A multiple linear regression that describes the variation in the variable y in terms of the k (explanatory) variables  $x_1, x_2, \ldots, x_k$  has the form

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon.$$

- $\beta_0, \beta_1, \ldots, \beta_k$  are population parameters.
- Assumptions about the model include
  - (i) The error  $\varepsilon$  has a normal distribution with mean 0 for each fixed set of values  $x_1 = \xi_1, x_2 = \xi_2, \dots, x_k = \xi_k$ .
  - (ii) The error is independent of the  $x_j$  s, so the standard deviation of  $\varepsilon$  is fixed. I.e.,  $\varepsilon \sim N(0, \sigma^2)$ . (Homoskedasticity)
- These assumptions imply that

$$E(y|x_1 = \xi_1, x_2 = \xi_2, \dots, x_k = \xi_k) = \beta_0 + \beta_1 \xi_1 + \beta_2 \xi_2 + \dots + \beta_k \xi_k,$$

where  $E(y|x_1 = \xi_1, x_2 = \xi_2, \dots, x_k = \xi_k)$  is the expected (mean) y value for all observations satisfying  $x_1 = \xi_1, x_2 = \xi_2, \dots, x_k = \xi_k$ .

Using sample data, we compute the estimated regression equation

$$\hat{y}_i = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_k x_k.$$

• The sample coefficients  $b_0, b_1, \ldots, b_k$  are sample statistics (that estimate the population coefficients  $\beta_0, \ldots, \beta_k$ ). As such,  $b_0, \ldots, b_k$  are all random variables, and the assumptions above imply that they all have *normal* distributions. In particular...

... for each i between 0 and k:

$$\frac{b_i - \beta_i}{SE(b_i)} \sim t$$
-distribution, with  $n - (k+1) d.f.$ 

We can use this to test individual coefficients for statistical significance:

$$H_0: \beta_i = 0$$
 vs.  $H_a: \beta_i \neq 0.$ 

Rejecting  $H_0$  means that there is significant linear correlation between  $x_i$  and y, and  $b_i$  is a reliable measure of the marginal change in y for a one-unit change in  $x_i$ , assuming that all other variables are held fixed. We can also test the *overall significance* of the regression:

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

 $H_a: \beta_j \neq 0$ , for at least one j between 1 and k.

The test statistic is

$$F^* = \frac{MSS_m}{MSS_e} = \frac{(\sum(\hat{y}_i - \overline{y})^2)/k}{(\sum(y_i - \hat{y}_i)^2)/(n - k + 1)}$$

which follows the F-distribution with k numerator d.f. and n-k+1 denominator d.f.

 $H_0$  is rejected at the  $\alpha$  level of significance if  $F^* > F_{\alpha}$  (the critical *F*-value). Alternatively, reject  $H_0$  if the *P*-value,  $Prob(F > F^*)$ , is smaller than  $\alpha$ .

## **Comments:**

- 1. For a simple regression  $\overline{y}(x) = \beta_0 + \beta_1 x$ , testing  $H_0: \beta_1 = 0$  (t-test) yields the same conclusion as the F-test for overall significance.
- 2. Most software packages (including JMP) compute the *t*-scores for all of the regression coefficients and their *p*-values, as well as the *F*-score of overall significance and its *p*-value.
- 3. One can also test other hypotheses about individual coefficients, e.g.,

$$H_0: \beta_3 = 2 \text{ vs. } H_a: \beta_3 < 2.$$

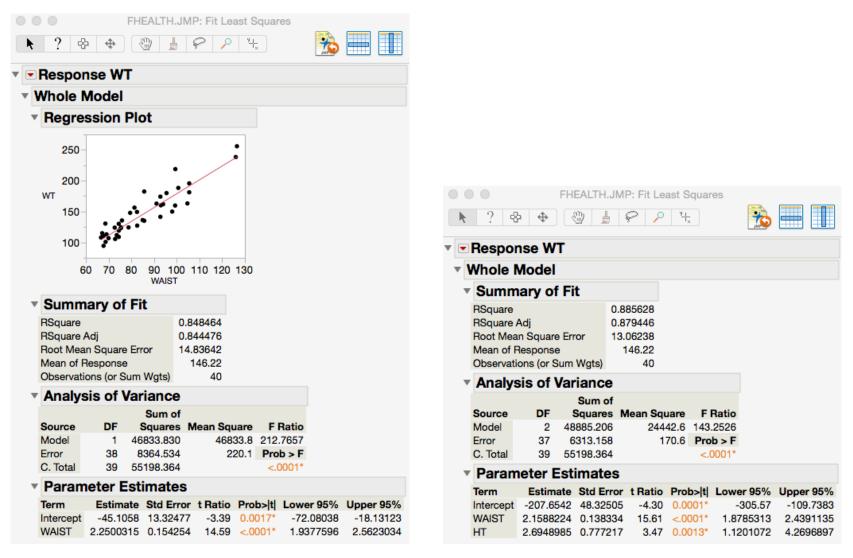
With this null hypothesis,

$$\frac{b_3-2}{SE(b_3)} \sim t$$
-distribution, with  $n - (k+1) d.f.$ ,

and the mechanics of the test are the same as any other t-test.

**Example:** Estimating weight with a measuring tape.

Response variable: weight (lbs). Effect variables: waist circumfrence (cm), height (inches), age (years), arm circumfrence (cm).



► ? €			P /		ares	<b>%</b>		k	? 4		00	.JMP: Fit L	east So	quares		
Responsible	nse WT							• •	Respor	nse W1	Г					
Whole I	Model							-	Nhole N	Nodel						
Summary of Fit									Summ		Ei+					
RSquare Adj 0.879 Root Mean Square Error 13.0		0.888886 0.879626 13.0526 146.22 40	6 6 2					RSquare RSquare Adj Root Mean Square Error Mean of Response Observations (or Sum Wgts)		0.934142 0.928654 10.04885 146.22 s) 40						
<ul> <li>Analysis of Variance</li> </ul>								Analys	sis of V	/arianc	е					
Source Model Error	DF Sc 3 490 36 61	65.035 33.329		355.0 9 70.4 P	F Ratio 95.9969 rob > F				Source Model Error	<b>DF</b> 3	Sum	of es Mean S 03 17	187.7	F Ratio 170.2098 Prob > F		
C. Total		98.364		•	<.0001*				C. Total	39	55198.36	64		<.0001*		
<ul> <li>Param</li> </ul>	neter Estin							•	Param	neter E	stimat	es				
Term Intercept WAIST HT AGE		48.3395 ).18261 ).77667	6 -4.34 1 12.49 7 3.46	0.0001 <.0001 0.0014	-307 1.911 1.111	r 95% 7.9658 10662 13492 82158	Upper 95% -111.8915 2.6517691 4.2616984 0.2234156		Term Intercept WAIST HT ARM	-278.45 1.057 3.6014	ate         Std E           36         39.63           52         0.238           71         0.623           21         0.689	704 -7.0 874 4.4 289 5.7	3 <.00 3 <.00 8 <.00	01* -35 01* 0.5 01* 2.3	er 95% 58.8413 730619 373827 538359	Upper 95% -198.066 1.5419781 4.8655593 4.9525283