

# Module 9: Logistic Regression

Rebecca C. Steorts

# Agenda

- ▶ 1986 Challenger explosion
- ▶ What happened?
- ▶ Background
- ▶ Logistic regression

# The Challenger Case Study

On 28 January 1986, the Space Shuttle Challenger broke apart, 73 seconds into flight. All seven crew members died. The cause of the disaster was the failure of an o-ring on the right solid rocket booster.

# O-rings

- ▶ O-rings help seal the joints of different segments of the solid rocket boosters.
- ▶ We learned after this fatal mission that o-rings can fail at extremely low temperature.

# Motivations and goals

- ▶ In 1986, the Challenger space shuttle exploded as it took off.
- ▶ The question of interest was what happened and could it have been prevented?
- ▶ We will revisit not just the challenger data, but other missions to understand the relationship between o-ring failure and temperature.
- ▶ To understand this, we need to learn about logistic regression.

## Loading the Faraway Package

```
library(faraway)
data("orings")
orings[1,] <- c(53,1)
head(orings)
```

##	temp	damage
## 1	53	1
## 2	57	1
## 3	58	1
## 4	63	1
## 5	66	0
## 6	67	0

# Space Shuttle Missions

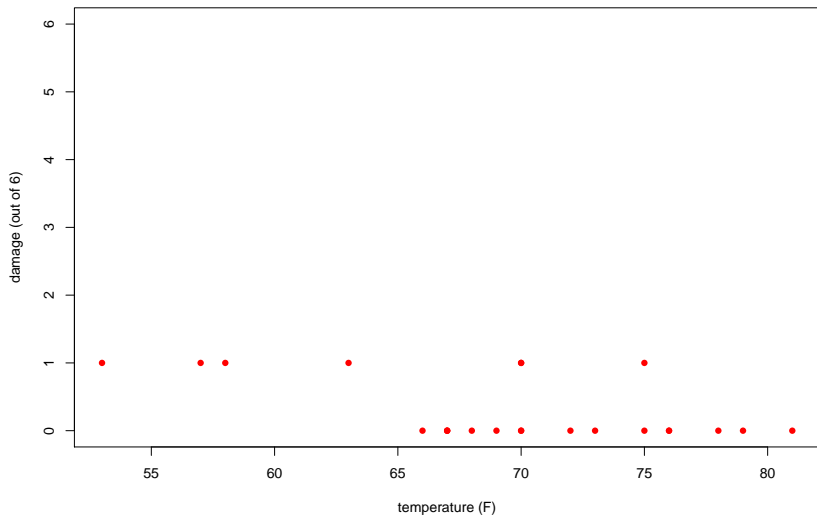
The 1986 crash of the space shuttle Challenger was linked to failure of o-ring seals in the rocket engines.

Data was collected on the 23 previous shuttle missions, where the following variables were collected:

- ▶ temperate for each mission
- ▶ damage to the number of o-rings (out of a total of six)

# Plot

```
plot(damage~temp, data=orings, xlab="temperature (F)",  
     ylab="damage (out of 6)",  
     pch=16, col="red", ylim=c(0,6))
```





# Plot

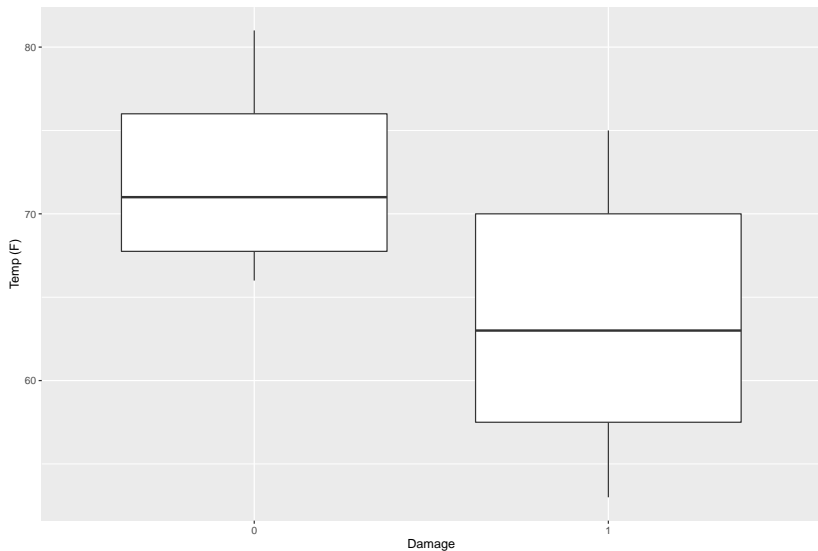
```
library(ggplot2)
geom_boxplot(outlier.colour="black", outlier.shape=14,
             outlier.size=2, notch=FALSE)

## geom_boxplot: outlier.colour = black, outlier.fill = NULL
## stat_boxplot: na.rm = FALSE, orientation = NA
## position_dodge2

damage <- as.factor(orings$damage)
temp <- orings$temp
head(damage)

## [1] 1 1 1 1 0 0
## Levels: 0 1
```

## Boxplot of temperature versus o-ring failure



## Response and covariate

- ▶ The response is the damage to the o-ring (in each shuttle launch).
- ▶ The covariate is the temperature (F) in each shuttle launch.

