

Data Analysis
Spring Semester, 2023
April 6, 2023
Lecture 18

Third Midterm on Thursday, April 27. It will focus on Chapters 8 and 9.

Chapter Eight
Inferences about More than Two Population Central Values

Context

The procedures in this chapter generalize the test of the equality of means of two independent populations. This generalization is often called the one-way layout. While this design has somewhat limited value in practice, the material in this chapter is fundamental for further generalizations. The key ideas that are first developed in the one-way analysis of variance are: the generalization of the t-test, the expected mean square calculation (which is described in Chapter 14 and is crucial for power calculations), and the introduction to multiple testing of hypotheses in Chapter 9.

Analysis of Variance Table

The results from a one-way layout are conventionally displayed in an analysis of variance table

Analysis of Variance Table
Complete Randomized Experiment

Source	Degrees of Freedom	Sum of Squares	Mean Square	F
Treatment	$I - 1$	$\sum_{i=1}^I J_i (Y_{i\cdot} - Y_{..})^2$	$SS_{Treatment} / (I - 1)$	$\frac{MS_{Treatment}}{MSE}$
Error	$n - I$	$\sum_{i=1}^I \sum_{j=1}^{J_i} (Y_{ij} - Y_{i\cdot})^2 = \sum_{i=1}^I (J_i - 1) S_i^2$	$SSE / (n - I)$	
Total	$n - 1$	$\sum_{i=1}^I \sum_{j=1}^{J_i} (Y_{ij} - Y_{..})^2$		

As in Chapters 11 and 12, the statistical estimate of the variance parameter in the model is the mean squared error. The model is $Y_{ij} = \mu + \alpha_i + \sigma_{1W}Z_{ij}$, for $i = 1, \dots, I$ (where I is the number of treatment settings), $j = 1, \dots, J_i$, and $\sum_{i=1}^I J_i \alpha_i = 0$. Then $\hat{\sigma}_{1W}^2 = MSE$.

Tests of hypotheses

The most common issue is whether the expected value of the outcome variable is the same for each setting of the treatment. The usual null hypothesis is then $H_0 : \mu_1 = \mu_2 = \dots = \mu_I$. The alternative hypothesis is $H_1 : \mu_i \neq \mu_{i'}, i \neq i'$; that is, there is at least one pair of treatment settings with unequal means. An equivalent statement using the effects model is that $H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_I = 0$. The equivalent alternative hypothesis is $H_1 : \alpha_i \neq \alpha_{i'}, i \neq i'$ for at least one pair of settings. The test statistic for this null hypothesis is $F = \frac{MS_{Treatment}}{MSE}$. The distribution of the test statistic under $H_0 : \mu_1 = \mu_2 = \dots = \mu_I$ is a central F distribution with $I - 1$ numerator and $n - I$ denominator degrees of freedom. I call this test the “overall F-test” or “global F-test.”

When the null hypothesis is true, the test statistic should be one, modulo statistical variability. When the alternative hypothesis is true, the test statistic should be greater than one modulo statistical variability. That is, the test of the null hypothesis is a right sided test.

Balanced one-way layouts

Typically, the questions that I ask on examinations used a balanced one-way layout. That is, $J_1 = J_2 = \dots = J_I = J$. This simplifies the calculations and permits natural and more complex issues discussed in Chapter Nine.

Example Past Examination Questions

A research team wishes to specify a manufacturing process so that Y , the area in a product affected by surface flaws is as small as possible. They have four levels of concentration of a chemical used to wash the product before the final manufacturing step and want to determine whether the concentration level causes a

change in $E(Y)$. They run a balanced one-way layout with 6 observations for each concentration with level 1 set at 10%, level 2 set at 15%, level 3 set at 20%, and level 4 set at 25%. They run a balanced one-way layout with 6 observations for each treatment. They observe that $y_1 = 264.5$, $y_2 = 255.9$, $y_3 = 216.2$, and $y_4 = 263.8$, where y_i is the average of the observations taken on the i th level. They also observe that $s_1^2 = 411.9$, $s_2^2 = 522.2$, $s_3^2 = 631.8$, and $s_4^2 = 521.9$, where s_i^2 is the unbiased estimate of the variance for the observations taken on the i th level.

- Complete the analysis of variance table for these results; that is, be sure to specify the degrees of freedom, sum of squares, mean square, and F-test.
- What is your conclusion? Use significance levels set to 0.10, 0.05, and 0.01. Make sure that you discuss the optimal setting of the concentration level and how you could document it.
- In examinations, I add parts asking for the decomposition into linear, quadratic, and cubic components. See Chapter Nine Problems.

Answers:

d. Analysis of Variance Table

Source	Sum of Squares	Degrees of Freedom	Mean Square	
Treatments	9467.4	3	3155.8	F=6.046
Error	10439.0	20	521.95	
Total	19906.4	23		

Answer: The critical values for the F-test are 2.38 for the 0.10 level, 3.10 for the 0.05 level, and 4.94 for the 0.01 level. Reject the null hypothesis that the mean responses of the four treatments are equal. The optimum setting is to use 20% as the concentration setting. This can be confirmed with Fisher's protected confidence intervals.

