

Binomial Regression I

Learning goals

Understand binomial regression

- know when you should use binormal regression.
- be able to write binomial regression model and its likelihood.
- be able to obtain estimators of parameters or function of parameters using R script.
- be able to quantify uncertainty of the estimators (e.g., computing CI).
- be able to test hypothesis.
- be able to do model selection.

Understand asymptotic properties of MLEs (maximum likelihood estimators)

- use asymptotic normality of MLEs to quantify uncertainty of the estimators (e.g., Wald CI).

Understand Wald test and likelihood ratio test (LRT)

- use them to test hypothesis in binomial regression

Understand (scaled) deviance

- use it to test model adequacy or perform LRT and model comparison.

Challenger example

Challenger disaster

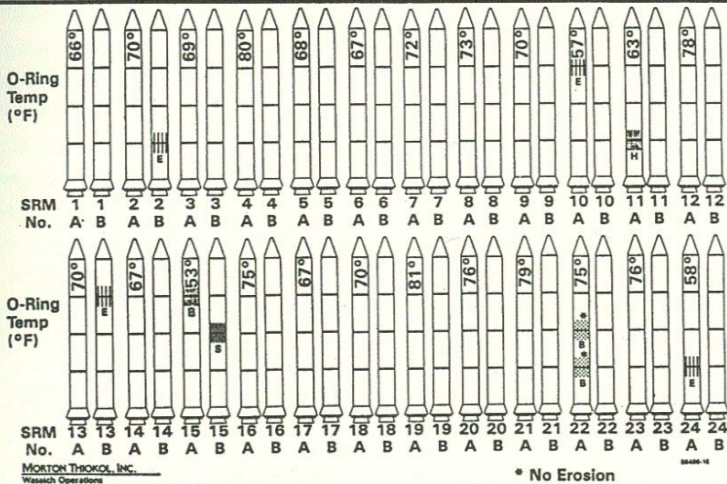
On the 28th of January 1986 the Space Shuttle Challenger broke apart after an O-ring seal failed at liftoff, leading to the deaths of its seven crew members.

Despite concerns about the O-rings failing due to the cold—the forecast temperature was $29 \pm 3^{\circ}\text{F}$ —no one was able to provide an analysis that convinced NASA (who were under pressure to launch following several previous delays) not to go ahead.

The way the data was presented didn't help matters.

Challenger disaster: data

History of O-Ring Damage in Field Joints (Cont)



MORTON THOKOL, INC.
Wassatch Operations

INFORMATION ON THIS PAGE WAS PREPARED TO SUPPORT AN ORAL PRESENTATION
AND CANNOT BE CONSIDERED COMPLETE WITHOUT THE ORAL DISCUSSION

Challenger disaster: data

Summarise data using number of damaged O-rings (out of 6) per launch.

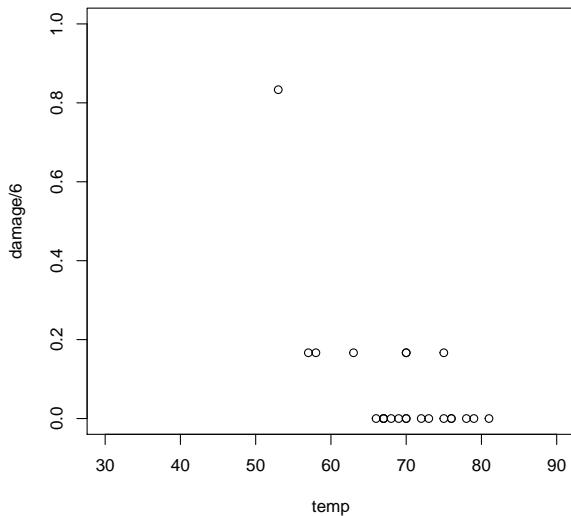
Data is in dataframe `orings`.

Response damage is number of damaged O-rings (out of 6).

