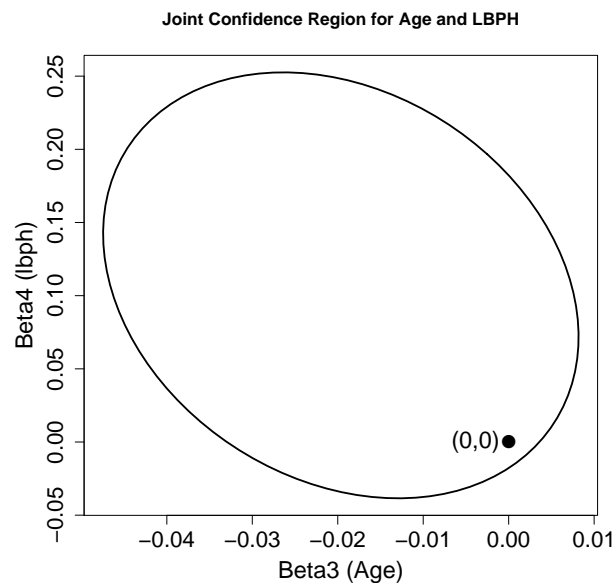


1. Faraway 3.1. For the `prostate` data, fit a model with `lpsa` as the response and the other variables as predictors.

(a) Compute 90 and 95% CIs for the parameter associated with `age`. Using just these intervals, what could we have deduced about the p-value for `age` in the regression summary?

The 90% CI for the age coefficient is $(-0.0382, -0.0011)$, and the 95% CI is $(-0.0418, 0.0026)$. Since 0 is in the first but not the second, we conclude that a two-sided t -test for $H_0 : \beta_{age} = 0$ is between .05 and .1 providing moderate evidence of significance.

(b) Compute and display a 95% joint confidence region for the parameters associated with `age` and `lbph`. Plot the origin on this display. The location of the origin on the display tells us the outcome of a certain hypothesis test. State that test and its outcome.



Whether the origin is within the confidence region tests the joint significance of age and lbph in the model - specifically, $H_0 : \beta_3 = \beta_4 = 0$. Since $(0, 0)$ is within the region, we conclude (at $\alpha = .05$ level) that there is no statistical evidence to suggest that age and lbph are jointly significant in the model.

(c) Suppose a new patient with the following values arrives:

lcavol	lweight	age	lbph	svi	lcp	gleason	pgg45
1.44692	3.62301	65.00000	0.30010	0.00000	-0.79851	7.00000	15.00000

Predict the `lpsa` for this patient along with an appropriate 95% CI.

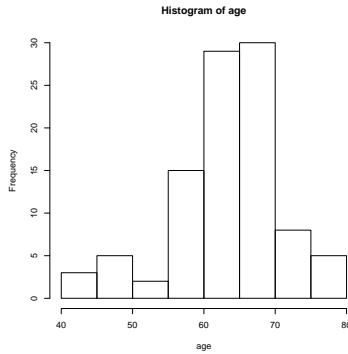
Let \mathbf{x}_0 be the $1 \times p$ vector of values given in the table above (with a one prepended). Then the 95% CI is given by

$$\mathbf{x}_0 \hat{\boldsymbol{\beta}} \pm t_{.975}(n-p) \sqrt{\text{MSE} \cdot \mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0} \approx 2.389 \pm 0.109 \text{ or } (2.299, 2.479).$$

(d) Repeat the last question for a patient with the same values except that he or she is age 20. Explain why the CI is wider.

In this case, let \mathbf{x}_1 denote the explanatory variables, then

$$\mathbf{x}_1 \hat{\boldsymbol{\beta}} \pm t_{.975}(n-p) \sqrt{\text{MSE} \cdot \mathbf{x}'_1 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_1} \approx 3.273 \pm 0.509 \text{ or } (1.880, 2.898).$$



The increase in uncertainty is due to the lack of data points for ages near 20. This is evident in the histogram of ages to the left.