Chapter Nine Multiple Comparisons

9.2 Linear Contrasts

Hypotheses are typically stated using contrasts linear in the treatment means. For example, in a balanced one-way layout with $I = 4$, let $\lambda = a_1\mu_1 + a_2\mu_2 + a_3\mu_3 + a_4\mu_4$,

where $\mu_i = E(Y_{ij}) = E(Y_{i\bullet})$. The $a_i, i = 1, \dots, I$ are fixed constants with $\sum_{i=1}^I a_i = 0$. Here

$\lambda = a_1\mu_1 + a_2\mu_2 + a_3\mu_3 + a_4\mu_4$ is a linear combination of parameters and hence is itself a parameter. We need a statistic $\hat{\lambda}$ such that $E(\hat{\lambda}) = \lambda$.

That statistic is $\hat{\lambda} = a_1Y_{1\bullet} + a_2Y_{2\bullet} + a_3Y_{3\bullet} + a_4Y_{4\bullet}$. The expected value is

$E(\hat{\lambda}) = a_1E(Y_{1\bullet}) + a_2E(Y_{2\bullet}) + a_3E(Y_{3\bullet}) + a_4E(Y_{4\bullet}) = a_1\mu_1 + a_2\mu_2 + a_3\mu_3 + a_4\mu_4 = \lambda$. The variance

is $\text{var}(\hat{\lambda}) = a_1^2 \text{var}(Y_{1\bullet}) + a_2^2 \text{var}(Y_{2\bullet}) + a_3^2 \text{var}(Y_{3\bullet}) + a_4^2 \text{var}(Y_{4\bullet}) = \sum_{i=1}^I a_i^2 \frac{\sigma_{1W}^2}{J}$, where J is the

number of observations in each treatment. For example,

$\lambda_L = (-3)\mu_1 + (-1)\mu_2 + (1)\mu_3 + (3)\mu_4$ is a contrast that is called the “linear contrast.” A vector space description of the linear contrast would be as the inner product of two vectors $(-3, -1, 1, 3)$ and $(\mu_1, \mu_2, \mu_3, \mu_4)$

The statistic estimating the linear contrast is $\hat{\lambda}_L = (-3)Y_{1\bullet} + (-1)Y_{2\bullet} + (1)Y_{3\bullet} + (3)Y_{4\bullet}$. Then

$E(\hat{\lambda}_L) = \lambda_L$, and $\text{var}(\hat{\lambda}_L) = \sum_{i=1}^I a_i^2 \frac{\sigma_{1W}^2}{J} = [(-3)^2 + (-1)^2 + (1)^2 + (3)^2] \frac{\sigma_{1W}^2}{J} = \frac{20\sigma_{1W}^2}{J}$.

The sum of squares due to a contrast is $SS_\lambda = \frac{(\hat{\lambda})^2}{\sum_{i=1}^I \frac{a_i^2}{J}}$. For example, the sum of

squares due to the linear contrast for $I = 4$ is $SS_L = \frac{(\hat{\lambda}_L)^2}{20/J}$, with one degree of freedom.

The set of linear contrasts in I mean parameters is a vector space with dimension $I - 1$. Another contrast for a balanced one way layout with $I = 4$ is

$\lambda_Q = (1)\mu_1 + (-1)\mu_2 + (-1)\mu_3 + (1)\mu_4$, which is called the “quadratic contrast.” Of course, the quadratic contrast is the inner product of the vector $(1, -1, -1, 1)$ and $(\mu_1, \mu_2, \mu_3, \mu_4)$

. The inner product of the vector of coefficients of the linear contrast $(-3, -1, 1, 3)$ and the vector of coefficients of the quadratic contrast $(1, -1, -1, 1)$ is 0. Considered

as vectors, these two contrasts are orthogonal. The statistic estimating the quadratic contrast is $\hat{\lambda}_Q = (1)Y_{1\bullet} + (-1)Y_{2\bullet} + (-1)Y_{3\bullet} + (1)Y_{4\bullet}$, and

$$SS_Q = \frac{(\hat{\lambda}_Q)^2}{[1^2 + (-1)^2 + (-1)^2 + (1)^2]/J} = \frac{(\hat{\lambda}_Q)^2}{4/J}, \text{ with one degree of freedom.}$$

The dimension of the set of linear contrasts in 4 mean parameters is still a vector space with dimension $4 - 1 = 3$, so that there is a third orthogonal contrast

$\lambda_C = (-1)\mu_1 + (3)\mu_2 + (-3)\mu_3 + (1)\mu_4$, called the cubic contrast. The cubic contrast is the inner product of the vector $(-1, 3, -3, 1)$ and $(\mu_1, \mu_2, \mu_3, \mu_4)$. The statistic estimating the cubic contrast is $\hat{\lambda}_C = (-1)Y_{1\bullet} + (3)Y_{2\bullet} + (-3)Y_{3\bullet} + (1)Y_{4\bullet}$, and

$$SS_C = \frac{(\hat{\lambda}_C)^2}{[(-1)^2 + (3)^2 + (-3)^2 + (1)^2]/J} = \frac{(\hat{\lambda}_C)^2}{20/J}, \text{ with one degree of freedom.}$$

Each pair of these three contrasts is orthogonal. Each has an associated sum of squares based on 1 degree of freedom, and $SS_L + SS_Q + SS_C = SS_{Treatments}$.

