

Data Analysis
Spring Semester, 2023
February 16, 2023
Lecture 9

Thursday: Examination One: Chapters 3, 4, 5, 6, and 7. I will be available on Zoom from 4:30 pm until 6:15 in case there are problems.

Chapter Six
Inferences Comparing Two Population Central Values

Two Independent Sample Test

Let X_1, X_2, \dots, X_n be a random sample of size n from the random variable X , which is $N(\mu_X, \sigma_X^2)$. Let B_1, B_2, \dots, B_m be a random sample of size m from the random variable B , which is $N(\mu_B, \sigma_B^2)$. The two samples are independent of each other.

The random variable \bar{X}_n is $N(\mu_X, \frac{\sigma_X^2}{n})$, and \bar{B}_m is $N(\mu_B, \frac{\sigma_B^2}{m})$. The two sample averages are independent. We seek to use this data to test $H_0 : E(X - B) = 0$ against the alternative hypothesis $H_1 : E(X - B) \neq 0$ at the α level of significance. Our test statistic is $TS = \bar{X}_n - \bar{B}_m$.

Distribution of the Test Statistic

The distribution of $TS = \bar{X}_n - \bar{B}_m$ is normal; and
 $E(TS) = E(\bar{X}_n - \bar{B}_m) = E(\bar{X}_n) - E(\bar{B}_m) = \mu_X - \mu_B$, and
 $\text{var}(TS) = \text{var}(\bar{X}_n - \bar{B}_m) = \text{var}(\bar{X}_n) + \text{var}(\bar{B}_m) = \frac{\sigma_X^2}{n} + \frac{\sigma_B^2}{m}$.

Null Distribution of the Test Statistic

Since $H_0 : E(X - B) = 0$ which is equivalent to $H_0 : \mu_X = \mu_B$,
 $E_0(TS) = E_0(\bar{X}_n - \bar{B}_m) = \mu_X - \mu_B = 0$.

Under the null hypothesis that the distribution of X and B are identical, not only does $\mu_X = \mu_B$ under the null, but also $\sigma_X^2 = \sigma_B^2 = \sigma^2$. Using this assumption,

$$\text{var}_0(TS) = \text{var}_0(\bar{X}_n - \bar{B}_m) = \frac{\sigma^2}{n} + \frac{\sigma^2}{m} = \sigma^2 \left(\frac{1}{n} + \frac{1}{m} \right).$$

Test when variances unknown but equal

Just as in Chapter 5, we use an estimate of σ^2 and stretch the critical values an amount determined by the degrees of freedom of our estimate. There are many

estimates of σ^2 . For example, $S_X^2 = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{n-1}$ and $S_B^2 = \frac{\sum_{i=1}^m (B_i - \bar{B}_m)^2}{m-1}$ are unbiased estimates of σ^2 , with $n-1$ and $m-1$ degrees of freedom respectively. That is, $E(S_X^2) = E(S_B^2) = \sigma^2$. We use both estimates. Let $S_P^2 = \frac{(n-1)S_X^2 + (m-1)S_B^2}{n+m-2}$.

This estimator has $n+m-2$ degrees of freedom. Then our studentized statistic is

$$T_{n+m-2} = \frac{\bar{X}_n - \bar{B}_m - 0}{\sqrt{S_P^2 \left(\frac{1}{n} + \frac{1}{m} \right)}}.$$

If $\alpha = 0.10$, we reject H_0 when $|T_{n+m-2}| \geq t_{1.645, n+m-2}$. Similarly, if $\alpha = 0.05$, we reject H_0 when $|T_{n+m-2}| \geq t_{1.960, n+m-2}$. If $\alpha = 0.01$, we reject H_0 when $|T_{n+m-2}| \geq t_{2.576, n+m-2}$.

Sample size for two sample test:

Problem 6 in Chapter 6 Study Guide

In a clinical trial, $2J$ patients suffering from an illness will be randomly assigned to one of two groups so that J will receive an experimental treatment and J will receive the best available treatment. The random variable X is the response of a patient to the experimental medicine, and the random variable B is the response of a patient to the best currently available treatment. The random variables X and B are normally distributed. The null hypothesis to be tested is that $E(X) - E(B) = 0$ against the alternative that $E(X) - E(B) > 0$ at the α , $\alpha \leq 0.5$, level of significance. When the null hypothesis is true, $\text{var}(X) = \text{var}(B) = \sigma_0^2$. When the alternative hypothesis is true, $\text{var}(B) = \sigma_0^2$, but $\text{var}(X) = \sigma_1^2 > \sigma_0^2$. What is the number J in each group that would have to be taken so that the probability of a Type II error for the test of the null hypothesis specified in the common section is β , $\beta \leq 0.5$, when $E(X) - E(B) = \Delta > 0$?

