

# Regression Analysis

## Multiple Linear Regression

**Nicoleta Serban, Ph.D.**

*Professor*

School of Industrial and Systems Engineering

Testing for Subsets of  
Regression Parameters



1

## About This Lesson



2

# Testing Overall Regression

**Analysis of Variance (ANOVA) for multiple regression:**

Variability Source	DF	Sum of Squares	Mean SS	F-Statistic
Regression	$p$	SSReg	SSReg / $p$	MSSReg / MSE
Residual	$n-p-1$	SSE	SSE / $(n-p-1)$	
Total	$n-1$	SST		

$$\text{SSReg} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad \text{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{SST} = \sum_{i=1}^n (y_i - \bar{y})^2$$

Null hypothesis: All predictor coefficients are 0, i.e.,  $\mathbf{H}_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$ .

Reject  $\mathbf{H}_0$  if F-statistic is large ( $> F_{\alpha, p, n-p-1}$  for  $\alpha$  significance level,  $p$  and  $n-p-1$  df).

- At least one of the coefficients is different from zero at the  $\alpha$  significance level.

p-value = Prob( $F_{p, n-p-1} > \text{F-statistic}$ ) for F-distribution with  $p$  and  $n-p-1$  df.

- Reject  $\mathbf{H}_0$  if p-value is small.



3

# Testing Subsets of Coefficients

**Analysis of Variance (ANOVA):**

$$\text{SST}(X_1, \dots, X_p) = \text{SSReg}(X_1, \dots, X_p) + \text{SSE}(X_1, \dots, X_p)$$

$$\text{SSReg}(X_1, \dots, X_p) = \text{SSReg}(X_1) + \text{SSReg}(X_2|X_1) + \text{SSReg}(X_3|X_1, X_2) + \dots + \text{SSReg}(X_p|X_1, \dots, X_{p-1})$$

**SSReg( $X_1$ ):** Sum of squares (SS) explained using only  $X_1$

**SSReg( $X_2|X_1$ ):** **Extra** SS explained using  $X_2$  in addition to  $X_1$

**SSReg( $X_3|X_1, X_2$ ):** **Extra** SS explained using  $X_3$  in addition to  $X_1$  and  $X_2$

**SSReg( $X_p|X_1, \dots, X_{p-1}$ ):** **Extra** SS explained using  $X_p$  in addition to  $X_1, X_2, \dots, X_{p-1}$



4

# Testing Subsets of Coefficients

- Does  $X_1$  alone significantly aid in predicting  $Y$ ?
  - $SSReg(X_1)$  vs.  $SSE(X_1)$
- Does the addition of  $X_2$  significantly contribute to the prediction of  $Y$  after accounting (controlling) for the contribution of  $X_1$ ?
  - $SSReg(X_2 | X_1)$  vs.  $SSE(X_1, X_2)$
- Does the addition of  $X_3$  significantly contribute to the prediction of  $Y$  after accounting (controlling) for the contribution of  $X_1$  and  $X_2$ ?
  - $SSReg(X_3 | X_1, X_2)$  vs.  $SSE(X_1, X_2, X_3)$
- Does the addition of  $X_p$  significantly contribute to the prediction of  $Y$  after accounting (controlling) for the contribution of  $X_1, \dots, X_{p-1}$ ?
  - $SSReg(X_p | X_1, \dots, X_{p-1})$  vs.  $SSE(X_1, X_2, \dots, X_p)$

# Testing Subsets of Coefficients

## Partial F-test:

- Consider a full model with two sets of predictors,  $X_1, \dots, X_p$  (perhaps controlling factors) and  $(Z_1, \dots, Z_q)$  (perhaps additional explanatory factors):

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \alpha_1 Z_1 + \dots + \alpha_q Z_q + \varepsilon$$

- Test whether any of the  $Z$  factors add explanatory power to the model:

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_q = 0 \quad \text{vs.} \quad H_a: \alpha_i \neq 0 \text{ for at least one } \alpha_i, i = 1, \dots, q$$

$$F\text{-statistic} = F_{\text{partial}} = \frac{SSReg(Z_1, \dots, Z_q | X_1, \dots, X_p) / q}{SSE(Z_1, \dots, Z_q, X_1, \dots, X_p) / (n - p - q - 1)}$$

- Reject  $H_0$  if F-statistic is large ( $F\text{-statistic} > F_{\alpha, q, n-p-q-1}$ )
  - At least one coefficient is different from zero at the  $\alpha$  significance level

# Testing for Statistical Significance

- Consider a full model with the set of predictors,  $X_1, \dots, X_p$  and an additional predicting variable Z:

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \alpha Z + \varepsilon$$

- Test whether Z has explanatory or predictive power:

$$H_0: \alpha = 0 \text{ vs } H_a: \alpha \neq 0$$

$$F\text{-statistic} = F_{\text{partial}} = \frac{SS\text{Reg}(Z|X_1, \dots, X_p)/1}{SSE(Z, X_1, \dots, X_p)/(n-p-2)}$$

- Reject  $H_0$  if F-statistic is large ( $F\text{-statistic} > F_{\alpha, 1, n-p-2}$ )

**This is equivalent to testing for statistical significance using the t-test**



7

# Testing for Statistical Significance

- Consider a full model with the set of predictors,  $X_1, \dots, X_p$  and an additional predicting variable Z:

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \alpha Z + \varepsilon$$

- Test whether Z has explanatory or predictive power:

$$H_0: \alpha = 0 \text{ vs } H_a: \alpha \neq 0$$

$$F\text{-statistic} = F_{\text{partial}} = \frac{SS\text{Reg}(Z|X_1, \dots, X_p)/1}{SSE(Z, X_1, \dots, X_p)/(n-p-2)}$$

- Reject  $H_0$  if F-statistic is large ( $F\text{-statistic} > F_{\alpha, 1, n-p-2}$ )

**This is equivalent to testing for statistical significance using the t-test**

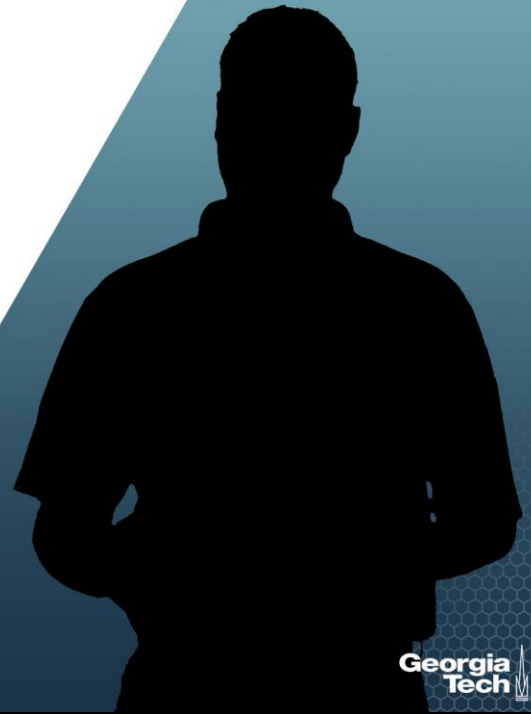
- Interpretation of the t-test for statistical significance is conditional on other predicting variables being in the model.
- The relationship between Y and X is statistically significant given all other predicting variables being in the model.

**Do not perform variable selection based on the p-values of the t-tests!**



8

# Summary



Georgia  
Tech