Lack of fit test

There was a problem in the discussion of Chapter 8 that I wish to complete. That problem was:

A research team wishes to specify a manufacturing process so that Y , the area in a product affected by surface flaws is as small as possible. They have four levels of concentration of a chemical used to wash the product before the final manufacturing step and want to determine whether the concentration level causes a change in $E(Y)$. They run a balanced one-way layout with 6 observations for each concentration with level 1 set at 10%, level 2 set at 15%, level 3 set at 20%, and level 4 set at 25%. They run a balanced one-way layout with 6 observations for each treatment. They observe that $y_{1\bullet} = 264.5$, $y_{2\bullet} = 255.9$, $y_{3\bullet} = 216.2$, and $y_{4\bullet} = 263.8$, where $y_{i\bullet}$ is the average of the observations taken on the i th level. They also observe that $s_1^2 = 411.9$, $s_2^2 = 522.2$, $s_3^2 = 631.8$, and $s_4^2 = 521.9$, where s_i^2 is the unbiased estimate of the variance for the observations taken on the i th level.

- Complete the analysis of variance table for these results; that is, be sure to specify the degrees of freedom, sum of squares, mean square, and F-test.

- b. What is your conclusion? Use significance levels set to 0.10, 0.05, and 0.01. Make sure that you discuss the optimal setting of the concentration level and how you could document it.

Answer:

Analysis of Variance Table

Source	Sum of Squares	Degrees of Freedom	Mean Square	
Treatments	9467.4	3	3155.8	F=6.046
Error	10439.0	20	521.95	
Total	19906.4	23		

The critical values for the F-test are 2.38 for the 0.10 level, 3.10 for the 0.05 level, and 4.94 for the 0.01 level. Reject the null hypothesis that the mean responses of the four treatments are equal. The optimum setting is to use 20% as the concentration setting. This can be confirmed with Fisher's protected confidence intervals.

New questions:

- Compute the sums of squares for the linear, quadratic, and cubic contrasts.
- Complete the analysis of variance table for the linear regression of the dependent variable (area affected by surface flaws) on the chemical concentration. That is, use the sum of squares for the linear contrast as the sum of squares due to the linear regression. Test the null hypothesis that the dependent variable has no linear association with the independent variable. Use the 0.10, 0.05, and 0.01 levels of significance.
- Compute the sum of squares for lack of fit for the linear regression model. Complete the analysis of variance table that includes the sum of squares due to the linear regression, the sum of squares due to lack of fit for the linear regression, and the sum of squares for pure error. What is your conclusion? Use the 0.10, 0.05, and 0.01 levels of significance.

Answers:

- c. For the contrasts:

$$\hat{\lambda}_L = (-3)Y_{1\bullet} + (-1)Y_{2\bullet} + (1)Y_{3\bullet} + (3)Y_{4\bullet} = (-3) \times 264.5 + (-1) \times 255.9 + (1) \times 216.2 + (3) \times 263.8 = -41.8;$$

$$\hat{\lambda}_Q = (1)Y_{1\bullet} + (-1)Y_{2\bullet} + (-1)Y_{3\bullet} + (1)Y_{4\bullet} = (1) \times 264.5 + (-1) \times 255.9 + (-1) \times 216.2 + (1) \times 263.8 = 56.2;$$

$$\hat{\lambda}_c = (-1)Y_{1\bullet} + (3)Y_{2\bullet} + (-3)Y_{3\bullet} + (1)Y_{4\bullet} = (-1) \times 264.5 + (3) \times 255.9 + (-3) \times 216.2 + (1) \times 263.8 = 118.4.$$

Then the sums of squares of the contrasts are:

$$SS_L = \frac{(\hat{\lambda}_L)^2}{20/J} = \frac{(-41.8)^2}{20/6} = 524.172; \quad SS_Q = \frac{(\hat{\lambda}_Q)^2}{4/J} = \frac{(56.2)^2}{4/6} = 4737.66; \text{ and}$$

$$SS_C = \frac{(\hat{\lambda}_C)^2}{20/J} = \frac{(118.4)^2}{20/6} = 4205.568. \text{ An important check is to calculate}$$

$$SS_L + SS_Q + SS_C = 524.172 + 4737.66 + 4205.568 = 9467.4 = SS_{Treatments}.$$

d. The analysis of variance table based on the calculations of Chapter 11 is:

Analysis of Variance Table

Linear Regression

Source	Sum of Squares	Degrees of Freedom	Mean Square	
Regression on Chemical Concentration	524.172	1	524.172	F=0.595
Error	19382.228	22	881.010	
Total	19906.4	23		

Since the Chapter 11 F-test has the value 0.595, which is less than 1, we will accept the null hypothesis that there is no linear relation between dependent and independent variables. More formally, the critical value for a central F distribution with 1 numerator and 22 denominator degrees of freedom and level of significance 0.10 is 2.949. Since the F statistic (which has value 0.595) is less than 2.949, the null hypothesis is accepted at the 0.10 level (and of course at the 0.05 and 0.01 levels).

