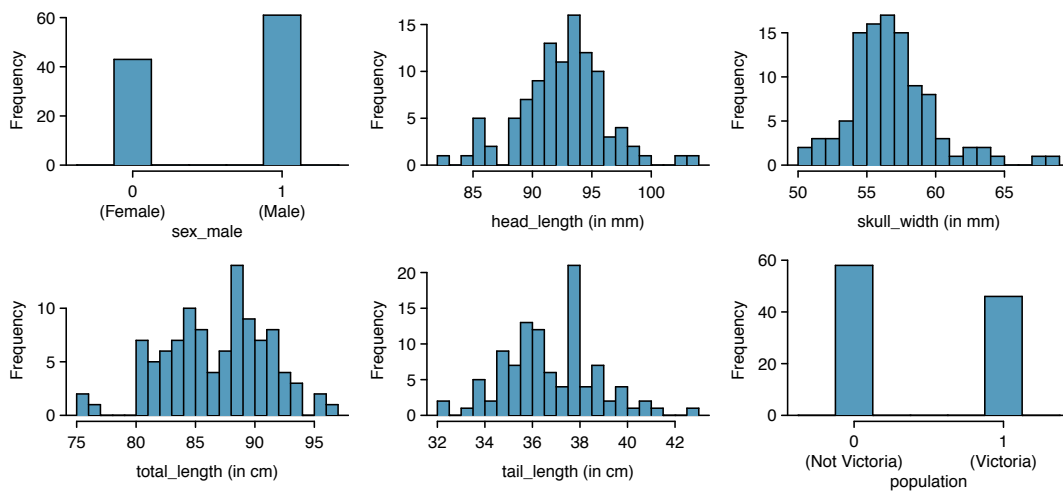### 8.5.4 Introduction to logistic regression

**8.15 Possum classification, Part I.** The common brushtail possum of the Australia region is a bit cuter than its distant cousin, the American opossum (see Figure 7.5 on page 334). We consider 104 brushtail possums from two regions in Australia, where the possums may be considered a random sample from the population. The first region is Victoria, which is in the eastern half of Australia and traverses the southern coast. The second region consists of New South Wales and Queensland, which make up eastern and northeastern Australia.

    We use logistic regression to differentiate between possums in these two regions. The outcome variable, called `population`, takes value 1 when a possum is from Victoria and 0 when it is from New South Wales or Queensland. We consider five predictors: `sex_male` (an indicator for a possum being male), `head_length`, `skull_width`, `total_length`, and `tail_length`. Each variable is summarized in a histogram. The full logistic regression model and a reduced model after variable selection are summarized in the table.



|  | | Full Model | | | | Reduced Model | | |
|---|---|---|---|---|---|---|---|---|
|  | Estimate | SE | Z | Pr(>\|Z\|) | Estimate | SE | Z | Pr(>\|Z\|) |
| (Intercept) | 39.2349 | 11.5368 | 3.40 | 0.0007 | 33.5095 | 9.9053 | 3.38 | 0.0007 |
| sex_male | -1.2376 | 0.6662 | -1.86 | 0.0632 | -1.4207 | 0.6457 | -2.20 | 0.0278 |
| head_length | -0.1601 | 0.1386 | -1.16 | 0.2480 |  |  |  |  |
| skull_width | -0.2012 | 0.1327 | -1.52 | 0.1294 | -0.2787 | 0.1226 | -2.27 | 0.0231 |
| total_length | 0.6488 | 0.1531 | 4.24 | 0.0000 | 0.5687 | 0.1322 | 4.30 | 0.0000 |
| tail_length | -1.8708 | 0.3741 | -5.00 | 0.0000 | -1.8057 | 0.3599 | -5.02 | 0.0000 |

(a) Examine each of the predictors. Are there any outliers that are likely to have a very large influence on the logistic regression model?

(b) The summary table for the full model indicates that at least one variable should be eliminated when using the p-value approach for variable selection: `head_length`. The second component of the table summarizes the reduced model following variable selection. Explain why the remaining estimates change between the two models.

**8.16 Challenger disaster, Part I.** On January 28, 1986, a routine launch was anticipated for the Challenger space shuttle. Seventy-three seconds into the flight, disaster happened: the shuttle broke apart, killing all seven crew members on board. An investigation into the cause of the disaster focused on a critical seal called an O-ring, and it is believed that damage to these O-rings during a shuttle launch may be related to the ambient temperature during the launch. The table below summarizes observational data on O-rings for 23 shuttle missions, where the mission order is based on the temperature at the time of the launch. *Temp* gives the temperature in Fahrenheit, *Damaged* represents the number of damaged O-rings, and *Undamaged* represents the number of O-rings that were not damaged.

| Shuttle Mission | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Temperature | 53 | 57 | 58 | 63 | 66 | 67 | 67 | 67 | 68 | 69 | 70 | 70 |
| Damaged | 5 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| Undamaged | 1 | 5 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 5 | 6 |

| Shuttle Mission | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Temperature | 70 | 70 | 72 | 73 | 75 | 75 | 76 | 76 | 78 | 79 | 81 |
| Damaged | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| Undamaged | 5 | 6 | 6 | 6 | 6 | 5 | 6 | 6 | 6 | 6 | 6 |

(a) Each column of the table above represents a different shuttle mission. Examine these data and describe what you observe with respect to the relationship between temperatures and damaged O-rings.