## Notation and Setup

- ▶ Let  $p_i$  be the probability that o-ring  $i$  fails.
- ▶ The corresponding **odds of failure** are

$$\frac{p_i}{1 - p_i}.$$

## Notation and Setup

- ▶ The probability of failure  $p_i$  is between  $[0, 1]$
- ▶ The odds of failure is any real number.

# Logistic Regression

The response

$$Y_i \mid p_i \sim \text{Bernoulli}(p_i) \quad (1)$$

for  $i = 1, \dots, n$ .

The logistic regression model writes that the logit of the probability  $p_i$  is a linear function of the predictor variable(s)  $x_i$ :

$$\text{logit}(p_i) := \log\left(\frac{p_i}{1 - p_i}\right) = \beta_0 + \beta_1 x_i. \quad (2)$$

# Interpretation of Co-efficients

- ▶ The regression co-efficients  $\beta_0, \beta_1$  are directly related to the log odds  $\log(\frac{p_i}{1-p_i})$  and not  $p_i$ .
- ▶ For example, the intercept  $\beta_0$  is the  $\log(\frac{p_i}{1-p_i})$  for observation  $i$  when the predictor takes a value of 0.
- ▶ The slope  $\beta_1$  refers to the change in the expected log odds of failure of an o-ring for a decrease in temperature.

## Intuition of Model

We assume our 23 data points are **conditionally independent**.

$$\Pr(\text{failure} = 1) = \frac{\exp\{\beta_0 + \beta_1 \times \text{temp}\}}{1 + \exp\{\beta_0 + \beta_1 \times \text{temp}\}}$$

$$\text{failure}_1, \dots, \text{failure}_{23} \mid \beta_0, \beta_1, \text{temp}_1, \dots, \text{temp}_{23} \quad (3)$$

$$\sim \prod_i \left( \frac{\exp\{\beta_0 + \beta_1 \times \text{temp}_i\}}{1 + \exp\{\beta_0 + \beta_1 \times \text{temp}_i\}} \right)^{\text{failure}_i} \quad (4)$$

$$\times \left( \frac{1}{1 + \exp\{\beta_0 + \beta_1 \times \text{temp}_i\}} \right)^{1 - \text{failure}_i} \quad (5)$$



## Exercise

Assume that  $\log(\frac{p_i}{1-p_i}) = \beta_0 + \beta_1 x_i$ .

Show that

$$p_i = \frac{e^{\beta_0 + \beta_1 x_i}}{e^{\beta_0 + \beta_1 x_i} + 1}.$$

This shows that logit function guarantees that the probability  $p_i$  lives in  $[0, 1]$ .

# Bayesian Logistic Regression

Recall that

$$Y_i \mid p_i \sim \text{Bernoulli}(p_i) \quad (6)$$

for  $i = 1, \dots, n$ .

$$\text{logit}(p_i) := \log\left(\frac{p_i}{1 - p_i}\right) = \beta_0 + \beta_1 x_i. \quad (7)$$

**How can we build minimal Bayesian prior knowledge?**

## Priors on $\beta_0$ and $\beta_1$

Conjugate priors do not exist on  $\beta_0$  and  $\beta_1$ .

We will consider the following weakly informative priors:

$$\beta_0 \sim \text{Normal}(0, 1000) \quad (8)$$

$$\beta_1 \sim \text{Normal}(0, 1000) \quad (9)$$

$$(10)$$

## Posterior sampling

Since we cannot find the posterior in closed form, we will resort to MCMC to approximate inference regarding  $\beta_0, \beta_1$ .

We can do this easily using the `logitMCMC` function in the `MCMCpack` R package.

This package implements a random walk Metropolis algorithm.