For part e, the lack of fit test analysis of variance table is below:

Analysis of Variance Table

Linear Regression and Lack of Fit Sums of Squares

Source	Sum of Squares	Degrees of Freedom	Mean Square	F test
Regression on Chemical Concentration	524.172	1	524.172	
Lack of (linear) fit	8943.228	2	4471.614	8.57
Pure Error	10439.0	20	521.95	
Total	19906.4	23		

There are two ways of calculating the sum of squares for lack of fit. The first is to add the sum of squares for the quadratic contrast and the sum of squares for the cubic contrast (that is, $4737.66 + 4205.568 = 8943.228$ with 2 degrees of freedom). The second is to subtract the sum of squares for the linear contrast from the sum of squares for treatments (that is, $9467.4 - 524.172 = 8943.228$ with 3-1 degrees of freedom). The lack of fit null hypothesis is that the linear model is adequate. The test statistic is the ratio of the mean square for lack of fit divided by the mean square for pure error. When the linear model is adequate, this test statistic should be equal to 1, modulo statistical variability. When the test statistic is larger than one, the linear model is not adequate. That is a more complex model is required. Here that more complex model may have both a linear and quadratic term.

Here, $F_{LOF} = \frac{MS_{LOF}}{MS_{PE}} = \frac{4471.614}{521.95} = 8.57$, with 2 numerator and 20 denominator degrees of freedom. The critical value for a central F distribution with 2 numerator and 20 denominator degrees of freedom and level of significance 0.01 is 5.850. Since $F_{LOF} = \frac{MS_{LOF}}{MS_{PE}} = 8.57 > 5.850$, we reject the null hypothesis that the linear model is adequate at the 0.01 level.

Comments

One should always take advantage of observations from repeated independent variable values. The lack of fit test is an objective data-based indicator of the adequacy of a model. A research team should design its study so that there are repeated independent variable values. One can also group observations into sets that have approximately the same independent variable values so that one can calculate an approximate lack of fit test.

We can continue our example with the techniques from Chapter 12. The linear model is that $Y_i = \beta_0 + \beta_1 x_i + \sigma_{L|x} Z_i$. We accepted $H_{0L} : \beta_1 = 0$ rather than $H_{1L} : \beta_1 \neq 0$ because the relevant F statistic was 0.595, and we rejected the null hypothesis that the linear model was adequate (because the lack of fit F statistic was 8.57, significant at the 0.01 level). It is natural to consider a quadratic model $Y_i = \gamma_0 + \gamma_1 x_i + \gamma_2 x_i^2 + \sigma_{Q|x} Z_i$. The analysis of variance table is shown below. The F-test for the quadratic coefficient given that the linear coefficient is in the model is

$$F_{QL} = \frac{SS_Q / 1}{MSE} = \frac{4737.66}{697.36} = 6.79 \text{ with 1 numerator and 21 denominator degrees of}$$

freedom. The critical values are 2.961 for the 0.10 level, 4.325 for the 0.05 level and 8.017 for the 0.01 level. We would reject $H_{0Q} : \gamma_2 = 0$ rather than $H_{1Q} : \gamma_2 \neq 0$ at the 0.10 and 0.05 levels. We would accept $H_{0Q} : \gamma_2 = 0$ at the 0.01 level.

Analysis of Variance Table
Linear and Quadratic Regression

Source	Sum of Squares	Degrees of Freedom	Mean Square	F test
Linear Regression on Chemical Concentration	524.172	1	524.172	
Quadratic regression given linear regression	4737.66	1	4737.66	6.79
Error	14644.568	21	697.36	
Total	19906.4	23		

In the event that we use the linear and quadratic model, we can test for lack of fit. The analysis of variance table is shown below.

Analysis of Variance Table
Linear and Quadratic Regression and Lack of Fit Sums of Squares

Source	Sum of Squares	Degrees of Freedom	Mean Square	F test
Linear Regression on Chemical Concentration	524.172	1	524.172	
Quadratic regression given linear regression	4737.66	1	4737.66	
Lack of (quadratic) fit	4205.568	1	4205.568	8.057
Pure Error	10439.0	20	521.95	
Total	19906.4	23		

The lack of fit F-test for the quadratic model is $F_{Q\ LOF} = \frac{SS_{LOF\ Q}/1}{MSE} = \frac{4205.568}{521.95} = 8.057$ with 1 numerator and 20 denominator degrees of freedom. The critical values are 2.975 for the 0.10 level, 4.351 for the 0.05 level and 8.096 for the 0.01 level. We would reject the adequacy of the quadratic model at the 0.10 and 0.05 levels. We would accept the adequacy of the quadratic model at the 0.01 level (just barely).