Predictor `temp` is temperature ($^{\circ}F$)

```
> library(faraway)
> data(orings)
> str(orings)
data.frame: 23 obs. of  2 variables:
 $ temp  : num  53 57 58 63 66 67 67 67 68 69 ...
 $ damage: num  5 1 1 1 0 0 0 0 0 0 ...
> plot(damage/6 ~ temp, data = orings,
+       xlim = c(30, 90), ylim = c(0, 1))
```

Challenger disaster: visualize data



Challenger disaster: model

Make the assumption that Y_i , the number of damaged O-rings on the i -th launch, has distribution

$$Y_i \sim \text{bin}(6, p_i)$$

where p_i depends on the temperature t_i . We also assume that the Y_i are independent.

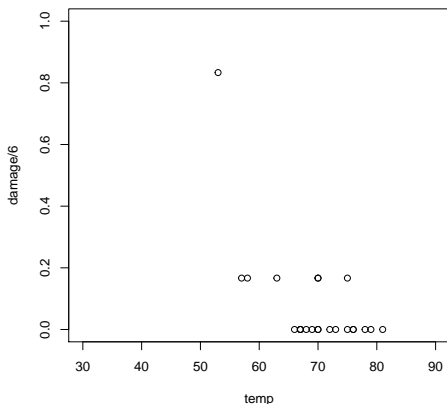
How to relate p_i with t_i .

- **Question:** just use a linear relationship? $p_i = \beta_0 + \beta_1 t_i$?

No. We need parameter bounds so $0 \leq p_i \leq 1$. We need alternate function to relate p_i with t_i .

Alternative function to relate p_i with t_i

For a single launch, best estimate of p_i is just $y_i/6$. From plot of $y_i/6$ against t_i it is reasonable to assume that $p_i = p(t_i)$ where p is a smooth function of the temperature, decreasing from 1 down to 0 as the temperature increases.



We choose logistic function: suppose that for some β_0 and β_1

$$p(t) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 t)}} = \frac{e^{\beta_0 + \beta_1 t}}{1 + e^{\beta_0 + \beta_1 t}}$$

- Restricts $0 \leq p(t) \leq 1$
- Monotonic increasing/decreasing function
- Can be linearized through the *logit* transformation

$$\log \frac{p(t)}{1 - p(t)} = \beta_0 + \beta_1 t$$

We choose logistic function: suppose that for some β_0 and β_1

$$p(t) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 t)}} = \frac{e^{\beta_0 + \beta_1 t}}{1 + e^{\beta_0 + \beta_1 t}}$$

- $p(t) = 1/2$ for $t = -\beta_0/\beta_1$, so $-\beta_0/\beta_1$ controls the location of the curve.
- $p'(-\beta_0/\beta_1) = \beta_1/4$, so β_1 controls the steepness of the curve.

Example: see R script and result in the section “logistic function with different values for beta0 and beta1” of Challenger.pdf.

Challenger disaster: model

Make the assumption that Y_i , the number of damaged O-rings on the i -th launch, has distribution

$$Y_i \sim \text{bin}(6, p_i)$$

where

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

We also assume that the Y_i are independent.

Challenger disaster: model fitting

Make the assumption that Y_i , the number of damaged O-rings on the i -th launch, has distribution

$$Y_i \sim \text{bin}(6, p_i)$$

where

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

We also assume that the Y_i are independent.

Maximum likelihood estimators (MLE) $\hat{\beta}_0, \hat{\beta}_1$: values for β_0, β_1 which maximize the log-likelihood.

Challenger disaster: model fitting

The log-likelihood is

$$\begin{aligned}l(\beta_0, \beta_1) &= \log \mathcal{L}(\beta_0, \beta_1) \\&= \log \mathbb{P}(Y = y \mid \beta_0, \beta_1) = \log \prod_i \mathbb{P}(Y_i = y_i \mid \beta_0, \beta_1) \\&= \sum_i \log \left(\binom{6}{y_i} p_i^{y_i} (1 - p_i)^{6 - y_i} \right) \\&= c + \sum_i (y_i \log p_i + (6 - y_i) \log(1 - p_i)) \\&= c + \sum_i \left(y_i \log \frac{p_i}{1 - p_i} + 6 \log(1 - p_i) \right)\end{aligned}$$

Put $\eta_i = \beta_0 + \beta_1 t_i$, then $\log p_i / (1 - p_i) = \eta_i$ and $\log(1 - p_i) = -\log(1 + e^{\eta_i})$.

