

Chapter 9 - Multiple and Logistic Regression

David Simbandumwe

Baby weights, Part I. (9.1, p. 350) The Child Health and Development Studies investigate a range of topics. One study considered all pregnancies between 1960 and 1967 among women in the Kaiser Foundation Health Plan in the San Francisco East Bay area. Here, we study the relationship between smoking and weight of the baby. The variable *smoke* is coded 1 if the mother is a smoker, and 0 if not. The summary table below shows the results of a linear regression model for predicting the average birth weight of babies, measured in ounces, based on the smoking status of the mother.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	123.05	0.65	189.60	0.0000
smoke	-8.94	1.03	-8.65	0.0000

The variability within the smokers and non-smokers are about equal and the distributions are symmetric. With these conditions satisfied, it is reasonable to apply the model. (Note that we don't need to check linearity since the predictor has only two levels.)

- (a) Write the equation of the regression line.

$$\widehat{weight} = 123.05 - 8.94 \times smoke$$

- (b) Interpret the slope in this context, and calculate the predicted birth weight of babies born to smoker and non-smoker mothers.

```
(weight_no <- 123.05)
```

```
## [1] 123.05
```

```
(weight_yes <- 123.05 - 8.94*1)
```

```
## [1] 114.11
```

intercept: if a mother does not smoke the predicted birth weight is 123.05 **slope:** if a mother smokes the birth weight will be 8.94 ounces predicted weight if a mother does not smoke is 123.05 predicted weight if a mother smokes is 114.11

- (c) Is there a statistically significant relationship between the average birth weight and smoking?

yes the p-values are 0 indicating that there is significant evidence to reject the null hypothesis

Absenteeism, Part I. (9.4, p. 352) Researchers interested in the relationship between absenteeism from school and certain demographic characteristics of children collected data from 146 randomly sampled students in rural New South Wales, Australia, in a particular school year. Below are three observations from this data set.

	eth	sex	lrn	days
1	0	1	1	2
2	0	1	1	11
\vdots	\vdots	\vdots	\vdots	\vdots
146	1	0	0	37

The summary table below shows the results of a linear regression model for predicting the average number of days absent based on ethnic background (**eth**: 0 - aboriginal, 1 - not aboriginal), sex (**sex**: 0 - female, 1 - male), and learner status (**lrn**: 0 - average learner, 1 - slow learner).

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	18.93	2.57	7.37	0.0000
eth	-9.11	2.60	-3.51	0.0000
sex	3.10	2.64	1.18	0.2411
lrn	2.15	2.65	0.81	0.4177

- (a) Write the equation of the regression line.

$$\widehat{abs} = 18.93 - 9.11 \times eth + 3.10 \times male + 2.15 \times slow$$

- (b) Interpret each one of the slopes in this context.

intercept: model predicts 18.93 absence days for non aboriginal, female, average learners not aboriginal - results in a decrease of 9.11 absence days male - increase the absence days by 3.11 slow learners - increase the absence days by 2.15

- (c) Calculate the residual for the first observation in the data set: a student who is aboriginal, male, a slow learner, and missed 2 days of school.

```
eth <- 1
male <- 1
slow <- 1
abs <- 18.93 - 9.22*eth + 3.10*male + 2.15*slow

abs1 <- 2
(residual <- abs1 - abs)
```

```
## [1] -12.96
```

residual = -12.96

- (d) The variance of the residuals is 240.57, and the variance of the number of absent days for all students in the data set is 264.17. Calculate the R^2 and the adjusted R^2 . Note that there are 146 observations in the data set.

$$R^2 = 1 - \frac{Var(e_i)/(n - k - 1)}{Var(y_i)/(n - 1)}$$

```
var_e <-240.57
var_y <- 264.17
n <-146
k <- 3
```

```
(r_sq_adj = 1 - ( (var_e/var_y)))
```

```
## [1] 0.08933641
```

```
(r_sq_adj = 1 - ( (var_e/var_y) * ((n-1)/(n-k-1)) ))
```

```
## [1] 0.07009704
```

r-squared = 0.089 r-squared-adj = 0.070

Absenteeism, Part II. (9.8, p. 357) Exercise above considers a model that predicts the number of days absent using three predictors: ethnic background (**eth**), gender (**sex**), and learner status (**lrn**). The table below shows the adjusted R-squared for the model as well as adjusted R-squared values for all models we evaluate in the first step of the backwards elimination process.

	Model	Adjusted R^2
1	Full model	0.0701
2	No ethnicity	-0.0033
3	No sex	0.0676
4	No learner status	0.0723

Which, if any, variable should be removed from the model first?

remove ethnicity

Challenger disaster, Part I. (9.16, p. 380) 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.

lower temperatures correspond to more o rings failures

- (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.

intercept: 11.663 if the temperature is 0 it is expected that 11.663 o rings will fail slope: for each 0.216 decrease temperature an oring will fail

	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.

$$\widehat{o_rings} = 11.663 - 0.216 \times temp$$

intercept: at a temperature of 0 11.66 o rings will fail slope: for each 1 degree decrease in temperature 0.216 o rings will fail

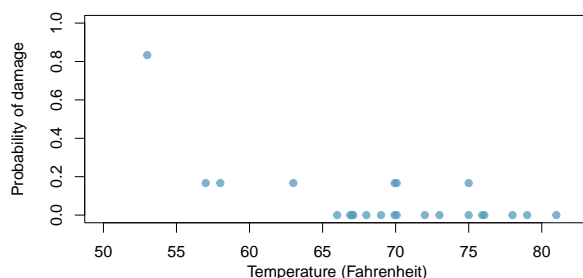
```
11.6630 - 0.2162 * 53
```

```
## [1] 0.2044
```

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

The p-value for the model is 0 so there is enough evidence to reject the null hypothesis justifying the concern for the o-rings however the reviewing the data set launch 1 seems like an outlier. No other launch has more than 1 failure. Since the sample size is so small it would warrant additional investigation before reaching a conclusion.

