimator b_1 does not depend on X_2
 - Type I SS and Type II SS are the same
 - In other words, the contribution of each predictor is the same regardless of whether or not the other predictor is in the model
- Consider $X_1 = a + bX_2$

 $- r_{12} = \pm 1$

- Estimator b_1 does not exist
- Type II SS are zero
- In other words, there is no contribution of the predictor if the other predictor is already in the model

Extreme Case in SAS

• Consider the following data set

```
data a1;
    input case x1 x2 y;
    cards;
    1 3 3 5
    2 4 5 8
    3 1 -1 7
    4 6 9 15
;
```

- Notice $x_2 = 2x_1 3$
- Will generate 3-D plot and run regression

```
/* Generate 3-D Scatterplot */
proc g3d data=a1;
    scatter x2*x1=y / rotate=30;
run;
```



7-21

```
proc reg data=a1;
    model y=x2 x1;
run; quit;
```

		Analysis	of Variance		
		Sum of	Mean		
Source	DF	Squares	Square	F Value	Pr > F
Model	1	55.59211	55.59211	96.02	0.0103
Error	2	1.15789	0.57895		
Corrected Total	3	56.75000			
Root MSE		0.76089	R-Square	0.97	96
Dependent Mean		8.75000	Adj R-Sq	0.96	94
Coeff Var		8.69584			

- NOTE: Model is not full rank. Least-squares solutions for the parameters are not unique. Some statistics will be misleading. A reported DF of 0 or B means that the estimate is biased.
- NOTE: The following parameters have been set to 0, since the variables are a linear combination of other variables as shown.

x1 = 1.5 * Intercept + 0.5 * x2

		Paramete	er Estimate	S	
		Parameter	Standard		
Variable	DF	Estimate	Error	t Value	Pr > t
Intercept	В	-0.65789	1.03271	-0.64	0.5893
x2	В	1.71053	0.17456	9.80	0.0103
x1	0	0	•	•	•

- In this example, no inverse exists so X_1 dropped
- In practice, we are concerned with less extremal cases
- General results still hold
 - Regression coefficients are not well estimated
 - Regression coefficients may be scientifically meaningless
 - Type I SS and II SS will differ substantially
 - R^2 and predicted values are usually ok

Prelim Diagnose: Pairwise Correlations

- Assesses "pairwise collinearity" but not complicated multicollinearity
- Consider our body fat example

```
proc reg data=a1 corr;
    var skinfold thigh midarm fat;
    model midarm = skinfold thigh;
run; quit;
```

Correlation

Variable	skinfold	thigh	midarm	fat
skinfold	1.0000	0.9238	0.4578	0.8433
thigh	0.9238	1.0000	0.0847	0.8781
midarm	0.4578	0.0847	1.0000	0.1424
fat	0.8433	0.8781	0.1424	1.0000

- relatively strong correlation between thigh and skinfold.

• "MODEL midarm = skinfold thigh" reported $R^2 = 0.9904$

- All three
$$\rightarrow r = \sqrt{0.9904} = .995$$

- Should not use model with all three predictors

Coefficient Estimation

• Page 284 summarizes coefficients

Variables in Model	b_1	b_2
skinfold	0.8572	-
thigh	-	0.8565
skinfold, thigh	0.2224	0.6594
skinfold, thigh, midarm	4.3340	-2.857

- skinfold and thigh similar info, hence are exchangeable.
- Coeffs change when both are included (sum ≈ 0.86)
- Very dramatic change when midarm is in (but this change is still dominated by the multicollinearity)
- Reflected in std errors too

Chapter Review

- Extra Sums of Squares
- Partial correlations
- Standardized regression coefficients
- Multicollinearity
 - Effects
 - Remedies