Data Analysis of the One-Way Layout

When the global null hypothesis is rejected, researchers want to know which settings of the treatment variable are associated with larger expected values and which are associated with smaller expected values. Such questions lead to the issues of multiple comparisons, which is covered more deeply in Chapter Nine. A relatively simple approach is to use Fisher's protected t confidence intervals (which is also called Fisher's Least Significant Difference). Fisher's protected confidence intervals are calculated only when the global null hypothesis is rejected. Then one calculates a confidence interval for $E(Y_{ij}) - E(Y_{i'j}) = \mu_i - \mu_{i'}$ for each pair of treatment settings (i, i') using a procedure analogous to the procedures in Chapter 6. These comparisons are called *post hoc* comparisons. For example, the 99% confidence interval for

$\mu_1 - \mu_2$ would be $y_{1\bullet} - y_{2\bullet} \pm t_{2,576,n-1} \sqrt{MSE(\frac{1}{J_1} + \frac{1}{J_2})}$. When the experiment is

underpowered, it might well happen that the global null hypothesis is rejected with no protected confidence interval excluding zero.

8.4 Checking on the AOV Conditions

The most important assumption is that of independence. This is guaranteed in a randomized experiment in which the experimental units are randomly assigned to treatment. When the data do not come from a randomized experiment, this assumption should be checked carefully. Common problems occur when time series data (for example, an exchange rate on successive days as the dependent variable) is used. Also data describing a geographical area such as a census tract have spatial autocorrelation. Data on students from the same class will be correlated because of the common instruction.

The analysis procedures for balanced analyses of variances are not sensitive to violations of the normality assumption and the homogeneity of variance assumption. The residuals in a one-way AOV are $r_{ij} = y_{ij} - y_{i\bullet}$. Residual analysis is

simple for this model. One can and should generate a probability plot of the residuals. Closeness of the plot to a straight suggests that the assumption of normality appears to be true. Hartley's $F_{\max} = \frac{S_{\max}^2}{S_{\min}^2}$ is sensitive to normality. The

Brown-Forsythe-Levene test for homogeneity of variance is more robust to violations of the assumption of normality. Many statistical packages will calculate this test, and you should use it routinely. There is another more robust test of this null hypothesis that corrects for the estimated kurtosis of the sampled random variables. Some statistical packages report this test as well or instead of Levene's test. If the hypothesis of constant variance is rejected, there are two common next steps. One is to use weighted least squares (with weights reflecting the difference in variance of observations), and the other approach is to transform the data to lessen the differences in variance. These transformations are called variance stabilizing transformations and are commonly used. They are helpful, especially when predictions of future values are to be made.

8.5 An Alternative Analysis: Transformations of the Data

Lecture Material on the "Delta Method":

The objective is to calculate the approximate mean and variance of a random variable $W = f(Y)$. The random variable Y has expected value μ_Y and variance σ_Y^2 , and the function f has finite derivatives. The delta method approximates the value of W using the first term of the Taylor series: $W \cong f(\mu_Y) + f'(\mu_Y)(Y - \mu_Y)$. Then, $E(W) \cong E[f(\mu_Y) + f'(\mu_Y)(Y - \mu_Y)] = E[f(\mu_Y)] + E[f'(\mu_Y)(Y - \mu_Y)]$. Now $E[f(\mu_Y)] = f(\mu_Y)$, because $f(\mu_Y)$ is a constant. For the second term, $E[f'(\mu_Y)(Y - \mu_Y)] = f'(\mu_Y)E[(Y - \mu_Y)] = 0$. The conclusion is that $E(W) \cong f(\mu_Y)$. The result that $E[F(r)] \cong F(\rho)$ is an application of this result.

The deviation $W - E(W) \cong f(\mu_Y) + f'(\mu_Y)(Y - \mu_Y) - f(\mu_Y) = f'(\mu_Y)(Y - \mu_Y)$. Then $E\{[W - E(W)]^2\} \cong E\{[f'(\mu_Y)(Y - \mu_Y)]^2\} = [f'(\mu_Y)]^2 E[(Y - \mu_Y)^2]$. That is, $\text{var}(W) \cong [f'(\mu_Y)]^2 \text{var}(Y)$.

Example Problem

The random variable Y has the Poisson distribution with expected value μ . Find the approximate mean and variance of \sqrt{Y} .

Solution

The expectation calculation is easy: $E(\sqrt{Y}) \cong \sqrt{\mu}$. Next, the first derivative calculation: $f(y) = y^{0.5}$ so that $f'(y) = 0.5y^{-0.5}$. From Chapter 4, $\text{var}(Y) = \mu$. Finally, $\text{var}(\sqrt{Y}) \cong [f'(\mu_Y)]^2 \text{var}(Y) = [0.5\mu^{-0.5}]^2 \mu = 0.25$.

Comment

Since $\text{var}(\sqrt{Y}) \cong 0.25$ independently of the expected value of the Poisson random variable, the square root transformation is said to be a variance stabilizing transformation for a Poisson random variable. Another transformation of a Poisson random variable is

$W = \sqrt{Y} + \sqrt{Y+1}$ has approximate expected value $E(W) \cong \sqrt{\mu} + \sqrt{\mu+1}$ and approximate variance $\text{var}(W) \cong 1$. This transformation is an example of a Freeman-Tukey deviate.