Challenger disaster: model fitting

There is no closed form solution for MLE $\hat{\beta}_0, \hat{\beta}_1$ in this model.
Numerical search procedures are required to find MLE.

- the `glm` function uses the iterative weighted least squares (IWLS) algorithm - I will cover this later.
- For now, let's use the `optim` function.

Example:

- See R script and result in “maximum likelihood fitting” of Challenger.pdf
- $\hat{\beta}_0 = 11.667$ and $\hat{\beta}_1 = -0.216$

Challenger disaster: questions

- Forecast probability of an O-ring being damaged when the launch temperature is 29°F .
- How good is our forecast? Can we provide a confidence interval?
- Is temperature useful to predict the O-ring failing?

Binomial regression

Binomial regression model

We suppose that we observe $Y_i \sim \text{bin}(m_i, p_i)$, $i = 1, \dots, n$, independent.

The m_i are known and we suppose that for some **link function** g ,

$$g(p_i) = \eta_i = \mathbf{x}_i^T \boldsymbol{\beta} = \beta_0 + \beta_1 x_{i1} + \dots + \beta_q x_{iq}$$

where \mathbf{x}_i are known predictors and $\boldsymbol{\beta}$ are unknown parameters.

Binomial regression model: link function

Usual choices for g :

logit

$$\eta = \log \frac{p}{1-p}, \quad p = \frac{1}{1 + \exp(-\eta)} = \frac{\exp(\eta)}{1 + \exp(\eta)}$$

complementary log-log

$$\eta = \log(-\log(1-p)), \quad p = 1 - \exp(-e^\eta)$$

probit

$$\eta = \Phi^{-1}(p), \quad p = \Phi(\eta),$$

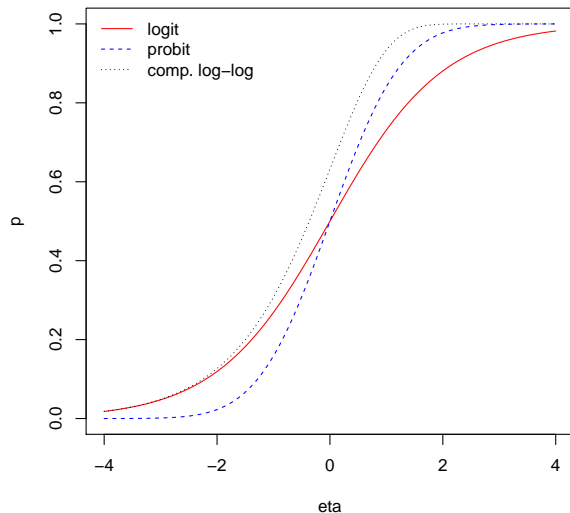
where Φ is the cumulative distribution function (cdf) of the standard normal distribution.

Binomial regression model: link function

```
> curve(1/(1+exp(-x)), -4, 4, ylim=c(0,1),  
+       xlab="eta", ylab="p", col="red",  
+       main="binomial link functions")  
> curve(pnorm(x), -4, 4, add=TRUE, col="blue", lty=2)  
> curve(1-exp(-exp(x)), -4, 4, add=TRUE, col="black", lty=3)  
> legend("topleft", c("logit", "probit", "comp. log-log"),  
+       col=c("red", "blue", "black"), lty=c(1,2,3), bty="n")
```