Solution: The test statistic is $TS = \bar{X}_J - \bar{B}_J$, and TS is

$N(E(X) - E(B), \frac{\text{var}(X)}{J} + \frac{\text{var}(B)}{J})$. The null distribution of TS is then

$N(0, \frac{\sigma_0^2}{J} + \frac{\sigma_0^2}{J})$. Hence, we reject $H_0 : E(X) - E(B) = 0$ against $H_1 : E(X) - E(B) > 0$

at the α level of significance when $TS \geq 0 + |z_\alpha| \sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_0^2}{J}}$. When

$E(X) - E(B) = \Delta > 0$ and $\text{var}(X) = \sigma_1^2 > \sigma_0^2$, $\text{var}(B) = \sigma_0^2$, the (alternative)

distribution of TS is then $N(\Delta, \frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J})$. Then, the probability of a Type II

error is $\beta = \Pr\{\text{Accept } H_0\} = \Pr\{TS < 0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}}\}$. That is,

$$\beta = \Pr\{TS < 0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}}\} = \Pr\{Z = \frac{TS - \Delta}{\sigma_1(\bar{X}_J - \bar{B}_J)} < \frac{0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}} - \Delta}{\sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J}}}\}. \text{ Since}$$

$\beta = \Pr\{\text{Accept } H_0\} \leq 0.5$, it is true that $\beta = \Pr\{Z < -|z_\beta|\}$. We now have two equations:

$$\beta = \Pr\{Z < \frac{0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}} - \Delta}{\sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J}}}\}, \text{ and}$$

$$\beta = \Pr\{Z < -|z_\beta|\}.$$

The problem is to choose J so that the probability of a Type II error is a specified value. That is, we should choose J so that the right-hand sides of the

two equations are equal:
$$\frac{0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}} - \Delta}{\sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J}}} = -|z_\beta|.$$

We solve for J in the equation above. This reduces to:

$$0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}} - \Delta = -|z_\beta| \sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J}}. \text{ That is, } |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}} + |z_\beta| \sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J}} = \Delta.$$

Next, solve for J to get $\sqrt{J} = \frac{|z_\alpha| \sqrt{2\sigma_0^2} + |z_\beta| \sqrt{\sigma_0^2 + \sigma_1^2}}{\Delta}$. Since J has to be an integer, we increase J to the next integer value.

Chapter 7

Inferences about Population Variances

Probability Theory Facts

Let Z be $N(0,1)$. Then Z^2 has the (central) chi-squared distribution with 1 degree of freedom. This is denoted χ_1^2 .

Let Z_1, Z_2, \dots, Z_n be $NID(0,1)$. Then $S_n = Z_1^2 + Z_2^2 + \dots + Z_n^2 = \sum_{i=1}^n Z_i^2$ follows the (central) chi-square distribution with n degrees of freedom, denoted χ_n^2 . The expected value of a χ_n^2 is n : $E(S_n) = E(Z_1^2) + E(Z_2^2) + \dots + E(Z_n^2)$. Since $\text{var}(Z) = 1 = E(Z^2) - [E(Z)]^2 = 1$, then $E(Z^2) - [0]^2 = 1$. Using this in $E(S_n) = E(Z_1^2) + E(Z_2^2) + \dots + E(Z_n^2) = n$. Further, the variance of a chi-square distribution with n degrees of freedom is $2n$: $\text{var}(S_n) = 2n$

Let Y be $N(\mu_Y, \sigma_Y^2)$. Then, $\frac{Y - \mu_Y}{\sigma_Y} = Z$ is $N(0,1)$. Let Y_1, Y_2, \dots, Y_n be a random sample

from Y , which is $N(\mu_Y, \sigma_Y^2)$. Then $\sum_{i=1}^n \left(\frac{Y_i - \mu_Y}{\sigma_Y}\right)^2$ is χ_n^2 . After factoring out σ_Y^2 ,

$$\frac{\sum_{i=1}^n (Y_i - \mu_Y)^2}{\sigma_Y^2} \text{ is also } \chi_n^2.$$

Since μ_Y is not known in applications, it must be estimated. An important property

of a sample from a normal distribution is that $\frac{\sum_{i=1}^n (Y_i - \bar{Y}_n)^2}{\sigma_Y^2}$ is distributed as χ_{n-1}^2 .

That is, using the sample mean has reduced the degrees of freedom by one. From AMS 310, the unbiased estimator of the sample variance is $S^2 = \frac{\sum (Y_i - \bar{Y}_n)^2}{n-1}$. Since

$$(n-1)S^2 = (n-1) \frac{\sum (Y_i - \bar{Y}_n)^2}{n-1} = \sum_{i=1}^n (Y_i - \bar{Y}_n)^2, \quad \frac{(n-1)S^2}{\sigma_Y^2} \text{ has a central chi-squared}$$

distribution with $n-1$ degrees of freedom when Y_1, Y_2, \dots, Y_n is a sample of size n from a $N(\mu_Y, \sigma_Y^2)$ distribution. This is our first important use of the chi-squared distribution.

The F-distribution

Let X be $N(\mu_X, \sigma_X^2)$. Using the standard score transformation, $\frac{X - \mu_X}{\sigma_X} = Z$ is $N(0,1)$.

Let X_1, X_2, \dots, X_n be a random sample of size n from X , which is $N(\mu_X, \sigma_X^2)$. Then,

$$\frac{\sum_{i=1}^n (X_i - \mu_X)^2}{\sigma_X^2} \text{ is } \chi_n^2.$$

Let Y be $N(\mu_Y, \sigma_Y^2)$. Then, $\frac{Y - \mu_Y}{\sigma_Y} = Z$ is also $N(0,1)$. Let Y_1, Y_2, \dots, Y_m be a random

sample of size m from Y , which is $N(\mu_Y, \sigma_Y^2)$. Then, $\frac{\sum_{i=1}^m (Y_i - \mu_Y)^2}{\sigma_Y^2}$ is χ_m^2 . The

definition of the central F distribution is that the random variable

$$F_{n,m} = \frac{\{[\sum_{i=1}^n (X_i - \mu_X)^2] / [n\sigma_X^2]\}}{\{[\sum_{i=1}^m (Y_i - \mu_Y)^2] / [m\sigma_Y^2]\}} \text{ has a (central) F distribution with } n \text{ numerator and } m$$

denominator degrees of freedom.