Example examination problem

The random variable Y , $Y > 0$, has $E(Y) = \theta$ and $\text{var}(Y) = \theta^3, \theta > 0$. Find the approximate mean and variance of $W = \ln(Y)$.

Answer: $E(W) \cong \ln(\theta)$, and $\text{var}(W) \cong [f'(\theta)]^2 \text{var}(Y) = [1/\theta]^2 \theta^3 = \theta$.

Exploratory Data Analysis Tool to Identify Variance Stabilizing Transformation

When the dependent variable in an analysis of variance is always positive,

calculate $y_{i\cdot}$ and $s_i^2 = \frac{\sum_{j=1}^{n_i} (y_{ij} - y_{i\cdot})^2}{n_i - 1}$, where i indexes the treatments in the analysis of variance. Plot $\log(s_i)$ against $\log(y_{i\cdot})$ and fit a straight line to the data. Call the slope m . When $m \cong 0$, no transformation is necessary. When $m \neq 0$, then analyze the transformed values $t_{ij} = y_{ij}^{1-m}$, with $m \neq 1$. When $m = 1$, use $t_{ij} = \log(y_{ij})$. There is a related set of techniques called the Box-Cox transformations that is also helpful.

Probability theory behind tool

The random variable Y has mean μ_Y and standard deviation σ_Y such that $\ln(\sigma_Y) = a + m \ln(\mu_Y)$. Then $\sigma_Y = \exp[\ln(\sigma_Y)] = \exp[a + m \ln(\mu_Y)] = \exp[a + \ln(\mu_Y^m)]$. Further, $\sigma_Y = \exp[\ln(\sigma_Y)] = \exp[a + m \ln(\mu_Y)] = \exp[a + \ln(\mu_Y^m)]$. This means that the standard deviation of Y is related to the expected value of Y in the equation $\sigma_Y = \exp[a + \ln(\mu_Y^m)] = \exp(a) \times \exp[\ln(\mu_Y^m)] = c \mu_Y^m$. For $W = f(Y)$, $\text{var}(W) \cong [f'(\mu_Y)]^2 \text{var}(Y)$, which is equivalent to $\sigma_W \cong |f'(\mu_Y)| \sigma_Y$. For functions f that are monotonically increasing, $\sigma_W \cong f'(\mu_Y) \sigma_Y$.

We seek to choose f so that $\sigma_W \cong f'(\mu_Y) \sigma_Y = k$, where k is a constant. That is, we seek f so that $\sigma_W \cong f'(\mu_Y) c \mu_Y^m = k$. That is, we seek f such that $f'(\mu_Y) = (\frac{k}{c}) \mu_Y^{-m}$. For $m \neq 1$, anti-differentiation finds the function to be $f(\mu_Y) = c' \mu_Y^{1-m} + c''$. Typically, $c' = 1$ and $c'' = 0$. The required transformation of Y is $W = f(Y) = Y^{1-m}$, $m \neq 1$. For $m = 1$, $W = f(Y) = \ln(Y)$.

Example test problem

The random variable Y has $E(Y) = \theta$ and $\text{var}(Y) = \theta^{2.5}$, $\theta > 0$. Find the transformation W that makes the variance of W approximately constant. What are the approximate mean and variance of W ?

Solution: Since $\sigma_Y = \theta^{1.25}$, $\ln(\sigma_Y) = 1.25 \ln(\theta) = 1.25 \ln(\mu_Y)$. The rule for finding a variance stabilizing transformation is to use the transformation $W = Y^{1-1.25} = \frac{1}{Y^{0.25}}$.

From the delta method, $E(W) \cong \frac{1}{\theta^{0.25}}$. The derivative of the transformation function is $f'(y) = -0.25 y^{-1.25}$. Then $\text{var}(W) \cong (-0.25 \times \theta^{-1.25})^2 \theta^{2.5} = (0.25)^2$.

Bonferroni's (Boole's) Inequality (Chapter 9)

Let R_1 be the event that null hypothesis one (H_{01}) is rejected, and let R_2 be the event that null hypothesis one (H_{02}) is rejected. Then from Chapter 4,

$\Pr\{R_1 \cup R_2\} = \Pr\{R_1\} + \Pr\{R_2\} - \Pr\{R_1 \cap R_2\}$. When R_1 and R_2 are independent, $\Pr\{R_1 \cap R_2\} = \Pr\{R_1\} \times \Pr\{R_2\}$. Otherwise, it may be difficult to calculate $\Pr\{R_1 \cap R_2\}$. Boole's inequality compares $\Pr\{R_1 \cup R_2\} = \Pr\{R_1\} + \Pr\{R_2\} - \Pr\{R_1 \cap R_2\}$ to

$\Pr\{R_1\} + \Pr\{R_2\}$. The result is $\Pr\{R_1 \cup R_2\} \leq \Pr\{R_1\} + \Pr\{R_2\}$. Although this result was found by Boole, virtually all researchers refer to it as Bonferroni's inequality.

