

Analysis of Variance in Regression

$$(y_i - \overline{y}) = (y_i - \overline{y}_i) + (\overline{y}_i - \overline{y})$$

$$\sum (y_{i} - \overline{y})^{2} = \sum (y_{i} - y_{i})^{2} + \sum (y_{i} - \overline{y})^{2}$$

Total (*SST*) • $df_{\text{Total}} = n-1$ Error (SSE)Regression (SSR) $df_{\rm Error} = n-2$ $df_{\rm Regression} = 1$

Analysis of Variance in Regression

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	F
Model	SSR	1	MSR = SSR/1	F = MSR/MSE
Error	SSE	<i>n</i> -2	MSE = SSE/(n-2)	
Total	$SST=S_{yy}$	<i>n</i> -1		

• Analysis of Variance - F-test

•
$$H_0: \beta_1 = 0$$
 $H_A: \beta_1 \neq 0$
 $F_s = \frac{MSR}{MSE}$
 $RR: F_s \ge F_{\alpha,1,n-2}$
 $Pvalue: P(F \ge F_s)$

Multiple Regression

- Numeric Response variable (Y)
- *p* Numeric predictor variables
- Model:

 $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon$

Partial Regression Coefficients: β_i = effect on the mean response of increasing the *i*th predictor variable by 1 unit, holding all other predictors constant

Analysis of Variance

• Only changes from SLR are the d.f. $- df_{Model} = p \qquad df_{Error} = n-p-1$

Source of	Sum of	Degrees of	Mean	
Variation	Squares	Freedom	Square	F
Model	SSR	р	MSR = SSR/p	F = MSR/MSE
Error	SSE	<i>n-p-</i> 1	MSE = SSE/(n-p-1)	
Total	$SST=S_{yy}$	<i>n</i> -1		

$$R^{2} = \frac{SST - SSE}{SST} = \frac{SSR}{SST}$$

Testing the Overall Model

- Tests whether **any** of the explanatory variables are associated with the response
- $H_0: \beta_1 = \dots = \beta_p = 0$
- $H_{\rm A}$: Not all $\beta_{\rm i} = 0$

 $F_{s} = \frac{MSR}{MSE}$ $RR: \quad F_{s} \ge F_{\alpha,p,n-p-1}$ $Pvalue: P(F \ge F_{s})$

Analysis of Covariance

- Combination of 1-Way ANOVA and Linear Regression
- Goal: Comparing numeric responses among *k* groups, adjusting for numeric concomitant variable(s), referred to as **Covariate(s)**
- Clinical trial applications: Response is Post-Trt score, covariate is Pre-Trt score
- Epidemiological applications: Outcomes compared across exposure conditions (i.e., genotype), adjusted for other risk factors (age, BMI, sex,...)

Models with Dummy Variables

- Since <u>genotype</u> is a categorical variable, we need to recode its values in order to include it in the regression model.
- If a categorical variable has *k* levels, need to create *k*-1 dummy variables that take on the values 1 if the level of interest is present, 0 otherwise.
- The baseline level of the categorical variable for which all *k*-1 dummy variables are set to 0
- The regression coefficient corresponding to a dummy variable is the difference between the mean for that level and the mean for baseline group, controlling for all numeric predictors