Application of the F distribution

The problem with this random variable is that the expected values are not known. As before, we use the sample averages as estimates of the expected values. The penalty for using sample data rather than expected values is a one degree reduction in both the numerator and denominator degrees of freedom. That is,

$$\frac{\{[\sum_{i=1}^n (X_i - \bar{X}_n)^2] / [(n-1)\sigma_X^2]\}}{\{[\sum_{i=1}^m (Y_i - \bar{Y}_m)^2] / [(m-1)\sigma_Y^2]\}} = \frac{S_X^2 / \sigma_X^2}{S_Y^2 / \sigma_Y^2} = F_{n-1, m-1} \text{ has a central F distribution with } n-1$$

numerator and $m-1$ denominator degrees of freedom. Of course, there is still the issue of the unknown variances of X and Y that has to be dealt with.

One use of this random variable is to test the null hypothesis $H_0 : \sigma_X^2 = \sigma_Y^2$. The most common alternative hypothesis is $H_1 : \sigma_X^2 > \sigma_Y^2$. The test statistic for this hypothesis

is $TS = \frac{S_X^2 / \sigma_X^2}{S_Y^2 / \sigma_Y^2}$. Under the null hypothesis $H_0 : \sigma_X^2 = \sigma_Y^2$, $TS = \frac{S_X^2}{S_Y^2}$, and its null

distribution is central $F_{n-1, m-1}$. Under the null hypothesis $E(S_X^2) = \sigma_X^2$ and $E(S_Y^2) = \sigma_Y^2$,

so that $E_0(TS) \cong \frac{E(S_X^2)}{E(S_Y^2)} = 1$. Under the alternative hypothesis, $E_1(TS) \cong \frac{E(S_X^2)}{E(S_Y^2)} > 1$. That is, the test of $H_0 : \sigma_X^2 = \sigma_Y^2$ against the alternative $H_1 : \sigma_X^2 > \sigma_Y^2$ is a right-sided test. A value of TS near 1 (modulo statistical variation) supports the null hypothesis, and a value of TS much greater than 1 supports the alternative. The next problem illustrates the test.

Problem 3 from Chapter 7 Study Guide

A research team took a random sample of 9 observations from a normally distributed random variable Y and observed that $\bar{y}_9 = 91.2$ and $s_Y^2 = 229.6$, where \bar{y}_9 was the average of the nine observations sampled from Y and s_Y^2 was the unbiased estimate of $\text{var}(Y)$. A second research team took a random sample of 10 observations from a normally distributed random variable X and observed that $\bar{x}_{10} = 103.5$ and $s_X^2 = 917.6$, where \bar{x}_{10} was the average of the ten observations sampled from X and s_X^2 was the unbiased estimate of $\text{var}(X)$. Test the null hypothesis $H_0 : \text{var}(X) = \text{var}(Y)$ against the alternative $H_1 : \text{var}(X) > \text{var}(Y)$ at the 0.10, 0.05, and 0.01 levels of significance.

Solution: One has a choice of which sample variance to put in the numerator. When one puts the variance *hypothesized* to be larger in the numerator, then the test is right-sided. Here $ts = \frac{s_X^2}{s_Y^2} = \frac{917.6}{229.6} = 3.9965$, with 9 numerator and 8

denominator degrees of freedom. The critical value for the 0.10 level is 2.56; for the 0.05 level, 3.39; and for the 0.01 level, 5.91. The correct decision is to reject the null hypothesis at the 0.10 and 0.05 levels and accept it at the 0.01 level. As before, the sample means are not needed for the problem. Students who use the information given as the cue to their choice of statistical tests sometimes respond to a question like this with a two-sample t-test. This is incorrect.

Confidence interval for the ratio of variances $\frac{\sigma_X^2}{\sigma_Y^2}$

For this task, we use $TS = \frac{S_Y^2 / \sigma_Y^2}{S_X^2 / \sigma_X^2}$, which has an F-distribution with $m-1$ numerator and $n-1$ denominator degrees of freedom. This choice may be counter-intuitive, but is necessary. The percentage points in Table 8 are based on right sided tail areas, so

that $\Pr\{F_{1-\alpha/2, m-1, n-1} < \frac{S_Y^2 / \sigma_Y^2}{S_X^2 / \sigma_X^2} < F_{\alpha/2, m-1, n-1}\} = 1 - \alpha$, and

$$\Pr\{F_{1-\alpha/2, m-1, n-1} < \frac{\sigma_X^2}{\sigma_Y^2} \bullet \frac{S_Y^2}{S_X^2} < F_{\alpha/2, m-1, n-1}\} = 1 - \alpha.$$

$$\text{Then, } \Pr\{(F_{1-\alpha/2, m-1, n-1}) \frac{S_X^2}{S_Y^2} < \frac{\sigma_X^2}{\sigma_Y^2} < (F_{\alpha/2, m-1, n-1}) \frac{S_X^2}{S_Y^2}\} = 1 - \alpha.$$