## Posterior sampling

```
library(MCMCpack)
```

```
## Loading required package: coda
```

```
## Loading required package: MASS
```

```
## ##
```

```
## ## Markov Chain Monte Carlo Package (MCMCpack)
```

```
## ## Copyright (C) 2003-2020 Andrew D. Martin, Kevin M. Quinn
```

```
## ##
```

```
## ## Support provided by the U.S. National Science Foundation
```

```
## ## (Grants SES-0350646 and SES-0350613)
```

```
## ##
```

```
failure <- orings$damage
```

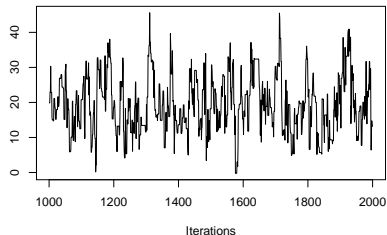
```
temperature <- orings$temp
```

```
output <- MCMClogit(failure~temperature,  
                    mcmc=1000, b0=0, B0=0.001)
```

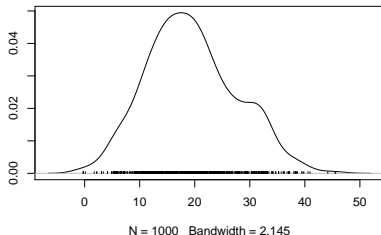
# Traceplots

`plot(output)`

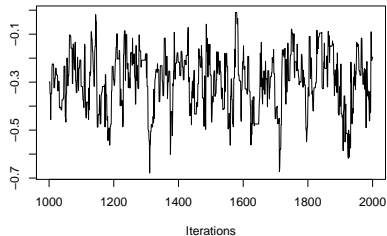
Trace of (Intercept)



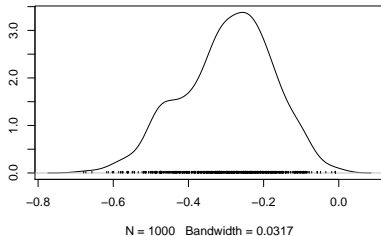
Density of (Intercept)



Trace of temperature



Density of temperature



# Summary

```
summary(output)
```

```
##
## Iterations = 1001:2000
## Thinning interval = 1
## Number of chains = 1
## Sample size per chain = 1000
##
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##              Mean      SD Naive SE Time-series SE
## (Intercept) 19.4239 8.1171 0.256684      0.88555
## temperature -0.2955 0.1191 0.003765      0.01309
##
## 2. Quantiles for each variable:
##
##              2.5%      25%      50%      75%      97.5%
## (Intercept)  5.3608 13.6196 18.7297 24.4156 36.08274
## temperature -0.5441 -0.3734 -0.2853 -0.2108 -0.09241
```

# Simulating Posterior Prediction

Given a certain temperature, we can simulate the results of future space shuttle launches using the posterior predictive distribution.

Suppose that on launch day, it's 80 degrees (F).

How would we simulate a predictive probability that a o-ring would fail?



# Simulating Posterior Prediction

```
library(boot)
```

```
##
```

```
## Attaching package: 'boot'
```

```
## The following objects are masked from 'package:faraway':
```

```
##
```

```
##      logit, melanoma
```

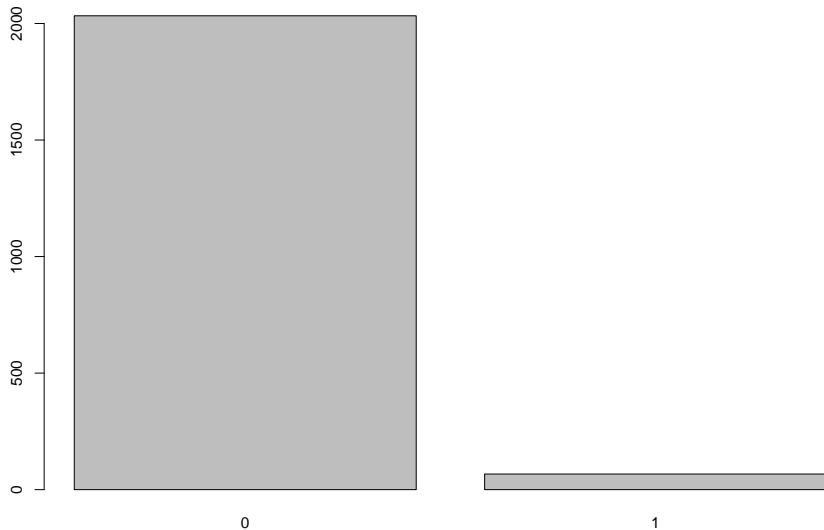
```
temp <- 80
```

```
fail.prob <- inv.logit(output[,1] + temp*output[,2])
```

```
y.pred <- rbinom(2100, size=1, prob=fail.prob)
```

# Simulating Posterior Prediction

```
barplot(table(y.pred))
```



## Your Turn

Suppose that it's very cold, 20 F.

How would we simulate a predictive probability that a o-ring would fail?

- ▶ What does your group think intuitively?
- ▶ Code up a simulation of a posterior prediction and what do you find?

# Summary

- ▶ Case Study on Challenger explosion
- ▶ Linear relationship between odds or log odds of failure and temperature seems reasonable implies use logistic regression
- ▶ What does logistic regression look like?
- ▶ We use the Metropolis algorithm via R
- ▶ What did we learn from the case study?
- ▶