### EXAMPLE 12.15

Refer to the output given in Example 12.14.

**a.** Test $H_0: \beta_1 = 0$ versus $H_a: \beta_1 \neq 0$ at the $\alpha = .10$ level.
**b.** Is the conclusion of the test compatible with the confidence interval?

**Solution**

**a.** The test statistic for $H_0: \beta_1 = 0$ versus $H_a: \beta_1 \neq 0$ is

$$t = \frac{\hat{\beta}_1}{s_{\hat{\beta}_1}} = \frac{.01291}{.00283} = 4.562$$

The .05 upper percentile for the $t$ distribution with df $= 54 - (4 + 1) = 49$ is 1.677. Because the computed value of the test statistic is greater than the tabled value, we conclude there is significant evidence to reject $H_0$. Thus, $x_1$ has additional predictive power in the presence of the other three explanatory variables.

**b.** The 90% confidence interval for $\beta_1$ did not include 0, which indicates that $H_0: \beta_1 = 0$ should be rejected at the $\alpha = .10$ level.

### EXAMPLE 12.16

Refer to Example 12.12. Locate the $t$ statistic for testing $H_0: \beta_3 \leq 0$ versus $H_a: \beta_3 > 0$ in the output given in Example 12.12. Do the data support $H_a: \beta_3 > 0$ at any of the usual values for $\alpha$?

**Solution** The $t$ statistics are shown under the heading *STUDENT'S T*. For $x_3$ (INCOME), the $t$ statistic is 2.62, which is computed as .26528/.10127. With df $= 17$, the tabled values from the $t$ distribution are 2.576 and 2.898 for $\alpha = .01$ and .005, respectively. Thus, $H_0$ would be rejected at the $\alpha = .01$ level but not at the $\alpha = .005$ level.

The output lists a $p$-value under the column heading $P$. This $p$-value is for a two-sided alternative hypothesis, $H_a: \beta_3 \neq 0$. The $p$-value for the 1-sided alternative $H_a: \beta_3 > 0$ is given by $p$-value $= Pr(t_{17} > 2.62) = .00896 < .01 = \alpha$.

The multiple regression $F$ and $t$ tests that we discuss in this chapter test different null hypotheses. It sometimes happens that the $F$ test results in the rejection of $H_0: \beta_1 = \beta_2 = \cdots = \beta_k = 0$, whereas no $t$ test of $H_0: \beta_j = 0$ is significant. In such a case, we can conclude that there is predictive value in the equation as a whole, but we cannot identify the specific variables that have predictive value. Remember that each $t$ test is testing the unique predictive value. Does this variable add predictive value, given all the other predictors? When two or more predictor variables are highly correlated among themselves, it often happens that no $x_j$ can be shown to have significant, unique predictive value, even though the $x$s together have been shown to be useful. If we are trying to predict housing sales based on gross domestic product and disposable income, we probably cannot prove that GDP adds value given DI, or that DI adds value given GDP.

## 12.5   Testing a Subset of Regression Coefficients

*F test for several $\beta_j$s*

In the last section, we presented an $F$ test for testing *all* the coefficients in a regression model and a $t$ test for testing *one* coefficient. Another $F$ test of the null hypothesis tests that *several* of the true coefficients are zero—that is, that several of the predictors have no value given the others. For example, if we try to predict the prevailing wage rate in

various geographical areas for clerical workers based on the national minimum wage, national inflation rate, population density in the area, and median apartment rental price in the area, we might well want to test whether the variables related to area (density and apartment price) added anything, given the national variables.

A null hypothesis for this situation would say that the true coefficients of density and apartment price were zero. According to this null hypothesis, these two independent variables together have no predictive value once minimum wage and inflation are included as predictors.

The idea is to compare the SS(Regression) or $R^2$ values when density and apartment price are excluded and when they are included in the prediction equation. When they are included, the $R^2$ is automatically at least as large as the $R^2$ when they are excluded because we can predict at least as well with more information as with less. Similarly, SS(Regression) will be larger for the complete model. The $F$ test for this null hypothesis tests whether the gain is more than could be expected by chance alone. In general, let $k$ be the total number of predictors, and let $g$ be the number of predictors with coefficients not hypothesized to be zero ($g < k$). Then $k - g$ represents the number of predictors with coefficients that are hypothesized to be zero. The idea is to find SS(Regression) values using all predictors (the **complete model**) and using only the $g$ predictors that do not appear in the null hypothesis (the **reduced model**). Once these have been computed, the test proceeds as outlined next. The notation is easier if we assume that the reduced model contains $\beta_1, \beta_2, \ldots, \beta_g$, so that the variables in the null hypothesis are listed last.

**complete and reduced models**

**F Test of a Subset of Predictors**

$H_0: \beta_{g+1} = \beta_{g+2} = \cdots = \beta_k = 0$    $k - g$

$H_a: H_0$ is not true.

T.S.: $F = \dfrac{[\text{SS(Regression, complete)} - \text{SS(Regression, reduced)}]/(k - g)}{\text{SS(Residual, complete)}/[n - (k + 1)]}$

R.R.: $F > F_\alpha$, where $F_\alpha$ cuts off a right-tail of area $\alpha$ of the $F$ distribution with $df_1 = (k - g)$ and $df_2 = [n - (k + 1)]$.

Check assumptions and draw conclusions.

### EXAMPLE 12.17

A state fisheries commission wants to estimate the number of bass caught in a given lake during a season in order to restock the lake with the appropriate number of young fish. The commission could get a fairly accurate assessment of the seasonal catch by extensive "netting sweeps" of the lake before and after a season, but this technique is much too expensive to be done routinely. Therefore, the commission samples a number of lakes and records $y$, the seasonal catch (thousands of bass per square mile of lake area); $x_1$, the number of lakeshore residences per square mile of lake area; $x_2$, the size of the lake in square miles; $x_3 = 1$ if the lake has public access, 0 if not; and $x_4$, a structure index. (Structures are weed beds, sunken trees, drop-offs, and other living places for bass.) The data are shown in Table 12.13.

