Topic 4 - Analysis of Variance
 Outline

 Approach to Regression
 • Partitioning sums of squares

 STAT 525 - Fall 2013
 • Degrees of freedom

 • Expected mean squares
 • General linear test

 • R<sup>2</sup> and the coefficient of correlation
 • What if X random variable?

STAT 525

# **Partitioning Sums of Squares**

- Organizes results arithmetically
- Total sums of squares in Y is defined

$$SSTO = \sum (Y_i - \overline{Y})^2$$

- Can partition sum of squares into
  - Model (explained by regression)
  - Error (unexplained / residual)
- Rewrite the total sum of squares as

$$\sum (Y_i - \overline{Y})^2 = \sum (Y_i - \hat{Y}_i + \hat{Y}_i - \overline{Y})^2$$
$$= \sum (\hat{Y}_i - \overline{Y})^2 + \sum (Y_i - \hat{Y}_i)^2$$

STAT 525

# **Total Sum of Squares**

• If we ignored  $X_h$ , the sample mean  $\overline{Y}$  would be the best linear unbiased predictor

$$Y_i = \beta_0 + \varepsilon_i = \mu + \varepsilon_i$$

- SSTO is the sum of squared deviations for this predictor
- Sum of squares has n-1 degrees of freedom because we replace  $\beta_0$  with  $\overline{Y}$
- The total mean square is SSTO/(n-1) and represents an unbiased estimate of  $\sigma^2$  under the above model

Topic 4

STAT 525

# **Regression Sum of Squares**

• SAS calls this *model* sum of squares

$$SSR = \sum (\hat{Y}_i - \overline{Y})^2$$

- Degrees of freedom is 1 due to the addition of the slope
- SSR large when  $\hat{Y}_i$ 's are different from  $\overline{Y}$
- This occurs when there is a linear trend
- Under regression model, can also express SSR as

$$SSR = b_1^2 \sum \left( X_i - \overline{X} \right)^2$$

STAT 525

Topic 4

## Error Sum of Squares

SAS & Total Sum of Squares

• "Corrected" means that the sample mean has been

• SAS uses "Corrected Total" for SSTO

subtracted off before squaring

• Uncorrected total sum of squares is  $\sum Y_i^2$ 

• Error (or residual) sum of squares is equal to the sum of squared residuals

$$SSE = \sum (Y_i - \hat{Y}_i)^2 = \sum e_i^2$$

- Degrees of freedom is n 2 due to using (b<sub>0</sub>, b<sub>1</sub>) in place of (β<sub>0</sub>, β<sub>1</sub>)
- SSE large when |residuals| are large. This implies  $Y_i$ 's vary substantially around fitted line
- The MSE=SSE/(n-2) and represents an unbiased estimate of  $\sigma^2$  when taking X into account

STAT 525

Topic 4

#### ANOVA Table

• Table puts this all together

Source ofVariationdfSSMSRegression1 $b_1^2 \sum (X_i - \overline{X})^2$ SSR/1(Model)Errorn-2 $\sum (Y_i - \hat{Y})^2$ SSE/(n-2)Totaln-1 $\sum (Y_i - \overline{Y})^2$ 

Topic 4

#### **Expected Mean Squares** F test • All means squares are random variables • Can use this structure to test $H_0: \beta_1 = 0$ • Already showed $E(MSE) = \sigma^2$ • Consider $F^{\star} = \frac{\text{MSR}}{\text{MSE}}$ • What about the MSR? • If $\beta_1 = 0$ then $F^*$ should be near one $E(MSR) = E(b_1^2 \sum (X_i - \overline{X})^2)$ $= E(b_1^2) \sum (X_i - \overline{X})^2$ • Need sampling distribution of $F^*$ under $H_0$ $= (\operatorname{Var}(b_1) + E(b_1)^2) \sum (X_i - \overline{X})^2$ • By Cochran's Thm (pg 70) $= \sigma^2 + \beta_1^2 \sum (X_i - \overline{X})^2$ $F^{\star} = \frac{\underline{SSR}}{1} \div \frac{\underline{SSE}}{n-2}$ $F^{\star} \sim \frac{\chi_1^2}{1} \div \frac{\chi_{n-2}^2}{n-2}$ • If $\beta_1 = 0$ , MSR unbiased estimate of $\sigma^2$ $\sim F_{1,n-2}$ Topic 4 Topic 4STAT 525 STAT 525 F test Example • When $H_0$ is false, MSR > MSE data a1; infile 'C:\Textdata\CH01TA01.txt'; • P-value = $\Pr(F(1, n-2) > F^{\star})$ input size hours; • Reject when $F^*$ large, P-value small proc reg data=a1; • Recall t-test for $H_0: \beta_1 = 0$ model hours=size; • Can show $t_{n-2}^2 \sim F_{1,n-2}$ id size; • Obtain exactly the same result (P-value) run;

STAT 525

Topic 4

11

STAT	525
JIMI	040

	De	epend	ent Varia	ble: ho	urs			
		Anal	ysis of V	Variance				
			Sum of	1	Mean			
	Source	DF	Squares	Sq	uare F Va	lue	Pr > F	
	Model	1	252378	25	2378 105	5.88	<.0001	
	Error	23	54825	2383.7	1562			
	Cor Total	24	307203					
	Root MSE		48.82	2331	R-Square	0	.8215	
	Dependent	Mean	312.28	. 000	Adj R-Sq	0	.8138	
	Coeff Var		15.63	3447				
			Paramete	er Estim	ates			
			Parameter	Standa	rd			
	Variable	DF	Estimate	Err	or t Valu	le P	°r >  t	
	Intercept	1	62.36586	26.177	43 2.3	88	0.0259	
	size	1	3.57020	0.346	97 10.2	29	<.0001	
opic 4								13

