

## Illapani\_Homework\_HW8

Srini Illapani

November 30, 2015

### Multiple and logistic regression

#### 8.2 Baby weights, Part II.

Exercise 8.1 introduces a data set on birth weight of babies. Another variable we consider is parity, which is 0 if the child is the first born, and 1 otherwise. The summary table below shows the results of a linear regression model for predicting the average birth weight of babies, measured in ounces, from parity.

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	120.07	0.60	199.94	0.0000
parity	-1.93	1.19	-1.62	0.1052

(a) Write the equation of the regression line.

birth~weight = 120.07 - 1.93 \* parity

(b) Interpret the slope in this context, and calculate the predicted birth weight of first borns and others.

The estimated body weight of babies born that are the first born is 1.93 ounces lower than babies born that are not first born.

Predicted wt of First born = 120.07 - 1.93 \* 1 Predicted wt of Not first born = 120.07 - 1.93 \* 0

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

Ho: beta\_1 = 0 Ha: beta\_1 != 0 t = -1.62 p = 0.1052

The p value is large (> 0.05) and so the null hypothesis cannot be rejected. There is no strong evidence that there is an association between birth weight and parity.

#### 8.4 Absenteeism, Part I.

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
⋮	⋮	⋮	⋮	⋮
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).<sup>18</sup>

	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.

$$\text{absenteeism} = (18.93) + (-9.11 * \text{eth}) + (3.1 * \text{sex}) + (2.15 * \text{lrn})$$

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

The estimated days of absenteeism for an aboriginal is approx 9 days more than a non-aboriginal.

Males are likely to be approx 3 days more absent than a female.

Slow learners are likely to be absent approx 2 days more than an average learner.

(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.

```
absenteeism = 18.93 + (- 9.11 * 1) + (3.1 * 1) + (2.15 * 1)
absenteeism_case = 2
```

*# residual is the observed outcome minus the expected outcome.*

```
ei = 2 - 15.07
ei
```

```
## [1] -13.07
```

The residual is negative and this model over predicts absenteeism here.

(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<sup>2</sup> and the adjusted R<sup>2</sup>. Note that there are 146 observations in the data set.

```
# R2 = 1 - variability of residuals / variability of the outcome
R2 = 1 - (240.57/264.17)
R2
```

```
## [1] 0.08933641

# R2_adj = 1 - (variability of residuals / variability of the outcome) * (n-1 / n-k-1)
R2_adj = 1 - (240.57/264.17) * (146 - 1)/(146-3-1)
R2_adj
## [1] 0.07009704
```

## 8.8 Absenteeism, Part II.

Exercise 8.4 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. Which, if any, variable should be removed from the model first?

	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

"No learner status" should be removed because it has the greatest adjusted R2 value in the model.

## 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.

Based on the data from the table above, o-rings are prone to damage when the temperature is below 66 degrees. Above 66 degrees there is no visible pattern or correlation between damage and temperature.

- (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

The model shows an intercept of 11.66, the relevance of intercept is a bit diminished or irrelevant since we have 6 o-rings as part of the observation for each mission. The temperature, as it increases the chances of failure is decreased by 0.2162 for each degree increase in temperature.

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

$$\log(p / 1 - p) = 11.6630 - 0.2162 * \text{Temperature}$$

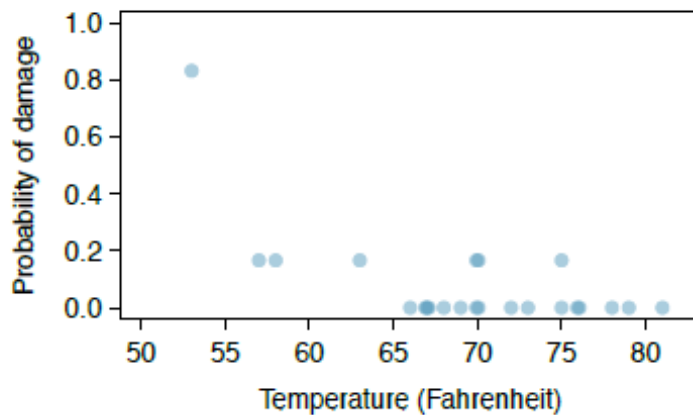
where p is the model-estimated probability that an O-ring will become damaged.

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

Yes, based on the data and the model I think that concerns regarding o-rings are justified. As increase in temperature at a certain point reduces the chances of damage to the o-ring.

## 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

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

where  $p$  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:

$$\hat{p}_{57} = 0.341$$

$$\hat{p}_{59} = 0.251$$

$$\hat{p}_{61} = 0.179$$

$$\hat{p}_{63} = 0.124$$

$$\hat{p}_{65} = 0.084$$

$$\hat{p}_{67} = 0.056$$

$$\hat{p}_{69} = 0.037$$

$$\hat{p}_{71} = 0.024$$

```
p_51 = exp (11.6630 - (.2162 * 51)) / (1 + exp (11.6630 - (.2162 * 51)))
p_51
## [1] 0.6540297

p_53 = exp (11.6630 - (.2162 * 53)) / (1 + exp (11.6630 - (.2162 * 53)))
p_53
## [1] 0.5509228

p_55 = exp (11.6630 - (.2162 * 55)) / (1 + exp (11.6630 - (.2162 * 55)))
p_55
## [1] 0.4432456
```

- (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.

```
# temperatures
x <- c(51, 53, 55, 57, 59, 61, 63, 65, 67, 69, 71)

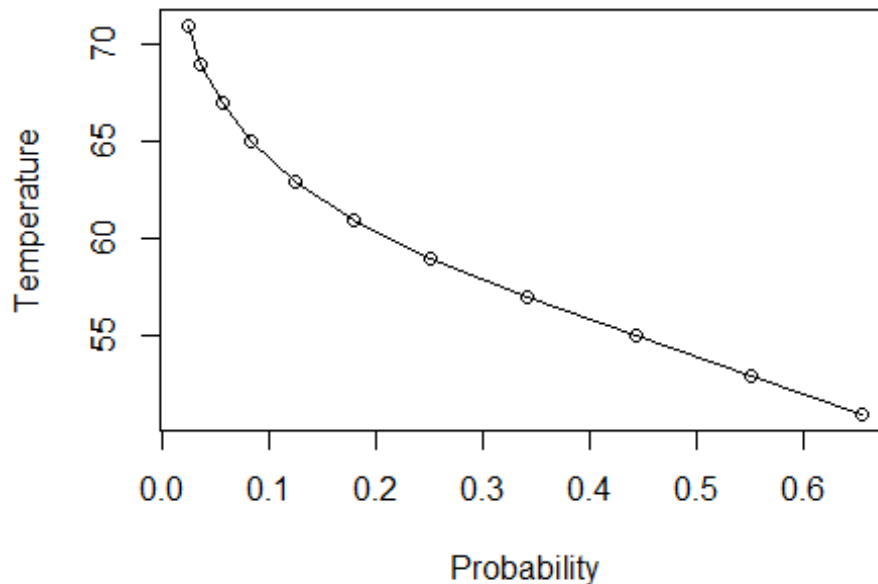
# probability
```

```

y <- c(0.654, 0.550, 0.443, 0.341, 0.251, 0.179, 0.124, 0.084, 0.056, 0.037,
0.024)

plot(x ~ y, xlab="Probability", ylab="Temperature")
lines(x ~ y)

```



- (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 for the number of o-rings is small but acceptable, the number of missions is a decent size though. One has to keep in mind that there could be many other variables that could cause the o-ring damage. We just found that temperature could be one such variable but to prove that this is the only source of cause, we need to rule out the impact of numerous other variables.