Binomial regression model: link function

binomial link functions



Binomial regression model: likelihood

Given observations y_i of $Y_i \sim \text{bin}(m_i, p_i = g^{-1}(\eta_i))$, where $\eta_i = \mathbf{x}_i^T \boldsymbol{\beta}$, the log-likelihood is

$$\begin{aligned} l(\boldsymbol{\beta}) &= \sum_{i=1}^n \log \mathbb{P}(Y_i = y_i) \\ &= \sum_{i=1}^n \log \left(\binom{m_i}{y_i} p_i^{y_i} (1 - p_i)^{m_i - y_i} \right) \\ &= c + \sum_{i=1}^n y_i \log(g^{-1}(\eta_i)) + (m_i - y_i) \log(1 - g^{-1}(\eta_i)) \end{aligned}$$

There is no closed form solution for MLE of $\boldsymbol{\beta}$. Numerical search procedures are required to find MLE.

Reminder: questions from Challenger disaster

- Forecast probability of an O-ring being damaged when the launch temperature is 29°F .
- How good is our forecast? Can we provide a confidence interval?
- Is temperature useful to predict the O-ring failing?

Forecast probability of an O-ring being damaged when the launch temperature is 29 °F.

$$\hat{p} = g^{-1}(\hat{\eta}) \text{ when } t = 29.$$

$$\hat{p} = \frac{\exp(\hat{\eta})}{1 + \exp(\hat{\eta})}, \text{ where } \hat{\eta} = \hat{\beta}_0 + \hat{\beta}_1 29$$

$\hat{p} = 0.995$: see R script and result in “prediction for temp of 29” of Challenger.pdf

Reminder: questions from Challenger disaster

- Forecast probability of an O-ring being damaged when the launch temperature is $29^{\circ}F$.
- How good is our forecast? Can we provide a confidence interval?
- Is temperature useful to predict the O-ring failing?

We need to know properties of MLE!!

Asymptotic properties MLE

Reminder: Maximum likelihood estimate (MLE)

Suppose that Y_i , $i = 1, \dots, n$, are independent, with densities/mass-functions $f_i(\cdot; \theta)$.

Given observations y_i of the Y_i , the log-likelihood is

$$l(\theta) = l(\theta; y) = \sum_i \log f_i(y_i; \theta).$$

MLE $\hat{\theta}$ is that value of θ which maximises $l(\theta)$.

Note that allowing f_i to depend on i means that we can include the case where the distribution of Y_i depends on some covariate x_i . That is, we can have $f_i(\cdot; \theta) = f(\cdot; x_i, \theta)$ for some common f .

Asymptotic properties of MLE

Under certain regularity conditions (to be introduced later), the MLE is

- asymptotically consistent,
- asymptotically normal,
- asymptotically efficient.

MLE: asymptotic consistency

Let θ^* denote a true value for θ . As $n \rightarrow \infty$, $\hat{\theta} \xrightarrow{P} \theta^*$.

That is for any $\epsilon > 0$

$$\mathbb{P}(|\hat{\theta} - \theta^*| > \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

MLE: asymptotic normality

Let θ^* denote a true value for θ .

$$\hat{\theta} \approx^d N(\theta^*, \mathcal{I}(\theta^*)^{-1}).$$

MLE: asymptotic normality

The **observed information** is

$$\mathcal{J}(\boldsymbol{\theta}) = -\frac{\partial^2 l(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} = -\mathbf{H}_{l(\boldsymbol{\theta})},$$

where the hessian matrix, $\mathbf{H}_{l(\boldsymbol{\theta})}$, is a square matrix of second-order partial derivatives of the log likelihood.

For binomial regression with one predictor (e.g., $\boldsymbol{\theta} = \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$):

$$\mathcal{J}(\boldsymbol{\theta}) = \mathcal{J}(\boldsymbol{\beta}) = \begin{pmatrix} -\frac{\partial^2 l(\boldsymbol{\beta})}{\partial \beta_0^2} & -\frac{\partial^2 l(\boldsymbol{\beta})}{\partial \beta_1 \partial \beta_0} \\ -\frac{\partial^2 l(\boldsymbol{\beta})}{\partial \beta_0 \partial \beta_1} & -\frac{\partial^2 l(\boldsymbol{\beta})}{\partial \beta_1^2} \end{pmatrix}$$

Clearly $\mathcal{J}(\boldsymbol{\theta}) = \mathcal{J}(\boldsymbol{\theta}; y)$ depends on y through $l(\boldsymbol{\theta}) = l(\boldsymbol{\theta}; y)$.