#### General Linear Test

- Reduced model  $\longrightarrow H_0: \beta_1 = 0$
- Can be shown that  $SSE(F) \leq SSE(R)$
- Idea: more parameters provide better fit
- If SSE(F) not much smaller than SSE(R), full model doesn't better explain Y

$$F^{\star} = \frac{(\text{SSE}(\text{R}) - \text{SSE}(\text{F}))/(df_R - df_F)}{\text{SSE}(\text{F})/df_F}$$

$$= \frac{(\text{SSTO} - \text{SSE})/1}{\text{SSE}/(n-2)}$$

• Same test as before but more general

STAT	525

# General Linear Test A different approach to the same problem Consider two models Full model : Y<sub>i</sub> = β<sub>0</sub> + β<sub>1</sub>X<sub>i</sub> + ε<sub>i</sub> Reduced model : Y<sub>i</sub> = β<sub>0</sub> + ε<sub>i</sub> Will compare models using SSE's Full model will be labeled SSE(F) Reduced model will be labeled SSE(R) Note: SSTO is the same under each model

STAT 525

Topic 4

# **Pearson Correlation**

• Number between -1 and 1 which measures the strength of the <u>linear</u> relationship between two variables

$$r = \frac{\sum (X_i - \overline{X})(Y_i - \overline{Y})}{\sqrt{\sum (X_i - \overline{X})^2 \sum (Y_i - \overline{Y})^2}}$$
$$= b_1 \sqrt{\frac{\sum (X_i - \overline{X})^2}{\sum (Y_i - \overline{Y})^2}}$$

• Test  $H_0: \beta_1 = 0$  similar to  $H_0: \rho = 0$ 

Topic 4

15

# **Coefficient of Determination**

• Defined as the proportion of total variation explained by the model utilizing X

$$R^2 = \frac{\text{SSR}}{\text{SSTO}} = 1 - \frac{\text{SSE}}{\text{SSTO}}$$

 $\bullet \ 0 \le R^2 \le 1$ 

• Often multiplied by 100 and described as a percent

STAT 525

# **Coefficient of Determination**

• Can show this is equal to  $r^2$ 

$$r^{2} = b_{1}^{2} \left( \frac{\sum (X_{i} - \overline{X})^{2}}{\sum (Y_{i} - \overline{Y})^{2}} \right)$$
$$= \frac{b_{1}^{2} \sum (X_{i} - \overline{X})^{2}}{\sum (Y_{i} - \overline{Y})^{2}}$$
$$= \frac{\text{SSR}}{\text{SSTO}}$$

- Relationship not true in multiple regression
- See page 75 for limitations of  $R^2/r$

Topic 4

#### STAT 525

## Normal Correlation Model

- So far, have assumed  $X_i$ 's are known constants
- In inference we've considered repeat sampling of error terms with the X<sub>i</sub>'s remaining fixed (Y<sub>i</sub>'s vary)
- What if this assumption is not appropriate?
- In other words, what if  $X_i$ 's are random?
- If interest still in relation between these two variables can use correlation model
- Normal correlation model assumes a bivariate normal distribution

#### STAT 525

Topic 4

17

## **Bivariate Normal Distribution**

- Consider random variables  $Y_1$  and  $Y_2$
- Distribution requires five parameters
  - $\mu_1$  and  $\sigma_1$  are the mean and std dev of  $Y_1$
  - $-\mu_2$  and  $\sigma_2$  are the mean and std dev of  $Y_2$
  - $-\rho_{12}$  is the coefficient of correlation
- Bivariate normal density and marginal distributions given on page 79
- Marginal distributions are normal
- Conditional distributions are also normal

19

Topic 4

- Consider the distribution of  $Y_1$  given  $Y_2$
- 1 Can show the distribution is normal
- 2 The mean can be expressed

$$\left(\mu_1 - \mu_2 \rho_{12} \frac{\sigma_1}{\sigma_2}\right) + \rho_{12} \frac{\sigma_1}{\sigma_2} Y_2 = \alpha_{1|2} + \beta_{12} Y_2$$

- 3 With constant variance  $\sigma_1^2 (1 \rho_{12}^2)$
- Similar properties of normal error regression model
- Can use regression to make inference about  $Y_1$  given  $Y_2$

#### So What if X is Random?

- Suppose  $X_i$ 's are random samples from  $g(X_i)$ ?
- Then the previous regression results hold if:
  - The conditional distributions of  $Y_i$  given  $X_i$  are normal and independent with conditional means  $\beta_0 + \beta_1 X_i$ and conditional variance  $\sigma^2$
  - The  $X_i$  are independent and  $g(X_i)$  does not involve the parameters  $\beta_0$ ,  $\beta_1$ , and  $\sigma^2$

Topic 4

STAT 525

#### Inference on $\rho_{12}$

- Point estimate using  $Y = Y_1$  and  $X = Y_2$  given on 4-15
- Interest in testing  $H_0: \rho_{12} = 0$
- Test statistic is

$$t^* = \frac{r_{12}\sqrt{n-2}}{\sqrt{1-r_{12}^2}}$$

- Same result as  $H_0: \beta = 0$
- Can also form CI using Fisher z transformation or large sample approx (pg 85)
- If X and Y are nonnormal, can use Spearman correlation (pg 87)

STAT 525

Topic 4

21

# Background Reading

- Appendix A
- KNNL Chapter 3

Topic 4

23