

Chapter 9 - Multiple and Logistic Regression

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.

$$weight_{baby} = -8.94 \times smoke_{category} + 123.05$$

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

The y-intercept refers to the average baby weight at birth from this data.

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

It appears as though the smoking preferences of the mother has a neagative effect on the birth weight of the baby. Given the very low p-value and initial assumptions met, we have a significant result here.

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.

$$absent = -9.11 \times ethnicity + 3.10 \times sex + 2.15 \times learning + 18.93$$

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

The intercept is the average absent days for any student.

(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.

$$absent = 3.10(1) + 2.15(1) + 18.93 = 24.18$$

24 days missed is the prediction so the residual is -22 days.

(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.

$$R^2 = 1 - \frac{var_{res}}{var_{outcome}} = 1 - \frac{240.57}{264.17} = 0.08934$$

$$R^2_{adj} = 1 - \frac{var_{res}}{var_{outcome}} \times \frac{n-1}{n-k-1} = 1 - \frac{240.57}{264.17} \times \frac{146-1}{146-3-1} = 0.07001$$

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?

Learner status should be dropped as it increases the adjusted R squared value.

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.

There does seem to be more failures at low temperatures than at higher temps even though failures exist at higher temps as well.

- (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

We see a negative relationship between temperature and damaged O-rings. The p-value suggests a meaningful relationship.

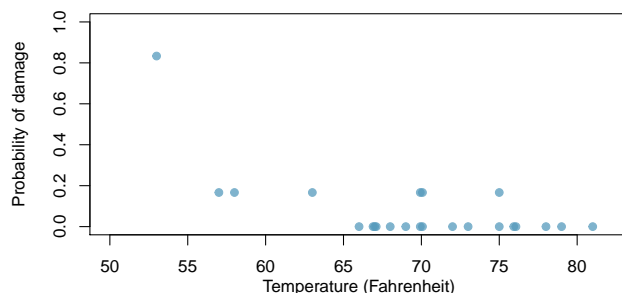
- (c) Write out the logistic model using the point estimates of the model parameters.

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

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

Given the error values and p-values for this data, it seems like a reasonable model in that we can be pretty confident that the slope of the line is not 0.

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:

```
temp_vec<- c(51,53,55)
for (i in temp_vec){
  linear<- exp(11.6630 - 0.2162*i)
  p<-exp(linear)/(1+exp(linear))
  print(paste("Probability of failure at temperature",i, "=",p))
}

## [1] "Probability of failure at temperature 51 = 0.868803621020442"
## [1] "Probability of failure at temperature 53 = 0.773256039910241"
## [1] "Probability of failure at temperature 55 = 0.689144812885923"
```

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

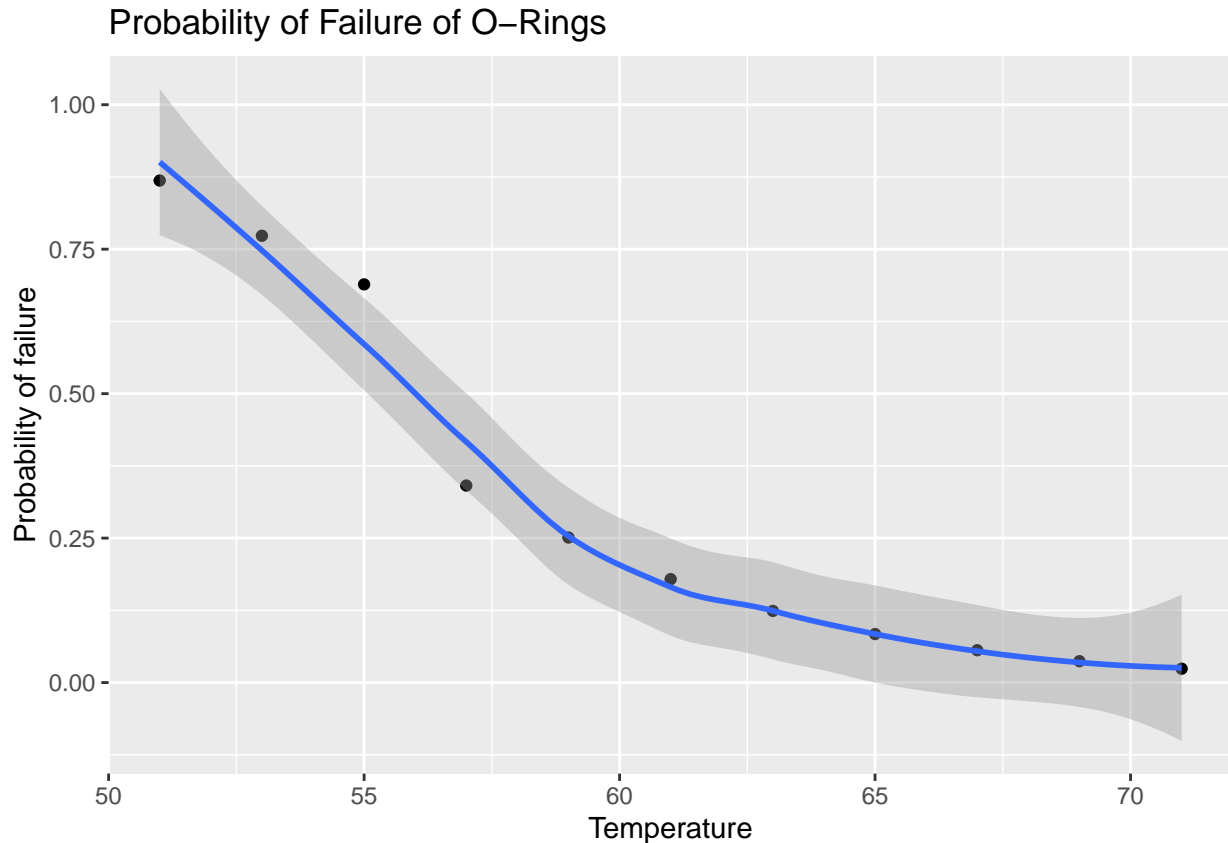
- (b) Add the model-estimated probabilities from part~(a) on the plot, then connect these dots using a smooth curve to represent the model-estimated probabilities.

```
library(ggplot2)

##
## Attaching package: 'ggplot2'
## The following object is masked from 'package:openintro':
##
##     diamonds
temp_vec<- c(51,53,55,57,59,61,63,65,67,69,71)
p_hat_vec<-c(.8688,.7732,.6891,.341,.251,.179,.124,.084,.056,.037,0.024)
ggplot() +
  geom_point(aes(x = temp_vec,y=p_hat_vec)) +
```

```
geom_smooth(aes(x = temp_vec,y=p_hat_vec)) +
labs(x="Temperature", y="Probability of failure",
title ="Probability of Failure of O-Rings")
```

`geom_smooth()` using method = 'loess' and formula 'y ~ x'



(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.

Fundamental assumptions necessary for Logistic Regression are independence of observations, which is met and Linear relationship of predictors. We have evidence to suggest that failure is linearly dependent on temp, but we should revisit that original linear model and verify that it is a solid model before going with the Logistic model. If that linear relationship holds, I think this would be a good prediction choice.