MLE: asymptotic normality

The Fisher information is

$$\mathcal{I}(\boldsymbol{\theta}) = \mathbb{E} \mathcal{J}(\boldsymbol{\theta}; \mathbf{Y}).$$

$\mathcal{I}(\boldsymbol{\theta})$ does not depend on \mathbf{y} .

Exercise: binomial regression with a logit link

We suppose that we observe $Y_i \sim \text{bin}(m_i, p_i)$, $i = 1, \dots, n$, independent, where $p_i = \frac{1}{1+\exp(-\eta_i)}$ and $\eta_i = \beta_0 + \beta_1 x_i$. The log-likelihood is

$$\begin{aligned} l(\beta_0, \beta_1) &= c + \sum_i \left[y_i \log \frac{p_i}{1 - p_i} + m_i \log(1 - p_i) \right] \\ &= c + \sum_i \left[y_i(\beta_0 + \beta_1 x_i) - m_i \log(1 + e^{\beta_0 + \beta_1 x_i}) \right]. \end{aligned}$$

Exercise: binomial regression with a logit link

Then,

$$\mathcal{J}(\beta) = \begin{pmatrix} \sum_i m_i p_i (1 - p_i) & \sum_i m_i x_i p_i (1 - p_i) \\ \sum_i m_i x_i p_i (1 - p_i) & \sum_i m_i x_i^2 p_i (1 - p_i) \end{pmatrix}.$$

See my derivation in “observed_information_binomial_regression.pdf”.

So, since there are no y_i terms left,

$$\mathcal{I}(\beta) = \begin{pmatrix} \sum_i m_i p_i (1 - p_i) & \sum_i m_i x_i p_i (1 - p_i) \\ \sum_i m_i x_i p_i (1 - p_i) & \sum_i m_i x_i^2 p_i (1 - p_i) \end{pmatrix}.$$

MLE: asymptotic normality

Let θ^* denote a true value for θ .

As $n \rightarrow \infty$,

$$\mathcal{I}(\theta^*)^{1/2}(\hat{\theta} - \theta^*) \xrightarrow{d} N(0, I).$$

That is,

$$\hat{\theta} \approx^d N(\theta^*, \mathcal{I}(\theta^*)^{-1}).$$

Any intuition for why variance is the inverse of fisher information which is related to hessian matrix?

MLE: asymptotic efficiency

- Asymptotic consistency: $\hat{\theta} \xrightarrow{P} \theta^*$
- Asymptotic normality: $\hat{\theta} \approx^d N(\theta^*, \mathcal{I}(\theta^*)^{-1})$

MLE: asymptotically unbiased estimator with smallest variance $\mathcal{I}(\theta^*)^{-1}$.

Wald CI for $\mathbf{t}^T \boldsymbol{\theta}$ (linear combination of parameters)

$$\mathbf{t} = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_m \end{pmatrix} \quad \boldsymbol{\theta} = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_m \end{pmatrix} \quad \boldsymbol{\theta}^* = \begin{pmatrix} \theta_1^* \\ \theta_2^* \\ \vdots \\ \theta_m^* \end{pmatrix} \quad \hat{\boldsymbol{\theta}} = \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\theta}_m \end{pmatrix}$$

$$\hat{\boldsymbol{\theta}} \approx N(\boldsymbol{\theta}^*, \mathcal{I}(\boldsymbol{\theta}^*)^{-1}).$$

