

## Post Hoc Tests

### Familywise Error

Also known as alpha inflation or cumulative Type I error. Familywise error (FWE) represents the probability that any one of a set of comparisons or significance tests is a Type I error. As more tests are conducted, the likelihood that one or more are significant just due to chance (Type I error) increases. One can estimate familywise error with the following formula:

$$\alpha_{FWE} \leq 1 - (1 - \alpha_{EC})^c$$

where  $\alpha_{FWE}$  is the familywise error rate,  $\alpha_{EC}$  is the alpha rate for an individual test (almost always considered to be .05), and  $c$  is the number of comparisons.  $c$  as used in the formula is an exponent, so the parenthetical value is raised to the  $c^{\text{th}}$  power.

### Bonferroni

The Bonferroni simply calculates a new pairwise alpha to keep the familywise alpha value at .05 (or another specified value). The formula for doing this is as follows:

$$\alpha_B = \frac{\alpha_{FWE}}{c}$$

where  $\alpha_B$  is the new alpha based on the Bonferroni test that should be used to evaluate each comparison or significance test,  $\alpha_{FWE}$  is the familywise error rate as computed in the first formula, and  $c$  is the number of comparisons (statistical tests).

The Bonferroni is probably the most commonly used post hoc test, because it is highly flexible, very simple to compute, and can be used with any type of statistical test (e.g., correlations)—not just post hoc tests with ANOVA. The traditional Bonferroni, however, tends to lack power. The loss of power occurs for several reasons: (1) the familywise error calculation depends on the assumption that, for all tests, the null hypothesis is true. This is unlikely to be the case, especially after a significant omnibus test; (2) all tests are assumed to be *orthogonal* (i.e., independent or nonoverlapping) when calculating the familywise error test, and this is usually not the case when all pairwise comparisons are made; (3) the test does not take into account whether the findings are consistent with theory and past research. If consistent with previous findings and theory, an individual result should be less likely to be a Type I error; and (4) Type II error rates are too high for individual tests. In other words, the Bonferroni overcorrects for Type I error.

### Modified Bonferroni Approaches

Several alternatives to the traditional Bonferroni have been developed, including those developed by Holm, Holland and Copenhaver, Hommel, Rom, and others (see Olejnik, Li, Supattathum, & Huberty, 1997 for a review). These tests have greater power than the Bonferroni while retaining its flexible approach that allows for use with any set of statistical tests (e.g., t-tests, correlations, chi-squares).

*Sidak-Bonferroni.* Sidak (1967) suggested a relatively simple modification of the Bonferroni formula that would have less of an impact on statistical power but retain much of the flexibility of the Bonferroni method (Keppel & Wickens, 2004 discuss this testing approach). Instead of dividing by the number of comparisons, there is a slightly more complicated formula:

$$\alpha_{S-B} = 1 - (1 - \alpha_{FWE})^{1/c}$$

where  $\alpha_{S-B}$  is the Sidak-Bonferroni alpha level used to determine significance (something less than .05),  $\alpha_{FWE}$  is the computed familywise error according to the formula at the top of the first page, and  $c$  is the number of comparisons or statistical tests conducted in the “family.” The p-values obtained from the

computer printout must be smaller than  $\alpha_{S-B}$  to be considered significant. One can also extend this test to other statistical tests, such as correlations. In the case of correlations, one could replace  $df_A$  with the number of variables that are used in the group of correlations tests.  $c$  would represent the number of correlations in the correlation matrix. This approach is convenient and easy to do but has not received any systematic study, and it is likely that a single, simple correction will not result in the most efficient balance of Type I and Type II errors.

*Hochberg's Sequential Method.* This test uses a specific sequential method called a "step-up" approach as a more powerful alternative to the Bonferroni procedure. Sequential methods use a series of steps in the correction, depending on the result of each prior step. Contrasts are initially conducted and then ordered according to p-values (from smallest to largest in the "step-up" approach). Each step corrects for the previous number of tests rather than all the tests in the set. This test is a good, high power alternative to the other modified Bonferroni approaches as long as confidence intervals are not needed.

### Approaches for Pairwise Comparisons with ANOVA Designs

*Dunn.* Identical to the Bonferroni correction.

*Scheffe.* The Scheffe test computes a new critical value for an F test conducted when comparing two groups from the larger ANOVA (i.e., a correction for a standard t-test). The formula simply modifies the F-critical value by taking into account the number of groups being compared:  $(a - 1) F_{crit}$ . The new critical value represents the critical value for the maximum possible familywise error rate. As you might suppose, this also results in a higher than desired Type II error rate, by imposing a severe correction.

*Fisher LSD.* The Fisher LSD test stands for the Least Significant Difference test (rather than what you might have guessed). The LSD test is simply the rationale that if an omnibus test is conducted and is significant, the null hypothesis is *incorrect*. (If the omnibus test is nonsignificant, no post hoc tests are conducted.) The reasoning is based on the assumption that if the null hypothesis is incorrect, as indicated by a significant omnibus F-test, Type I errors are not really possible (or less likely), because they only occur when the null is true. So, by conducting an omnibus test first, one is screening out group differences that exist due to sampling error, and thus reducing the likelihood that a Type I error is present among the means. Fishers LSD test has been criticized for not sufficiently controlling for Type I error. Still, the Fisher LSD is sometimes found in the literature.