The values of $F_{\alpha/2, m-1, n-1}$ are given in Table 8. These tables do not explicitly give $F_{1-\alpha/2, m-1, n-1}$. One needs to use a property of the F distribution to get this value. Since

$$\Pr\{F_{1-\alpha/2, m-1, n-1} < \frac{S_Y^2 / \sigma_Y^2}{S_X^2 / \sigma_X^2}\} = 1 - \frac{\alpha}{2}, \Pr\{\frac{S_Y^2 / \sigma_Y^2}{S_X^2 / \sigma_X^2} < F_{1-\alpha/2, m-1, n-1}\} = \frac{\alpha}{2}$$

$$\Pr\{[1 / (\frac{S_Y^2 / \sigma_Y^2}{S_X^2 / \sigma_X^2})] > [1 / F_{1-\alpha/2, m-1, n-1}]\} = \frac{\alpha}{2}. \text{ That is}$$

$$\Pr\{\frac{S_X^2 / \sigma_X^2}{S_Y^2 / \sigma_Y^2} > \frac{1}{F_{1-\alpha/2, m-1, n-1}}\} = \frac{\alpha}{2}. \text{ The distribution of } \frac{S_X^2 / \sigma_X^2}{S_Y^2 / \sigma_Y^2} \text{ is a central F with } n-1$$

numerator degrees of freedom and $m-1$ denominator degrees of freedom. From the

$$\text{definition of the F percentage points in Table 8, } \Pr\{\frac{S_X^2 / \sigma_X^2}{S_Y^2 / \sigma_Y^2} > F_{\alpha/2, n-1, m-1}\} = \frac{\alpha}{2}.$$

Equating the two right hand sides of these inequalities shows that

$$F_{\alpha/2, n-1, m-1} = \frac{1}{F_{1-\alpha/2, m-1, n-1}}, \text{ or equivalently, } F_{1-\alpha/2, m-1, n-1} = \frac{1}{F_{\alpha/2, n-1, m-1}}. \text{ Then}$$

$$\Pr\{(F_{1-\alpha/2, m-1, n-1}) \frac{S_X^2}{S_Y^2} < \frac{\sigma_X^2}{\sigma_Y^2} < (F_{\alpha/2, m-1, n-1}) \frac{S_X^2}{S_Y^2}\} = 1 - \alpha \text{ reduces to}$$

$$\Pr\{(\frac{1}{F_{\alpha/2, n-1, m-1}}) \frac{S_X^2}{S_Y^2} < \frac{\sigma_X^2}{\sigma_Y^2} < (F_{\alpha/2, m-1, n-1}) \frac{S_X^2}{S_Y^2}\} = 1 - \alpha. \text{ The interval that contains } \frac{\sigma_X^2}{\sigma_Y^2} \text{ with}$$

$$\text{probability } 1 - \alpha \text{ is } (\frac{1}{F_{\alpha/2, n-1, m-1}}) \frac{S_X^2}{S_Y^2} \text{ to } (F_{\alpha/2, m-1, n-1}) \frac{S_X^2}{S_Y^2}. \text{ We use the observed sample}$$

$$\text{variances in the } 1 - \alpha \% \text{ confidence interval for } \frac{\sigma_X^2}{\sigma_Y^2} : (\frac{1}{F_{\alpha/2, n-1, m-1}}) \frac{s_X^2}{s_Y^2} \text{ to } (F_{\alpha/2, m-1, n-1}) \frac{s_X^2}{s_Y^2}.$$

Problem 4 from Chapter 7 Study Guide

A research team took a random sample of 9 observations from a normally distributed random variable Y and observed that $\bar{y}_9 = 91.2$ and $s_Y^2 = 529.6$, where \bar{y}_9 was the average of the nine observations sampled from Y and s_Y^2 was the unbiased estimate of $\text{var}(Y)$. A second research team took a random

sample of 10 observations from a normally distributed random variable X and observed that $\bar{x}_{10} = 103.5$ and $s_X^2 = 894.3$, where \bar{x}_{10} was the average of the ten observations sampled from X and s_X^2 was the unbiased estimate of $\text{var}(X)$. Find the 95% confidence interval for $\text{var}(X)/\text{var}(Y)$.

Solution: The sample variance $s_X^2 = 894.3$ is based on 9 degrees of freedom, and the sample variance $s_Y^2 = 529.6$ is based on 8 degrees of freedom. From Table 8, $F_{\alpha/2, m-1, n-1} = F_{0.025, 8, 9} = 4.10$, and $F_{\alpha/2, n-1, m-1} = F_{0.025, 9, 8} = 4.36$. The ratio $\frac{s_X^2}{s_Y^2} = \frac{894.3}{529.6} = 1.689$. The left endpoint is given by $\frac{1}{4.36} \frac{s_X^2}{s_Y^2} = 0.229 \bullet 1.689 = 0.387$

The right endpoint is given by $4.10 \frac{s_X^2}{s_Y^2} = 4.10 \bullet 1.689 = 6.92$. The 95%

confidence interval for $\text{var}(X)/\text{var}(Y)$ is from 0.387 to 6.92. Since the confidence interval for $\text{var}(X)/\text{var}(Y)$ includes 1, we would accept the null hypothesis that the ratio of the variances was 1 at the two-sided 0.05 level of significance. The confidence interval for the ratio of the variances has a right endpoint that is a factor of roughly 18 times the left endpoint. The sample averages do not enter into the solution of this problem. Some students respond incorrectly with a 95% confidence interval for the difference in means. This is not correct.

