**Hypothesis of Interest: Example One**

Now let's look at an example of formulating the hypothesis of interest. Suppose, in the school data set problem you are interested in testing whether the average **Reading3** score for female students at Cottonwood and Dogwood is equal to the average **Reading3** score for female students at Maple and Pine. How do you proceed with the formulation of such a hypothesis?

You begin by writing the hypothesis of interest in terms of the cell means. This provides you with the coefficients for the LSMESTIMATE statement. An easy way to begin with a two-way ANOVA is to make a table that lists the two factors and their levels in the order in which SAS reads them. The order for the factors is determined by the order in which they are entered into the CLASS statement. Generally, the first factor should be used as the row variable and the second factor as the column variable. The order of the factor levels is alphanumeric.

The body of the table represents the specific levels of the **School*Gender** interaction term. For example, $\mu_{11}$ represents the average **Reading3** score for the first row and first column, which is the average **Reading3** scores for female students of Cottonwood school. $\mu_{12}$ represents the average **Reading3** score for male students of Cottonwood school, and so on.

The entries of the column on the right side of the table are the averages of the rows and represent the main effects of the schools. For example, $\mu_{1.}$ represents the average of the first row across all of the columns, that is, the average **Reading3** scores of the Cottonwood students across both male and female genders. The entries of the row at the bottom are the averages of the columns and represent the main effect of the genders. For example, $\mu_{.1}$ represents the average of column one across all the rows, that is, the average **Reading3** scores of females across all four schools.

The hypothesis in question, that is, whether the average **Reading3** score for female students at Cottonwood and Dogwood is the same as the average **Reading3** score for female students at Maple and Pine, can be formulated mathematically in terms of cell means as $1/2(\mu_{11} + \mu_{21}) = 1/2(\mu_{31} + \mu_{41})$. In this equation, $\mu_{11}$ represents the average **Reading3** score for female students at Cottonwood, $\mu_{21}$ represents the average **Reading3** score for female students at Dogwood, $\mu_{31}$ represents the average **Reading3** score for female students at Maple, and $\mu_{41}$ represents the average **Reading3** score for female students at Pine.

The hypothesis is simplified and rewritten with zero on one side of the equation, with the levels of each factor in the correct order. The coefficients are then available and the LSMESTIMATE statement can be included in your program. The fractional coefficients can be written as decimals if they are decimals that do not repeat. In this example, 1/2 can be written as 0.5. However, accuracy and precision would be lost by writing 1/3 as 0.33. In such situations, you can multiply all coefficients by the common denominator to clear the fractions. The zeros represent the positions for the male students who are not involved in the comparison.

Another approach to writing the LSMESTIMATE statement is to begin with the two-way table based on the variables listed in the CLASS statement. Then you proceed by writing the hypothesis of interest in terms of coefficients for the corresponding cell means. Also note that, in this problem, the hypothesis refers to female gender and all schools, and makes no mention of male gender. Hence, in the body of the table, you fill in zero for any combination of school by gender not involved in the hypothesis of interest. For hypotheses on interaction terms, the LSMESTIMATE statement does not require coefficients for the main effects. For this reason, you are not obliged to fill in the margins of the table for this type of hypothesis. This approach produces the same coefficient results as the previous approach.

[Close]