*Dunnet.* The Dunnet test is similar to the Tukey test (described below) but is used only if a set of comparisons are being made to one particular group. For instance, we might have several treatment groups that are compared to one control group. Since this is rarely the of interest, and the Tukey serves a much more general purpose, I recommend the Tukey test.

*Tukey a* (also known as Tukey's HSD for honest significant difference). Tukey's test calculates a new critical value that can be used to evaluate whether differences between any two pairs of means are significant. The critical value is a little different because it involves the mean difference that has to be exceeded to achieve significance. So one simply calculates one critical value and then the difference between all possible pairs of means. Each difference is then compared to the Tukey critical value. If the difference is larger than the Tukey value, the comparison is significant. The formula for the critical value is as follows:

$$\bar{d}_T = q_T \sqrt{\frac{MS_{s/A}}{n}}$$

$q_T$  is the studentized range statistic (similar to the t-critical values, but different), which one finds in a table (Table C.9 in the Myers & Well text),  $MS_{s/A}$  is the mean square error from the overall F-test, and  $n$  is the sample size for each group. *Error df* referred to in the table is the  $df_{s/A}$  used in the ANOVA test. *FWE* is the desired familywise error rate. This is the test I usually recommend, because studies show it has

greater power than the other tests under most circumstances and it is readily available in computer packages. The Tukey-Kramer test is used by SPSS when the group sizes are unequal. It is important to note that the power advantage of the Tukey test depends on the assumption that all possible pairwise comparisons are being made. Although this is usually what is desired when post hoc tests are conducted, in circumstances where not all possible comparisons are needed, other tests, such as the Dunnett or a modified Bonferroni method should be considered because they may have power advantages.

*Games-Howell.* This test is used with variances are unequal (see Unequal Variances below) and also takes into account unequal group sizes. Severely unequal variances can lead to increased Type I error, and, with smaller sample sizes, more moderate differences in group variance can lead to increases in Type I error. The Games-Howell test, which is designed for unequal variances, is based on Welch's correction to  $df$  with the  $t$ -test and uses the studentized range statistic. This test appears to do better than the Tukey HSD if variances are very unequal (or moderately so in combination with small sample size) or can be used if the sample size per cell is very small (e.g., <6).

### Comments

I have included only a subset of all the possible post hoc corrections for familywise error. And, believe it or not, familywise error correction procedures currently available in most statistical packages (only some of which I have focused on here) represent only a subset of the approaches which have been proposed and studied.

Klockars, Hancock, and McAweeney (1995) discuss many of the post hoc ANOVA procedures, some of which seem to advantages over the traditional approaches such as the Tukey currently available in statistical software packages. The alternative procedures are often classified as sequential vs. simultaneous, weighted vs. unweighted, and step-up vs. step-down, and they involve elaborate computational procedures which are inconvenient to do by hand especially for a large number of comparisons. Modified Bonferroni procedures have been designed for a broader array of statistical circumstances beyond post hoc ANOVA tests (e.g., correlations or chi-square tests). Olejnik and colleagues (1997) review the modified Bonferroni procedures and their computations. They conclude that most of the modified Bonferroni procedures have clear advantages over the traditional Bonferroni procedure, but small differences among the alternatives in the amount of power or control of Type I error. Their results suggest that Rom's (1990) procedure has the most power (not currently available in SPSS).

Other authors have reviewed post hoc tests with additional attention to unequal error variances (e.g., Kromrey & La Rocca, 1995; Seaman, Levin, & Serlin, 1991). How heterogeneous (i.e., unequal) the error variances must be in order to cause problems is difficult to discern, because their impact is greater with lower sample sizes. Unfortunately, tests such as the Levine tests for unequal variances have lower power when sample size is smaller, so they may be least likely to indicate a problem with unequal variances when it is most likely to affect Type I errors. In terms of post ANOVA tests, the Games-Howell is good if there are large differences in variances between groups.

Keppel, G., & Wickens, T.D. (2004). *Design and analysis: A researchers handbook* (4<sup>th</sup> Edition). Upper Saddle River, NJ: Pearson.

Klockars, A.J., Hancock, G.R., & McAweeney, M.J. (1995). Power of unweighted and weighted versions of simultaneous and sequential multiple-comparison procedures. *Psychological Bulletin*, 118, 300-307.

Kromrey, J.D., & La Rocca, M.A. (1995). Power and Type I error rates of new pairwise multiple comparison procedures under heterogeneous variances. *Journal of Experimental Education*, 63, 343-362.

Olejnik, S., Li, J., Supattathum, S., and Huberty, C.J. (1997). Multiple testing and statistical power with modified Bonferroni procedures. *Journal of educational and behavioral statistics*, 22, 389-406.

Seaman, M.A., Levin, J.R., & Serlin, R.C. (1991). New developments in pairwise multiple comparisons: Some powerful and practicable procedures. *Psychological Bulletin*, 110, 577-586.