In statistics, $\Pr_0\{R_1 \cup R_2\} = \Pr_0\{\text{At least one Type I error}\} = \alpha_{\text{Overall}}$. It is natural to let $\alpha_1 = \Pr_0\{R_1\}$ and let $\alpha_2 = \Pr_0\{R_2\}$. Then Boole's (Bonferroni's) inequality is that $\alpha_{\text{Overall}} \leq \alpha_1 + \alpha_2$. The bound on the overall level of significance is the sum of the levels of significance of each individual test. This property is described as inflation of the significance level.

Let R_3 be the event that null hypothesis three (H_{03}) is rejected with level of significance α_3 , ..., and let R_N be the event that null hypothesis N (H_{0N}) is rejected with level of significance α_N . Then the inequality generalizes, and

$\Pr_0\{R_1 \cup \dots \cup R_N\} = \Pr_0\{\text{At least one Type I error}\} = \alpha_{\text{Overall}} \leq \alpha_1 + \dots + \alpha_N$. Conventionally, researchers set the individual levels of significance to be less than some relatively small value $\varepsilon > 0$ such as 0.01 or 0.05. Then, the inequality is

$\Pr_0\{R_1 \cup \dots \cup R_N\} = \Pr_0\{\text{At least one Type I error}\} = \alpha_{\text{Overall}} \leq \alpha_1 + \dots + \alpha_N \leq \varepsilon$. Researchers may set $\alpha_1 = \dots = \alpha_N = \alpha$ so that

$\Pr_0\{R_1 \cup \dots \cup R_N\} = \Pr_0\{\text{At least one Type I error}\} = \alpha_{\text{Overall}} \leq N\alpha \leq \varepsilon$. That is, researchers set each individual significance level to $\frac{\varepsilon}{N}$. Equivalently, in a multiple regression analysis, one can multiply the p-value of a regression coefficient by N , the number of independent variables considered.

Scheffe Confidence Intervals for All Contrasts

In a balanced one-way layout with $I = 4$, let $\lambda = a_1\mu_1 + a_2\mu_2 + a_3\mu_3 + a_4\mu_4$, where

$\mu_i = E(Y_{ij}) = E(Y_{i\bullet})$. The $a_i, i = 1, \dots, I$ are fixed constants with $\sum_{i=1}^I a_i = 0$. Here

$\lambda = a_1\mu_1 + a_2\mu_2 + a_3\mu_3 + a_4\mu_4$ is a linear combination of parameters and hence is itself a parameter. We need a statistic $\hat{\lambda}$ such that $E(\hat{\lambda}) = \lambda$.

That statistic is $\hat{\lambda} = a_1Y_{1\bullet} + a_2Y_{2\bullet} + a_3Y_{3\bullet} + a_4Y_{4\bullet}$. The expected value is

$E(\hat{\lambda}) = a_1E(Y_{1\bullet}) + a_2E(Y_{2\bullet}) + a_3E(Y_{3\bullet}) + a_4E(Y_{4\bullet}) = a_1\mu_1 + a_2\mu_2 + a_3\mu_3 + a_4\mu_4 = \lambda$. The variance

is $\text{var}(\hat{\lambda}) = a_1^2 \text{var}(Y_{1\bullet}) + a_2^2 \text{var}(Y_{2\bullet}) + a_3^2 \text{var}(Y_{3\bullet}) + a_4^2 \text{var}(Y_{4\bullet}) = \sum_{i=1}^4 a_i^2 \frac{\sigma_{1W}^2}{J}$, where J is the

number of observations in each treatment. A 99% confidence interval for λ , with

σ_{1W} known, is then $\hat{\lambda} \pm 2.576 \sqrt{\sum_{i=1}^4 a_i^2 \frac{\sigma_{1W}^2}{J}}$. As usual, σ_{1W}^2 is not known, and we use the *MSE* to estimate it. The 99% confidence interval for λ , with σ_{1W} unknown, is then $\hat{\lambda} \pm t_{2.576, n-I} \sqrt{\sum_{i=1}^4 a_i^2 \frac{MSE}{J}}$, where *MSE* has $n-I = n-4$ degrees of freedom.

This calculation assumes that there is only one contrast that we want a confidence interval for. Scheffe's method gives confidence intervals for all contrasts. Mathematically, the set of contrasts using $\{\mu_1, \mu_2, \mu_3, \mu_4\}$ is a vector space with dimension $4-1=3$. The 99% Scheffe confidence interval for λ , with σ_{1W}

unknown, is $\hat{\lambda} \pm \sqrt{3F_{0.99, 3, n-4}} \sqrt{\sum_{i=1}^4 a_i^2 \frac{MSE}{J}}$, where *MSE* has $n-4$ degrees of freedom.

We can calculate a 99% Scheffe confidence interval for any linear contrast. The confidence is then 99% that all of the confidence intervals are simultaneously correct.