- (e) In the text, we made a permutation test corresponding to the F -test for the significance of all the predictors. Execute the permutation test corresponding to the t -test for age in this model. (Hint: `summay(g)$coef[4,3]` gets you the t -statistic you need if the model is galled `g`.)
2. Faraway 3.3. Using the `teengamb` data, fit a model with `gamble` as the response and the other variables as predictors.
- Which variables are statistically significant?
 - What interpretation should be given to the coefficient for `sex`?
 - Predict the amount that a male with average (given these data) status, income, and verbal score would gamble along with an appropriate 95% CI. Repeat the prediction for a male with maximal values (for this data) of status, income and verbal score. Which CI is wider and why is this result expected?
 - Fit a model with just `income` as a predictor and use an F -test to compare it to the full model.
3. Faraway 3.4. Using the `sat` data:
- Fit a model with `total` sat score as the response and `expend`, `ratio` and `salary` as predictors. Test the hypothesis that $\beta_{salary} = 0$. Test the hypothesis that $\beta_{salary} = \beta_{ratio} = \beta_{expend} = 0$. Do any of these predictors have an effect on the response?
 - Now add `takers` to the model. Test the hypothesis that $\beta_{takers} = 0$. Compare this model to the previous one using an F -test. Demonstrate that the F -test and t -test here are equivalent.
4. For each of the following regression models, indicate whether it is a general linear regression model. If it is not, state wheter it can be expressed in the general linear model form:

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_{p-1} x_{i,p-1} + \epsilon_i$$

(a) $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 \log_{10} x_{i2} + \beta_3 x_{i1}^2 + \epsilon_i$

This is a linear regression model.

(b) $y_i = \epsilon_i \exp(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}^2)$

If we log-transform the response, then we have the general linear regression model

$$\log y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}^2 + \log \epsilon_i.$$

(c) $y_i = \beta_0 + \log_{10}(\beta_1 x_{i1}) + \beta_2 x_{i2} + \epsilon_i$

This problem is *not* linear due to the term $\log_{10}(\beta_1 x_{i1})$, and there is no obvious transformation that will make the model linear.

$$(d) \quad y_i = \beta_0 \exp(\beta_1 x_{i1}) + \epsilon_i$$

This problem is also *not* linear due to the term $\beta_0 \exp(\beta_1 x_{i1})$, and a log transformation will not work as in (b) since the error is additive.

$$(e) \quad y_i = [1 + \exp(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}^2 + \epsilon_i)]^{-1}$$

Consider the transformation $g(y) = (\log y - 1)^{-1}$. Then, the transformed model is

$$g(y_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}^2 + \epsilon_i,$$

which is a general linear regression model.

5. The information below relates a response variable y , a second measurement on wood volume, to four explanatory variables defined as: x_1 = a first measurement on wood volume, x_2 = the number of trees, x_3 = the average age of trees, and x_4 = the average volume per tree. Note that $x_4 = x_1/x_2$. Some of the information in the coefficients table and Analysis of Variance table below has not been reported, so that you can figure it out on your own.

Coefficients Table					Analysis of Variance Table					
Predictor	$\hat{\beta}_k$	SE($\hat{\beta}_k$)	t	p-val	Source	df	SS	MS	F	p-val
Intercept	23.45	14.90		0.122	Regression	4	887994			0.000
x_1	0.93209	0.08602		0.000	Error					
x_2		0.4721	1.5554	0.126	Total		902773	16718.02		
x_3	-0.4982	0.1520		0.002	Sequential Sums of Squares					
x_4	3.486	2.274		0.132	Source	df	Seq. SS			
					x_1	1	883880			
					x_2	1	183			
					x_3	1	3237			
					x_4	1	694			

(a) How many observations are in these data?

$$n = \text{TSS}/\text{MST} + 1 = (902773/16718.02) + 1 = 55.$$

(b) What is R^2 for this model?

R^2 is given by the regression sum of squares (SSReg) over the total sum of squares (TSS). From the partial ANOVA table we calculate

$$\frac{\text{SSReg}}{\text{TSS}} = \frac{887994}{902773} \approx 0.9836.$$

(c) What is the mean squared error?

$$\text{MSE} = \frac{\text{RSS}}{n - p} = \frac{\text{TSS} - \text{SSReg}}{n - p} = \frac{902773 - 887994}{55 - 5} = 295.58$$

(d) Give a 95% confidence interval for β_2 .

The t -statistic in the Coefficients Table tests if the coefficient is 0, hence, is given by

$$t^* = \frac{\hat{\beta}_2}{\text{SE}(\hat{\beta}_2)} \iff \hat{\beta}_2 = \text{SE}(\hat{\beta}_2) \cdot t^* \approx 0.4721 \cdot 1.5554 \approx 0.734.$$

The confidence interval (with Bonferroni correction) is of the form $\hat{\beta}_2 \pm t_{1-0.5/(2.5)}(50) \cdot \text{SE}(\hat{\beta}_2)$. This calculation results in $(-0.530, 1.998)$.

(e) Test the null hypothesis $H_0 : \beta_3 = 0$.

The fourth row of the Coefficients Table tests this hypothesis, and the reported p-value is 0.002 which gives fairly strong evidence that the coefficient is 0.