Chapter 11

Linear Regression and Correlation

The research context is that two variables have been observed for each of n participants. The research team then has a spreadsheet with n pairs of observations $(x_i, y_i), i = 1, \dots, n$. One of the variables (here y) is the outcome variable or dependent variable. This is the variable hypothesized to be affected by the other variable in scientific research. The other variable (here x) is the independent variable. It may be hypothesized to predict the outcome variable or to cause a change in the outcome variable. The research task is to document the association between independent and dependent variables. An example of a research project seeking to document a causal association would be a clinical trial in which x_i was the dosage of a medicine randomly assigned to a participant (say simvastatin) and y_i was the participant's response after a specified period taking the medicine (say cholesterol reduction after 3 months). An example of a study seeking to document the value of a predictive association would be an observational study in which x_i was the score of a statistics student on the first examination in a course and y_i was the student's score on the final examination in the course.

A recommended first step is to create the scatterplot of observations, with the vertical axis representing the dependent variable and the horizontal axis representing the independent variable. The “pencil test” is to hold up a pencil to the scatterplot and examine whether that describes the data well. If so, then it is reasonable to assume that a **linear model** (such as $\beta_0 + \beta_1 x$) describes the data. The linear model is reasonable for many data sets in observational studies. A more object procedure is to use a “nonlinear smoother” such as LOWESS to estimate the association. If the LOWESS curve is not well approximated by a line, then the assumption of linearity is not reasonable.

Estimating the Linear Model Parameters (section 11.2)

OLS (ordinary least squares) is the most used method to estimate the parameters of the linear model. An arbitrary linear model $b_0 + b_1 x$ is used as a *fit* for the dependent variable values. The method uses the *residual* $y_i - b_0 - b_1 x_i$. The fitting model is judged by how small the set of residuals is. OLS uses each residual and focuses on the magnitude of the residuals by examining the sum of squares function

$SS(b_0, b_1) = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$. The OLS method is to find the arguments $(\hat{\beta}_0, \hat{\beta}_1)$ that

make $SS(b_0, b_1)$ as small as possible. This minimization is a standard calculus problem. Step 1 is to calculate the partial derivatives of $SS(b_0, b_1)$ with respect to each argument. First, the partial with respect to b_0 :

$$\begin{aligned}\frac{\partial SS(b_0, b_1)}{\partial b_0} &= \frac{\partial}{\partial b_0} \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2 = \sum_{i=1}^n \frac{\partial}{\partial b_0} (y_i - b_0 - b_1 x_i)^2 = \sum_{i=1}^n 2(y_i - b_0 - b_1 x_i) \frac{\partial (y_i - b_0 - b_1 x_i)}{\partial b_0} \\ \frac{\partial SS(b_0, b_1)}{\partial b_0} &= \sum_{i=1}^n (-2)(y_i - b_0 - b_1 x_i).\end{aligned}$$

Second, the partial with respect to b_1 :

$$\begin{aligned}\frac{\partial SS(b_0, b_1)}{\partial b_1} &= \frac{\partial}{\partial b_1} \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2 = \sum_{i=1}^n \frac{\partial}{\partial b_1} (y_i - b_0 - b_1 x_i)^2 = \sum_{i=1}^n 2(y_i - b_0 - b_1 x_i) \frac{\partial (y_i - b_0 - b_1 x_i)}{\partial b_1} \\ \frac{\partial SS(b_0, b_1)}{\partial b_1} &= \sum_{i=1}^n (-2x_i)(y_i - b_0 - b_1 x_i).\end{aligned}$$

Step 2 is to find the arguments $(\hat{\beta}_0, \hat{\beta}_1)$ that make the two partial derivatives zero. The resulting equations are called the *normal equations*:

$$\begin{aligned}\sum_{i=1}^n (-2)(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \text{ and} \\ \sum_{i=1}^n (-2)(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)x_i &= 0.\end{aligned}$$

These equations have a very important interpretation. Let $r_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i, i = 1, \dots, n$.

The first normal equation is equivalent to $\sum_{i=1}^n r_i = 0$, and the second is $\sum_{i=1}^n r_i x_i = 0$.

That is, there are two constraints on the n residuals. The OLS residuals must sum to zero, and the OLS residuals are orthogonal to the independent variable values. The n residuals then have $n-2$ degrees of freedom.

Step 3 is to solve this two linear equation system in two unknowns. Start by using the first normal equation to solve for $\hat{\beta}_0$:

$$\begin{aligned}\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= \sum_{i=1}^n y_i - \sum_{i=1}^n \hat{\beta}_0 - \sum_{i=1}^n \hat{\beta}_1 x_i = n\bar{y}_n - n\hat{\beta}_0 - \hat{\beta}_1(n\bar{x}_n) = 0. \text{ Solving for } \hat{\beta}_0 \text{ yields} \\ \hat{\beta}_0 &= \bar{y}_n - \hat{\beta}_1 \bar{x}_n. \text{ Next, insert the solution for } \hat{\beta}_0 \text{ in the second normal equation and} \\ \text{solve for } \hat{\beta}_1 &:\end{aligned}$$

$$0 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)x_i = \sum_{i=1}^n \{[y_i - (\bar{y}_n - \hat{\beta}_1 \bar{x}_n) - \hat{\beta}_1 x_i]x_i\} = \sum_{i=1}^n [(y_i - \bar{y}_n)x_i] - \sum_{i=1}^n [\hat{\beta}_1(x_i - \bar{x}_n)x_i],$$

The solution is $\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y}_n)x_i}{\sum_{i=1}^n (x_i - \bar{x}_n)x_i}$. There are several modifications of this

formula that are helpful. The first results from noting that

$$\sum_{i=1}^n (x_i - \bar{x}_n)^2 = \sum_{i=1}^n (x_i - \bar{x}_n)x_i - \sum_{i=1}^n (x_i - \bar{x}_n)\bar{x}_n = \sum_{i=1}^n (x_i - \bar{x}_n)x_i \text{ and}$$

$$\sum_{i=1}^n (y_i - \bar{y}_n)(x_i - \bar{x}_n) = \sum_{i=1}^n (y_i - \bar{y}_n)x_i - \sum_{i=1}^n (y_i - \bar{y}_n)\bar{x}_n = \sum_{i=1}^n (y_i - \bar{y}_n)x_i. \text{ The OLS solution is}$$

then $\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y}_n)(x_i - \bar{x}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2}$. This is a very commonly quoted formula.