The commission is convinced that residences and size are important variables in predicting catch because they both reflect how intensively the lake has been

ıum wage, ıent rental ɔd to area . ficients of ɔsis, these num wage and apart- equation. when they on as with ɔ F test for by chance number of represents ɔ. The idea ) and using **ed model**), notation is at the vari- (k − g) ıtion ɔaught in a ppropriate essment of ınd after a Therefore, onal catch ɔshore resi- iles; $x_3$ = 1 ıctures are ɔe data are ıt variables ɔ has been

**TABLE 12.13** Bass catch data

| Lake | Catch | Residence | Size | Access | Structure |
|---|---|---|---|---|---|
| 1 | 3.6 | 92.2 | .21 | 0 | 81 |
| 2 | .8 | 86.7 | .30 | 0 | 26 |
| 3 | 2.5 | 80.2 | .31 | 0 | 52 |
| 4 | 2.9 | 87.2 | .40 | 0 | 64 |
| 5 | 1.4 | 64.9 | .44 | 0 | 40 |
| 6 | .9 | 90.1 | .56 | 0 | 22 |
| 7 | 3.2 | 60.7 | .78 | 0 | 80 |
| 8 | 2.7 | 50.9 | 1.21 | 0 | 60 |
| 9 | 2.2 | 86.1 | .34 | 1 | 30 |
| 10 | 5.9 | 90.0 | .40 | 1 | 90 |
| 11 | 3.3 | 80.4 | .52 | 1 | 74 |
| 12 | 2.9 | 75.0 | .66 | 1 | 50 |
| 13 | 3.6 | 70.0 | .78 | 1 | 61 |
| 14 | 2.4 | 64.6 | .91 | 1 | 40 |
| 15 | .9 | 50.0 | 1.10 | 1 | 22 |
| 16 | 2.0 | 50.0 | 1.24 | 1 | 50 |
| 17 | 1.9 | 51.2 | 1.47 | 1 | 37 |
| 18 | 3.1 | 40.1 | 2.21 | 1 | 61 |
| 19 | 2.6 | 45.0 | 2.46 | 1 | 39 |
| 20 | 3.4 | 50.0 | 2.80 | 1 | 53 |

fished. However, the commission is uncertain whether access and structure are useful as additional predictor variables. Therefore, two regression models (with all four predictor variables entered linearly) are fitted to the data, the first model with all four variables and the second model without access and structure. The relevant portions of the Minitab output follow:

```
Full Model:

Regression Analysis: catch versus residenc, size, access, structur

The regression equation is
catch = - 2.78 + 0.0268 residenc + 0.504 size + 0.743 access + 0.0511 structur

Predictor      Coef     SE Coef       T       P
Constant    -2.7840      0.8157   -3.41   0.004
residenc   0.026794    0.009141    2.93   0.010
size         0.5035      0.2208    2.28   0.038
access       0.7429      0.2021    3.68   0.002
structur   0.051129    0.004542   11.26   0.000

S = 0.389498   R-Sq = 91.4%   R-Sq(adj) = 89.1%

Analysis of Variance

Source           DF        SS       MS       F       P
Regression        4   24.0624   6.0156   39.65   0.000
Residual Error   15    2.2756   0.1517
Total            19   26.3380
```

```
              Reduced Model:

              Regression Analysis: catch versus residenc, size

              The regression equation is
              catch = - 0.87 + 0.0394 residenc + 0.828 size

              Predictor      Coef      SE Coef       T       P
              Constant     -0.871       2.409      -0.36   0.722
              residenc     0.03941      0.02733     1.44    0.168
              size         0.8280       0.6372      1.30    0.211

              S = 1.17387   R-Sq = 11.1%   R-Sq(adj) = 0.6%

              Analysis of Variance

              Source          DF      SS       MS      F       P
              Regression       2     2.913    1.456   1.06    0.369
              Residual Error  17    23.425    1.378
              Total           19    26.338
```

a. Write the complete and reduced models.
b. Write the null hypothesis for testing that the omitted variables have no (incremental) predictive value.
c. Perform an $F$ test for this null hypothesis.

**Solution**

a. The complete and reduced models are, respectively,

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \varepsilon_i$$

and

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i$$

The corresponding multiple regression least-squares equations based on the sample data are

Complete: $\hat{y} = -2.78 + .0268x_1 + .504x_2 + .743x_3 + .0511x_4$
Reduced: $\hat{y} = -.87 + .0394x_1 + .828x_2$

b. The appropriate null hypothesis of no predictive power for $x_3$ and $x_4$ is $H_0: \beta_3 = \beta_4 = 0$.

c. The test statistic for the $H_0$ of part (b) makes use of SS(Regression, complete) = 24.0624, SS(Regression, reduced) = 2.913, SS(Residual, complete) = 2.2756, $k = 4$, $g = 2$, and $n = 20$:

$$\text{T.S.:} \quad F = \frac{[\text{SS(Regression, complete)} - \text{SS(Regression, reduced)}]/(4 - 2)}{\text{SS(Residual, complete)}/(20 - 5)}$$

$$= \frac{(24.0624 - 2.913)/2}{2.2756/(20 - 5)} = 69.705$$

The tabled value $F_{.01}$ for 2 and 15 df is 6.36. The value of the test statistic is much larger than the tabled value, so we have conclusive evidence that the access and structure variables add predictive value ($p < .0001$).