In practice, we do not know $\boldsymbol{\theta}^*$, the true value of $\boldsymbol{\theta}$. Thus, we approximate $\mathcal{I}(\boldsymbol{\theta}^*)^{-1}$ using $\mathcal{I}(\hat{\boldsymbol{\theta}})^{-1}$, leading to

$$\hat{\boldsymbol{\theta}} \approx N(\boldsymbol{\theta}^*, \mathcal{I}(\hat{\boldsymbol{\theta}})^{-1}).$$

Then,

$$\mathbf{t}^T \hat{\boldsymbol{\theta}} \approx N(\mathbf{t}^T \boldsymbol{\theta}^*, \mathbf{t}^T \mathcal{I}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{t}).$$

Wald CI for $\mathbf{t}^T \boldsymbol{\theta}$ (linear combination of parameters)

$$\mathbf{t}^T \hat{\boldsymbol{\theta}} \approx N(\mathbf{t}^T \boldsymbol{\theta}^*, \mathbf{t}^T \mathcal{I}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{t}).$$

An approximate $100(1 - \alpha)\%$ confidence interval for $\mathbf{t}^T \boldsymbol{\theta}$:

$$\mathbf{t}^T \hat{\boldsymbol{\theta}} \pm z_{\alpha} \sqrt{\mathbf{t}^T \mathcal{I}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{t}}, \quad \text{where } \Phi(z_{\alpha}) = 1 - \alpha/2$$

In particular, taking \mathbf{t} = standard unit vectors, $\begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix},$

we can obtain CI for each of parameters $\theta_1, \theta_2, \dots, \theta_n$.

An approximate $100(1 - \alpha)\%$ confidence interval for θ_i :

$$\hat{\theta}_i \pm z_{\alpha} \sqrt{(\mathcal{I}(\hat{\boldsymbol{\theta}})^{-1})_{i,i}}$$

If \mathcal{I} is unavailable then we can approximate it using the observed information \mathcal{J} , i.e., $\mathcal{I}(\hat{\boldsymbol{\theta}})^{-1} \approx \mathcal{J}(\hat{\boldsymbol{\theta}}; \mathbf{y})^{-1}$.

Reminder: questions from Challenger disaster

- Forecast probability of an O-ring being damaged when the launch temperature is 29°F .
- How good is our forecast? Can we provide a confidence interval?
- Is temperature useful to predict the O-ring failing?

How good is our forecast? Can we provide a confidence interval?

CI for $p = g^{-1}(\eta) = \frac{\exp(\eta)}{1+\exp(\eta)}$, where $\eta = \beta_0 + \beta_1 29$.

- Step 1: Compute a CI for η , (η_l, η_r) .
- Step 2: CI for $p = g^{-1}(\eta)$ is $(g^{-1}(\eta_l), g^{-1}(\eta_r))$.

CI for p : (0.864307, 0.9998686). See R script and result in “Confidence Interval for p” of Challenger.pdf

log likelihood ratio CI

We have

$$2l(\hat{\theta}) - 2l(\theta^*) \approx \chi_k^2$$

where k is the dimension of θ^* .

This result can also, in principle, be used to construct a $100(1 - \alpha)\%$ confidence region for θ :

$$\{\theta : 2l(\hat{\theta}) - 2l(\theta) \leq \chi_k^2(1 - \alpha)\}$$

where $\chi_k^2(1 - \alpha)$ is the $100(1 - \alpha)\%$ point for a χ_k^2 distribution.

One lab problem will be to plot CI using the log likelihood ratio.

This approximation is generally better than the normal approximation for $\hat{\theta}$. That is, it holds for smaller sample sizes.

MLE: regularity conditions

For maximum likelihood theory to hold we require

- l smooth enough with respect to θ (third derivatives exist and continuous)
- Third order derivatives of l have bounded expectations
- Support of Y_i does not depend on θ
- The domain Θ of θ is finite dimensional and doesn't depend on Y_i
- θ^* is not on the boundary of Θ .

References

- McCullagh & Nelder (1989), Appendix A.
- F.W. Scholz, Maximum likelihood estimation. *Encyclopedia of Statistical Sciences* Vol. 7, p.4629ff. Wiley, 2006.