The second shows the relation of $\hat{\beta}_1$ and the Pearson product moment correlation. The Pearson product moment correlation is a dimensionless measure of

association. The formula is $r(x, y) = \frac{\sum_{i=1}^n (y_i - \bar{y}_n)(x_i - \bar{x}_n)}{\sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y}_n)^2}}$. The Cauchy-

Schwartz inequality shows that $|r(x, y)| \leq 1$. A correlation of +1 or -1 shows a perfect linear association. A correlation of 0 means no linear association. The numerator of $\hat{\beta}_1$ and $r(x, y)$ are the same. Starting with $\hat{\beta}_1$,

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y}_n)(x_i - \bar{x}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} = \frac{\sum_{i=1}^n (y_i - \bar{y}_n)(x_i - \bar{x}_n)}{\sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y}_n)^2}} \cdot \frac{\sqrt{\sum_{i=1}^n (y_i - \bar{y}_n)^2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2}}. \text{ That is,}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y}_n)(x_i - \bar{x}_n)}{\sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y}_n)^2}} \cdot \frac{\sqrt{\sum_{i=1}^n (y_i - \bar{y}_n)^2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2}} = r(x, y) \cdot \frac{\sqrt{(n-1)s_Y^2}}{\sqrt{(n-1)s_X^2}} = \frac{s_Y}{s_X} \cdot r(x, y). \text{ The}$$

second formula is then $\hat{\beta}_1 = \frac{s_Y}{s_X} \cdot r(x, y)$.

The next variation will be used in calculating the distributional properties of $\hat{\beta}_1$ and uses the identity that

$$\sum_{i=1}^n (y_i - \bar{y}_n)(x_i - \bar{x}_n) = \sum_{i=1}^n [(x_i - \bar{x}_n)y_i] - \sum_{i=1}^n [(x_i - \bar{x}_n)\bar{y}_n] = \sum_{i=1}^n (x_i - \bar{x}_n)y_i. \text{ Then } \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)y_i}{\sum_{i=1}^n (x_i - \bar{x}_n)^2}$$

Fisher's Decomposition of the Total Sum of Squares

The total sum of squares of the dependent variable is defined to be

$$TSS = \sum_{i=1}^n (y_i - \bar{y}_n)^2 \text{ with } n-1 \text{ degrees of freedom. The } i\text{th residual was defined above to be } r_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i, i = 1, \dots, n. \text{ After substituting for } \hat{\beta}_0, \\ r_i = y_i - \bar{y}_n - \hat{\beta}_1(x_i - \bar{x}_n), i = 1, \dots, n.$$

Fisher's decomposition is a fundamental tool for the analysis of the linear model. It

$$\text{starts with } TSS = \sum_{i=1}^n (y_i - \bar{y}_n)^2 = \sum_{i=1}^n [y_i - \bar{y}_n - \hat{\beta}_1(x_i - \bar{x}_n) + \hat{\beta}_1(x_i - \bar{x}_n)]^2 = \sum_{i=1}^n [r_i + \hat{\beta}_1(x_i - \bar{x}_n)]^2.$$

$$\text{Next } TSS = \sum_{i=1}^n [r_i + \hat{\beta}_1(x_i - \bar{x}_n)]^2 = \sum_{i=1}^n [r_i^2 + \hat{\beta}_1^2(x_i - \bar{x}_n)^2 + 2\hat{\beta}_1 r_i(x_i - \bar{x}_n)], \text{ and}$$

$$TSS = \sum_{i=1}^n r_i^2 + \sum_{i=1}^n \hat{\beta}_1^2(x_i - \bar{x}_n)^2 + 2\hat{\beta}_1 \sum_{i=1}^n r_i(x_i - \bar{x}_n). \text{ The first sum } \sum_{i=1}^n r_i^2 = SSE, \text{ the sum of}$$

squared errors and has $n-2$ degrees of freedom. The second sum $\sum_{i=1}^n \hat{\beta}_1^2(x_i - \bar{x}_n)^2$ is

called the regression sum of squares and has 1 degree of freedom. It can be simplified:

$$\text{RegSS} = \sum_{i=1}^n \hat{\beta}_1^2(x_i - \bar{x}_n)^2 = \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x}_n)^2 = [r(x, y)]^2 \left[\frac{\sqrt{\sum_{i=1}^n (y_i - \bar{y}_n)^2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2}} \right]^2 \sum_{i=1}^n (x_i - \bar{x}_n)^2 \text{ and}$$

$$\text{RegSS} = [r(x, y)]^2 \left[\frac{\sqrt{\sum_{i=1}^n (y_i - \bar{y}_n)^2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2}} \right]^2 \sum_{i=1}^n (x_i - \bar{x}_n)^2 = [r(x, y)]^2 \sum_{i=1}^n (y_i - \bar{y}_n)^2 = [r(x, y)]^2 TSS.$$

$$\text{Finally, the third sum } 2\hat{\beta}_1 \sum_{i=1}^n r_i(x_i - \bar{x}_n) = 2\hat{\beta}_1 \left(\sum_{i=1}^n r_i x_i - \sum_{i=1}^n r_i \bar{x}_n \right) = 2\hat{\beta}_1(0 - 0) = 0.$$