(f) Test the null hypothesis $H_0 : \beta_1 = 1$. Why might this hypothesis be of interest?

We calculate

$$t^* = \frac{\hat{\beta}_1 - 1}{\text{SE}(\hat{\beta}_1)} \approx -0.0789$$

and calculate the p-value for two-sided test $p = 2P(t \leq t^*) \approx 0.937$. Hence, we conclude that there is no evidence to suggest that $\beta_1 = 1$. If we have a practical reason to believe that the response increases exactly as x_1 increases in the model, then there is not evidence to the contrary of this assumption.

(g) Give the F-statistic for testing $H_0 : \beta_3 = 0$ relative to the full model.

It can be shown [Renchner 204-205 and HW1.8a] that the F -statistic is given by squaring the associated t -statistic from the Coefficients Table. I.e.

$$F = t^2 = \left(\frac{\hat{\beta}_3}{\text{SE}(\hat{\beta}_3)} \right)^2 = \left(\frac{-0.4982}{0.1520} \right)^2 = 10.74287.$$

(h) Give $R(\beta_3|\beta_1, \beta_2)$ and find $R(\beta_3|\beta_1, \beta_2, \beta_4)$.

From the Sequential Sums of Squares table, we can read directly

$$R(\beta_3|\beta_1, \beta_2) = \text{SSReg}(x_3|x_1, x_2) = 3237.$$

On the other hand, we can recover $R(\beta_3|\beta_1, \beta_2, \beta_4)$ from the F -statistic calculated in (g). That is,

$$F = \frac{R(\beta_3|\beta_1, \beta_2, \beta_4)}{\text{MSE}} \iff R(\beta_3|\beta_1, \beta_2, \beta_4) = \text{MSE} \cdot F \approx 295.58 \cdot 10.74287 \approx 3175.377$$

(i) Test the model with only variables x_1 and x_2 against the model with all of the variables x_1, x_2, x_3, x_4 .

The relevant F -statistic is

$$F = \frac{\text{SSReg}(x_1, x_2)/2}{\text{MSE}(x_1, x_2, x_3, x_4)} = \frac{(\text{SSReg}(x_1) + \text{SSReg}(x_2|x_1))/2}{\text{MSE}(x_1, x_2, x_3, x_4)} = \frac{(883880 + 8183)/2}{295.5} = 1495.471.$$

The p-value for the test is practically 0 since $F \gg 1$, hence significant evidence exists for the joint explanatory power x_1 and x_2 .

(j) Test the model with only variables x_1 and x_2 against the model with variables x_1, x_2 , and x_3 .

Note that $\text{RSS}(x_1, x_2, x_3) = \text{RSS}(x_1, x_2, x_3, x_4) - \text{SSReg}(x_4|x_1, x_2, x_3) = 14085$. Hence the relevant F -statistic is

$$F = \frac{\text{SSReg}(x_1, x_2)/2}{\text{RSS}(x_1, x_2, x_3)/(n - p + 1)} = \frac{883880 + 8183}{14085} = 1753.461.$$

As before, the p-value is practically 0 providing significant evidence for the joint significance of x_1 and x_2 .

(k) Should the test in part (g) be the same as the test in part (j)? Why or why not?

These are distinct tests. The first tests the significance of x_1, x_2 in the context of the full model with x_1, x_2, x_3, x_4 , while the second tests only relative to x_1, x_2 , and x_3 .

(l) For estimating the point on the regression surface at $(x_1, x_2, x_3, x_4) = (100, 25, 50, 4)$, the standard error of the estimate for the point on the surface is 2.62. Give the estimated point on the surface, a 95% confidence interval for the point on the surface, and a 95% prediction interval for a new point with these x -values.

The prediction is given by

$$\widehat{y|\mathbf{x}_0} = (1, 100, 25, 50, 4)\widehat{\boldsymbol{\beta}} \approx 124.0506$$

The 95% CI for the mean response given $\mathbf{x} = (100, 25, 50, 4)$ is given by $\widehat{y|\mathbf{x}_0} \pm t_{.975}(n - p)2.62$ which is (118.788, 129.313). The prediction interval is given by $\widehat{y|\mathbf{x}_0} \pm t_{.975}(n - p)\sqrt{MSE \cdot (1 + 2.62^2/MSE)}$ which is (-483.427, 731.528).

(m) Test the null hypothesis $H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$.

This is the model significance test given by the F -statistic

$$F = 600.8484,$$

Which yields a p-value that is practically 0. Hence, the model is significant.