For the example that we have been studying, $\sqrt{3F_{0.99, 3, 24-4}} = \sqrt{(3 \times 4.938)} = 3.849$. The 99% Scheffe confidence interval for $\lambda_L = (-3)\mu_1 + (-1)\mu_2 + (1)\mu_3 + (3)\mu_4$ is

$$\hat{\lambda}_L \pm \sqrt{3F_{0.99, 3, 20}} \sqrt{\sum_{i=1}^4 a_i^2 \frac{MSE}{J}} = -41.8 \pm \sqrt{14.814} \sqrt{\frac{20 \times 521.95}{6}} = -41.8 \pm 160.55. \text{ Similarly, the 99\%}$$

Scheffe confidence interval for $\lambda_Q = (1)\mu_1 + (-1)\mu_2 + (-1)\mu_3 + (1)\mu_4$ is

$$\hat{\lambda}_Q \pm \sqrt{3F_{0.99, 3, 20}} \sqrt{\sum_{i=1}^4 a_i^2 \frac{MSE}{J}} = 56.2 \pm \sqrt{14.814} \sqrt{\frac{4 \times 521.95}{6}} = 56.2 \pm 71.80. \text{ The 99\% Scheffe}$$

confidence interval for $\lambda_C = (-1)\mu_1 + (3)\mu_2 + (-3)\mu_3 + (1)\mu_4$ is

$$\hat{\lambda}_C \pm \sqrt{3F_{0.99, 3, 20}} \sqrt{\sum_{i=1}^4 a_i^2 \frac{MSE}{J}} = 118.4 \pm \sqrt{14.814} \sqrt{\frac{20 \times 521.95}{6}} = 118.4 \pm 160.55. \text{ The confidence that}$$

these three intervals are simultaneously correct is 99%. There is no need to adjust the significance level for the multiple comparisons made, as in the Bonferroni inequality.

One can form Scheffe confidence intervals for any number of contrasts. For example, the 99% Scheffe confidence interval for $\mu_3 - \mu_2$ is

$$y_{3\bullet} - y_{2\bullet} \pm \sqrt{3F_{0.99, 3, 20}} \sqrt{\sum_{i=1}^4 a_i^2 \frac{MSE}{J}} = 216.2 - 255.9 \pm \sqrt{14.814} \sqrt{\frac{2MSE}{6}} = -39.7 \pm 50.7. \text{ This}$$

Scheffe interval includes zero, while the protected t-confidence interval excluded zero. This is an example of the greater strictness of the Scheffe confidence interval.

Tukey's W procedure

To use this procedure, one first calculates the Tukey sampling margin of error

$W = q_\alpha(I, \nu) \sqrt{\frac{MSE}{J}}$, where I is the number of treatments in the one-way layout,

$\nu = n - I$ is the degrees of freedom of the MSE, and $q_\alpha(I, \nu)$ is a percentile extracted from Table 10 (percentage points of the Studentized range). For the current

example with $\alpha = 0.01$, $W = q_\alpha(I, \nu) \sqrt{\frac{MSE}{J}} = 5.02 \sqrt{\frac{521.95}{6}} = 46.82$. The 99% Tukey

confidence interval for $\mu_3 - \mu_2$ is $y_{3\bullet} - y_{2\bullet} \pm W = 216.2 - 255.9 \pm 46.82 = -39.7 \pm 46.82$.

Since this includes zero, we would conclude that the expected value of the second treatment is equal to the expected value of the third treatment at the 0.01 level of significance. Similarly, the 99% Tukey confidence interval for $\mu_3 - \mu_1$ is

$y_{3\bullet} - y_{1\bullet} \pm W = 216.2 - 264.5 \pm 46.82 = -48.3 \pm 46.82$. The 99% Tukey confidence interval for $\mu_3 - \mu_1$ excludes zero, suggesting that these two treatments have different expected values. In general, Tukey confidence intervals for $\mu_i - \mu_j$ are narrower than the Scheffe confidence intervals.

Working-Hotelling Confidence Band for $E[\hat{Y}(x)]$

An important Chapter 11 result was that $\hat{Y}(x)$ has the normal distribution

$$N(\beta_0 + \beta_1 x, \sigma_{Y|x}^2 (\frac{1}{n} + \frac{(x - \bar{x}_n)^2}{\sum_{i=1}^n (x_i - \bar{x}_n)^2})).$$

The 95% confidence interval for $E[\hat{Y}(x)] = \beta_0 + \beta_1 x$, where x is a single value is then

$$\hat{Y}(x) \pm t_{1.960, n-2} \sqrt{MSE (\frac{1}{n} + \frac{(x - \bar{x}_n)^2}{\sum_{i=1}^n (x_i - \bar{x}_n)^2})}. \text{ The same multiple comparison issues hold}$$

when one seeks confidence intervals for $\hat{Y}(x)$ for multiple values of x . The application of the theory of the Scheffe confidence intervals to the problem of finding confidence intervals for multiple independent variable values generates the

"Working-Hotelling Confidence Region for $E[\hat{Y}(x)]$. This region is

$$\hat{Y}(x) \pm \sqrt{2F_{0.05, 2, n-2}} \sqrt{MSE (\frac{1}{n} + \frac{(x - \bar{x}_n)^2}{\sum_{i=1}^n (x_i - \bar{x}_n)^2})}, -\infty < x < \infty. \text{ There is then 95\% confidence}$$

that the regression line is completely covered by this region.