Challenger disaster, Part II. (9.18, p. 381) Exercise above 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

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

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:

$$\begin{array}{llll} \hat{p}_{57} = 0.341 & \hat{p}_{59} = 0.251 & \hat{p}_{61} = 0.179 & \hat{p}_{63} = 0.124 \\ \hat{p}_{65} = 0.084 & \hat{p}_{67} = 0.056 & \hat{p}_{69} = 0.037 & \hat{p}_{71} = 0.024 \end{array}$$

```
temp <- c(51,53,55)
logit_p <- 11.6630 - 0.2162 * temp
(p_hat <- exp(logit_p) / (1+exp(logit_p)))
```

```
## [1] 0.6540297 0.5509228 0.4432456
```

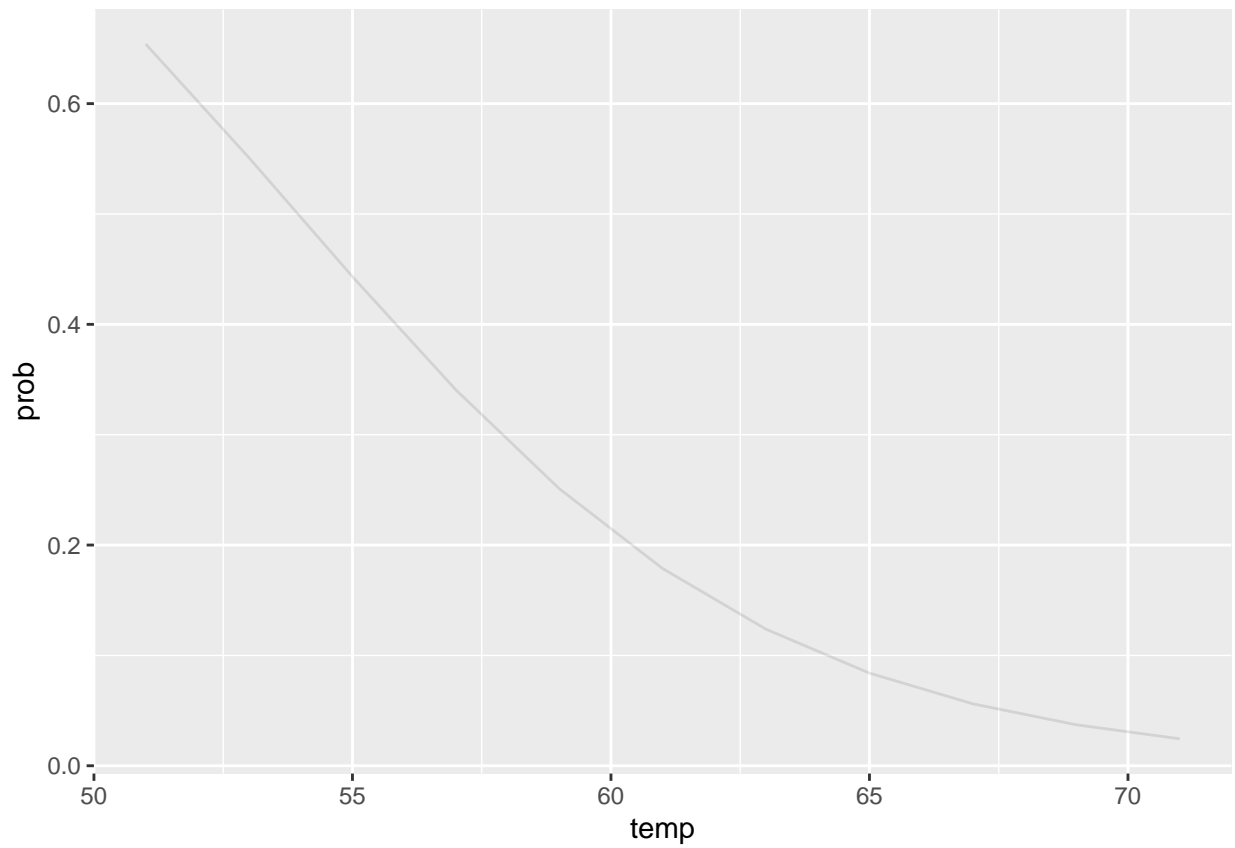
prob 51 = 0.654 prob 53 = 0.551 prob 55 = 0.443

- (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.

```
o_df <- tibble(temp=seq(51, 71, by = 2))
logit_p <- 11.6630 - 0.2162 * o_df$temp
(o_df$prob <- exp(logit_p) / (1+exp(logit_p)))
```

```
## [1] 0.65402974 0.55092283 0.44324565 0.34064976 0.25109139 0.17869707
## [7] 0.12372702 0.08393843 0.05612566 0.03715479 0.02443024
```

```
ggplot(data=o_df, aes(x=temp , y=prob)) +  
  geom_line(alpha=.1)
```



- (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.

The sample size is small with an obvious outlier that would bias the model.