# Regression Analysis
## Analysis of Variance

**Nicoleta Serban, Ph.D.**
*Professor*
School of Industrial and Systems Engineering

Parameter Estimation

Georgia Tech

1

# About This Lesson

Georgia Tech

2

# ANOVA: Model & Assumptions

**Data**: $Y_{ij}$ for $j = 1, \cdots, n_i; i = 1, \cdots, k$

**Model**: $Y_{ij} = \mu_i + \varepsilon_{ij}$ where $\varepsilon_{ij} = $ error term

**Assumptions**:

- ***Constant Variance Assumption:*** $\mathrm{Var}(\varepsilon_{ij}) = \sigma^2$

- ***Independence Assumption:*** $\{\varepsilon_{1j}, \cdots, \varepsilon_{kj}\}$ are independent random variables

- ***Normality Assumption:*** $\varepsilon_{ij} \sim$ Normal $(0, \sigma^2)$

**Georgia Tech**

3

# ANOVA: Variance Estimation

Comparing means from multiple populations assuming the variances are the same and equal to $\sigma^2$:

*Pooled Variance Estimator:*

$$S_{pool}^2 = \frac{\sum_{i=1}^{k}(n_i - 1)S_i^2}{\sum_{i=1}^{k}(n_i - 1)} = \frac{\sum_{i=1}^{k}\sum_{j=1}^{n_i}\left(Y_{ij} - \bar{Y}_i\right)^2}{N - k}$$

Where N = total number of samples = $(n_1 + \ldots + n_k)$

The degrees of freedom is N-k because we replace $\mu_i \leftarrow \bar{Y}_i$ for i = 1,…,k, thus losing k degrees of freedom

**Georgia Tech**

4

# ANOVA: Variance Estimation (cont'd)

$$\widehat{\sigma}^2 = \frac{\sum_{i=1}^{k}(n_i-1)S_i^2}{\sum_{i=1}^{k}(n_i-1)} = \frac{\sum_{i=1}^{k}\sum_{j=1}^{n_i}\left(Y_{i_j}-\bar{Y}_i\right)^2}{N-k} = \underline{\textbf{MSE}}$$

$$\sum_{i=1}^{k}\sum_{j=1}^{n_i}\left(Y_{i_j}-\bar{Y}_i\right)^2 = \underline{\textbf{S}}\text{um of }\underline{\textbf{S}}\text{quares of }\underline{\textbf{E}}\text{rror} = \underline{\textbf{SSE}}$$

We will use interchangeably Sum of Squared Errors and Sum of Squared Residuals.

**Georgia Tech**

5

# Mean Squared Error (MSE)

$S_1^2,...,S_k^2$     *The sum of independent chi-square random variables is also chi-square*

$$\frac{SSE}{\sigma^2} = \frac{(n_1-1)S_1^2}{\sigma^2} + \cdots + \frac{(n_k-1)S_k^2}{\sigma^2} \sim \chi_\nu^2 \text{ where } \nu = N-k$$

The sampling distribution of the pooled variance is a chi-square distribution with N-k degrees of freedom.

**Georgia Tech**

6

# Estimating Parameters in ANOVA

$$\hat{\mu}_i = \bar{Y}_i = \frac{1}{n_i}\sum_{j=1}^{n_i} Y_{ij}$$

What is the sampling distribution?

If $Y_{i1},...,Y_{in} \sim N(\mu_i, \sigma^2)$ ➡ $\hat{\mu}_i = \bar{Y}_i = \frac{Y_{i1}+...+Y_{in}}{n_i} \sim N(\mu_i, \sigma^2/n_i)$

But $\sigma^2$ is unknown.
So replace $\sigma^2$ with the pooled variance estimation: $\sigma^2 \longleftarrow$ MSE

$$\frac{\hat{\mu}_i - \mu_i}{\sqrt{MSE/n_i}} \sim t_{N-k}$$

Why $N - k$?

$$MSE = \hat{\sigma}^2 \sim \chi^2_{N-k}$$

**Georgia Tech**

# Confidence Intervals for the Means

We can use the estimated sample means

$$\hat{\mu}_i = \bar{Y}_i = \frac{1}{n_i}\sum_{j=1}^{n^i} Y_{ij} \text{ for } i = 1, \cdots, k$$
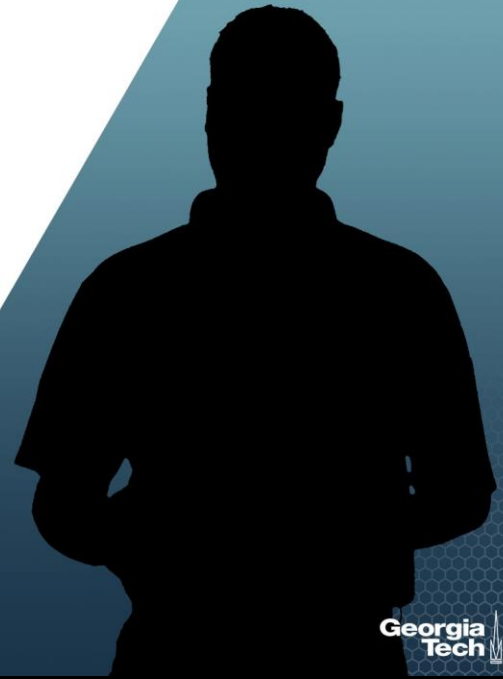
and the estimated variance

$$\hat{\sigma}^2 = MSE$$

to calculate $(1 - \alpha)$ confidence intervals for the treatment means:

$$\left(\hat{\mu}_i - t_{\alpha/2,\, N-k}\sqrt{MSE/n_i},\, \hat{\mu}_i + t_{\alpha/2,\, N-k}\sqrt{MSE/n_i}\right)$$

**Georgia Tech**

# Summary