(b) Failures have been coded as 1 for a damaged O-ring and 0 for an undamaged O-ring, and a logistic regression model was fit to these data. A summary of this model is given below. Describe the key components of this summary table in words.

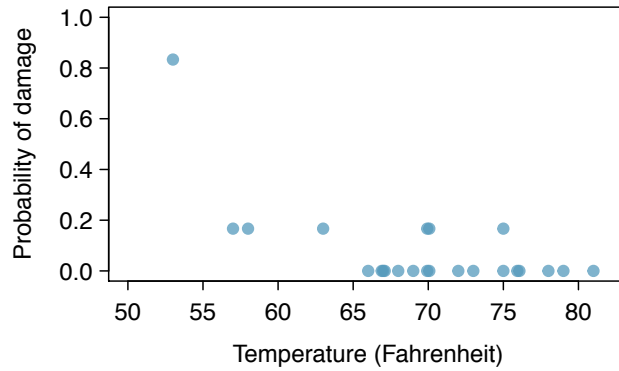| | Estimate | Std. Error | z value | Pr($>$|z|) |
|---|---|---|---|---|
| (Intercept) | 11.6630 | 3.2963 | 3.54 | 0.0004 |
| Temperature | -0.2162 | 0.0532 | -4.07 | 0.0000 |

(c) Write out the logistic model using the point estimates of the model parameters.

(d) Based on the model, do you think concerns regarding O-rings are justified? Explain.

**8.17 Possum classification, Part II.** A logistic regression model was proposed for classifying common brushtail possums into their two regions in Exercise 8.15. The outcome variable took value 1 if the possum was from Victoria and 0 otherwise.

| | Estimate | SE | Z | Pr($>$|Z|) |
|---|---|---|---|---|
| (Intercept) | 33.5095 | 9.9053 | 3.38 | 0.0007 |
| sex_male | -1.4207 | 0.6457 | -2.20 | 0.0278 |
| skull_width | -0.2787 | 0.1226 | -2.27 | 0.0231 |
| total_length | 0.5687 | 0.1322 | 4.30 | 0.0000 |
| tail_length | -1.8057 | 0.3599 | -5.02 | 0.0000 |

(a) Write out the form of the model. Also identify which of the variables are positively associated when controlling for other variables.

(b) Suppose we see a brushtail possum at a zoo in the US, and a sign says the possum had been captured in the wild in Australia, but it doesn't say which part of Australia. However, the sign does indicate that the possum is male, its skull is about 63 mm wide, its tail is 37 cm long, and its total length is 83 cm. What is the reduced model's computed probability that this possum is from Victoria? How confident are you in the model's accuracy of this probability calculation?

**8.18   Challenger disaster, Part II.** Exercise 8.16 introduced us to O-rings that were identified as a plausible explanation for the breakup of the Challenger space shuttle 73 seconds into takeoff in 1986. The investigation found that the ambient temperature at the time of the shuttle launch was closely related to the damage of O-rings, which are a critical component of the shuttle. See this earlier exercise if you would like to browse the original data.



(a) The data provided in the previous exercise are shown in the plot. The logistic model fit to these data may be written as

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = 11.6630 - 0.2162 \times Temperature$$

where $\hat{p}$ is the model-estimated probability that an O-ring will become damaged. Use the model to calculate the probability that an O-ring will become damaged at each of the following ambient temperatures: 51, 53, and 55 degrees Fahrenheit. The model-estimated probabilities for several additional ambient temperatures are provided below, where subscripts indicate the temperature:

| | | | |
|---|---|---|---|
| $\hat{p}_{57} = 0.341$ | $\hat{p}_{59} = 0.251$ | $\hat{p}_{61} = 0.179$ | $\hat{p}_{63} = 0.124$ |
| $\hat{p}_{65} = 0.084$ | $\hat{p}_{67} = 0.056$ | $\hat{p}_{69} = 0.037$ | $\hat{p}_{71} = 0.024$ |

(b) Add the model-estimated probabilities from part (a) on the plot, then connect these dots using a smooth curve to represent the model-estimated probabilities.

(c) Describe any concerns you may have regarding applying logistic regression in this application, and note any assumptions that are required to accept the model's validity.