This is conventionally displayed in an Analysis of Variance Table as below:

Analysis of Variance Table One Predictor Linear Regression

Source	DF	Sum of Squares	Mean Square	F
Regression	1	$\sum_{i=1}^n \hat{\beta}_1^2 (x_i - \bar{x}_n)^2 = [r(x, y)]^2 TSS$	$[r(x, y)]^2 TSS$	$\frac{(n-2)[r(x, y)]^2}{1 - [r(x, y)]^2}$
Error	$n-2$	$\sum_{i=1}^n r_i^2 = \{1 - [r(x, y)]^2\} TSS$	$\frac{\{1 - [r(x, y)]^2\} TSS}{(n-2)}$	
Total	$n-1$	$TSS = (n-1)s_{DV}^2$		

11.3 Inferences

There must be a probabilistic model for the data so that researchers can make inferences and find confidence intervals. The model for one predictor linear regression is $Y_i = \beta_0 + \beta_1 x_i + \sigma_{Y|x} Z_i$. The outcome or dependent (random) variables $Y_i, i = 1, \dots, n$ are each assumed to be the sum of the linear regression expected value $\beta_0 + \beta_1 x_i$ and a random error term $\sigma_{Y|x} Z_i$. The random variables $Z_i, i = 1, \dots, n$ are assumed to be independent standard normal random variables. The parameter β_0 is the intercept parameter and is fixed but unknown. The parameter β_1 is the slope parameter and is also fixed but unknown. This parameter is the focus of the statistical analysis. The parameter $\sigma_{Y|x}$ is also fixed but unknown. Another description of this model is that $Y_i, i = 1, \dots, n$ are independent normally distributed random variables with Y_i having the distribution $N(\beta_0 + \beta_1 x_i, \sigma_{Y|x}^2)$. That is, $E(Y_i | X = x_i) = \beta_0 + \beta_1 x_i$, and $\text{var}(Y_i | X = x_i) = \sigma_{Y|x}^2$. The assumption that $\text{var}(Y_i | X = x_i) = \sigma_{Y|x}^2$ is called the *homoscedasticity* assumption.

There are four assumptions. There are two important assumptions: the outcome variables $Y_i, i = 1, \dots, n$ are independent and $E(Y_i | X = x_i) = \beta_0 + \beta_1 x_i$ for $i = 1, \dots, n$. Homoscedasticity is less important. The assumption that $Y_i, i = 1, \dots, n$ are normally distributed random variables is least important.

Variance Calculations

The most complex variance formula in this course so far is:

$$\text{var}(aX + bY) = a^2 \text{var } X + b^2 \text{var } Y + 2ab \text{cov}(X, Y).$$

More complex calculations are required for the variance-covariance matrix of the OLS estimates. The easiest way is to use the variance-covariance matrix of a random vector. Let Y be an $n \times 1$ vector of random variables $(Y_1, Y_2, \dots, Y_n)^T$. That is, each component of the vector is a random variable. Then the expected value of vector Y is the $n \times 1$ vector whose components are the respective means of the random variables; that is, $E(Y) = (EY_1, EY_2, \dots, EY_n)^T$. The variance-covariance matrix of the random vector Y is the $n \times n$ matrix whose diagonal entries are the respective variances of the random variables and whose off-diagonal elements are the covariances of the random variables. That is,

$$\text{vcv}(Y) = \begin{bmatrix} \text{var}(Y_1) & \text{cov}(Y_1, Y_2) & \cdots & \text{cov}(Y_1, Y_n) \\ \text{cov}(Y_2, Y_1) & \text{var}(Y_2) & \cdots & \text{cov}(Y_2, Y_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(Y_n, Y_1) & \text{cov}(Y_n, Y_2) & \cdots & \text{var}(Y_n) \end{bmatrix}.$$

In terms of expectation operator calculations, $\text{vcv}(Y) = E[(Y - EY)(Y - EY)^T] = \Sigma$.

Variance of a Set of Linear Combinations

Let W be the $m \times 1$ random vector of linear combinations of Y given by $W = MY$, where M is a matrix of constants having m rows and n columns. Then

$E(W) = E(MY) = ME(Y)$, The definition of the variance-covariance matrix of W is

$\text{vcv}(W) = E[(W - EW)(W - EW)^T] = E[(MY - MEY)(MY - MEY)^T]$, and

$\text{vcv}(W) = E[(MY - MEY)(MY - MEY)^T] = E\{M(Y - EY)[M(Y - EY)]^T\}$

From matrix algebra, when A is an $n \times m$ matrix and B is an $m \times p$ matrix, then

$(AB)^T = B^T A^T$. Then $[M(Y - EY)]^T = (Y - EY)^T M^T$, and

$\text{vcv}(W) = \text{vcv}(MY) = E\{M(Y - EY)(Y - EY)^T M^T\} = M\{E[(Y - EY)(Y - EY)^T]\}M^T$ from the

linear operator property of E . Since $\text{vcv}(Y) = E[(Y - EY)(Y - EY)^T] = \Sigma$,

$\text{vcv}(W) = \text{vcv}(MY) = M \times \text{vcv}(Y) \times M^T = M\Sigma M^T$

Examples

The first use of this result is to find the variance of a linear combination of values from Y , an $n \times 1$ vector of random variables. Let a be an $n \times 1$ vector of constants, and let $W = a^T Y$. Then $\text{var}(a^T Y) = a^T \times \text{vcv}(Y) \times (a^T)^T = a^T \times \text{vcv}(Y) \times a$. This is the completely general form of $\text{var}(aX + bY) = a^2 \text{var}(X) + b^2 \text{var}(Y) + 2ab \text{cov}(X, Y)$.

The second example is fundamental to this chapter. The OLS estimates of the parameters are always the same functions of the observed data: $\hat{\beta}_0 = \bar{y}_n - \hat{\beta}_1 \bar{x}_n$ and

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}_n) y_i}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} .$$

It is then reasonable to study the random variables

$$\bar{Y}_n = \frac{\sum_{i=1}^n Y_i}{n} = \frac{1}{n} Y_1 + \frac{1}{n} Y_2 + \cdots + \frac{1}{n} Y_n$$

and

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}_n) Y_i}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} = \frac{(x_1 - \bar{x}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} Y_1 + \frac{(x_2 - \bar{x}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} Y_2 + \cdots + \frac{(x_n - \bar{x}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} Y_n .$$

Let $w_i = \frac{(x_i - \bar{x}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2}, i = 1, \dots, n$. Then $\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}_n) Y_i}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} = w_1 Y_1 + w_2 Y_2 + \cdots + w_n Y_n$.

Let $\begin{pmatrix} \bar{Y}_n \\ \hat{\beta}_1 \end{pmatrix} = \begin{bmatrix} 1/n & 1/n & \cdots & 1/n \\ w_1 & w_2 & \cdots & w_n \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$, which has the form MY , where

$M = \begin{bmatrix} 1/n & 1/n & \cdots & 1/n \\ w_1 & w_2 & \cdots & w_n \end{bmatrix}$. In the model $Y_i = \beta_0 + \beta_1 x_i + \sigma_{Y|x} Z_i$, $\text{vcv}(Y) = \sigma_{Y|x}^2 I_{n \times n}$. Then

$$\text{vcv} \begin{pmatrix} \bar{Y}_n \\ \hat{\beta}_1 \end{pmatrix} = \begin{bmatrix} 1/n & 1/n & \cdots & 1/n \\ w_1 & w_2 & \cdots & w_n \end{bmatrix} \times \sigma_{Y|x}^2 I_{n \times n} \times \begin{bmatrix} 1/n & w_2 \\ \vdots & \vdots \end{bmatrix} = \sigma_{Y|x}^2 \begin{bmatrix} 1/n & 1/n & \cdots & 1/n \\ w_1 & w_2 & \cdots & w_n \end{bmatrix} \times \begin{bmatrix} 1/n & w_2 \\ \vdots & \vdots \end{bmatrix} .$$

Then,

$$\text{vcv} \begin{pmatrix} \bar{Y}_n \\ \hat{\beta}_1 \end{pmatrix} = \sigma_{Y|x}^2 \begin{bmatrix} 1/n & 1/n & \cdots & 1/n \\ w_1 & w_2 & \cdots & w_n \end{bmatrix} \times \begin{bmatrix} 1/n & w_2 \\ \vdots & \vdots \end{bmatrix} = \sigma_{Y|x}^2 \begin{bmatrix} \sum_{i=1}^n \frac{1}{n^2} & \sum_{i=1}^n \frac{w_i}{n} \\ \sum_{i=1}^n \frac{w_i}{n} & \sum_{i=1}^n w_i^2 \end{bmatrix} .$$

Now $\sum_{i=1}^n \frac{w_i}{n} = \sum_{i=1}^n \frac{1}{n} \frac{(x_i - \bar{x}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} = \frac{1}{n \sum_{i=1}^n (x_i - \bar{x}_n)^2} \sum_{i=1}^n (x_i - \bar{x}_n) = 0$, and

$$\sum_{i=1}^n w_i^2 = \sum_{i=1}^n \left[\frac{(x_i - \bar{x}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \right]^2 = \frac{1}{\left[\sum_{i=1}^n (x_i - \bar{x}_n)^2 \right]^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2 = \frac{1}{\sum_{i=1}^n (x_i - \bar{x}_n)^2}.$$

The final result is that $\text{vcv} \begin{pmatrix} \bar{Y}_n \\ \hat{\beta}_1 \end{pmatrix} = \sigma_{Y|x}^2 \begin{bmatrix} \sum_{i=1}^n \frac{1}{n^2} & \sum_{i=1}^n \frac{w_i}{n} \\ \sum_{i=1}^n \frac{w_i}{n} & \sum_{i=1}^n w_i^2 \end{bmatrix} = \begin{bmatrix} \frac{\sigma_{Y|x}^2}{n} & 0 \\ 0 & \frac{\sigma_{Y|x}^2}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \end{bmatrix}$

To summarize this result, $\text{var}(\bar{Y}_n) = \frac{\sigma_{Y|x}^2}{n}$, $\text{var}(\hat{\beta}_1) = \frac{\sigma_{Y|x}^2}{\sum_{i=1}^n (x_i - \bar{x}_n)^2}$, and

$$\text{cov}(\bar{Y}_n, \hat{\beta}_1) = 0.$$