

Solution 8: Generalized linear regression

Yaqi Shi

07/23/2024

Part 1: Generalized linear model

Suppose that 2500 pregnant women are enrolled in a study and the outcome is the occurrence of preterm birth. Possible predictors of preterm birth include age of the woman, smoking, socioeconomic status, body mass index, bleeding during pregnancy, serum level of dde, and several dietary factors.

1. Formulate the problem of selecting the important predictors of preterm birth in a generalized linear model (GLM) framework.
2. Show the components of the GLM, including the link function and distribution (in exponential family form).
3. Describe (briefly) how estimation and inference could proceed via a frequentist approach.

Solution: $y_i = 1$ if woman i has preterm birth and $y_i = 0$ otherwise ($i = 1, \dots, n$) $y_i \sim \text{Bernoulli}(\pi_i)$
Probability density function:

$$\begin{aligned} f(y_i; \pi_i) &= \pi_i^{y_i} (1 - \pi_i)^{1-y_i} \\ &= \exp \{y_i \log \pi_i + (1 - y_i) \log (1 - \pi_i)\} \\ &= \exp \left\{ y_i \log \left(\frac{\pi_i}{1 - \pi_i} \right) + \log (1 - \pi_i) \right\} \\ &= \exp \left\{ \frac{y_i \theta_i - b(\theta_i)}{a(\phi)} + c(y_i, \phi) \right\} \end{aligned}$$

where

$$\theta_i = \log \left(\frac{\pi_i}{1 - \pi_i} \right), \quad b(\theta_i) = \log (1 + e^{\theta_i}), \quad a(\phi) = \phi = 1,$$

and $c(y_i, \phi) = 0$.

Link function: Any mapping from $\mathcal{R} \rightarrow [0, 1]$. A convenient choice is the canonical link,

$$\eta_i = \theta_i = \log \left(\frac{\pi_i}{1 - \pi_i} \right),$$

which is the logit. The probit and complementary log-log are alternatives. Frequentist Estimation: Maximum likelihood estimates can be obtained for a given model, say

$$\log \left(\frac{\pi_i}{1 - \pi_i} \right) = \mathbf{x}_i' \boldsymbol{\beta}$$

(where \mathbf{x}_i is a $p \times 1$ vector of predictors) by iterative weighted least squares

Frequentist Inference: One can select the important predictors to be included in the model by stepwise selection, using the AIC or BIC criterion.

Alternatively, one can just fit the model with all the predictors and then do inferences based on the MLEs and asymptotic standard errors. For example, for continuous predictors included as linear terms in the model, we can do a Wald test. Alternatively, we could do analysis of deviance (see notes for details) to test for significant differences in fit between the nested models with and without a particular predictor.

Part 2: GLMs in R (Logistic regression)

Consider the space shuttle data in the MASS library. Consider modeling the use of the autolander as the outcome (variable name use).

1. Fit a logistic regression model with autolander (variable auto) use (labeled as “auto” 1) versus not (0) as predicted by wind sign (variable wind).
2. Give the estimated odds ratio for autolander use comparing head winds, labeled as “head” in the variable headwind (numerator) to tail winds (denominator).
3. Give the estimated odds ratio for autolander use comparing head winds (numerator) to tail winds (denominator) adjusting for wind strength from the variable magn.
4. If you fit a logistic regression model to a binary variable, for example use of the autolander, then fit a logistic regression model for one minus the outcome (not using the autolander) what happens to the coefficients?

```
library(MASS)
?shuttle
data(shuttle)
head(shuttle)
```

```
##  stability error sign wind  magn vis  use
## 1      xstab  LX   pp head  Light  no auto
## 2      xstab  LX   pp head Medium no auto
## 3      xstab  LX   pp head Strong no auto
## 4      xstab  LX   pp tail Light  no auto
## 5      xstab  LX   pp tail Medium no auto
## 6      xstab  LX   pp tail Strong no auto
```

Solution

1-2.

```
str(shuttle)
```

```
## 'data.frame': 256 obs. of 7 variables:
## $ stability: Factor w/ 2 levels "stab","xstab": 2 2 2 2 2 2 2 2 2 2 ...
## $ error : Factor w/ 4 levels "LX","MM","SS",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ sign : Factor w/ 2 levels "nn","pp": 2 2 2 2 2 2 1 1 1 1 ...
## $ wind : Factor w/ 2 levels "head","tail": 1 1 1 2 2 2 1 1 1 2 ...
## $ magn : Factor w/ 4 levels "Light","Medium",...: 1 2 4 1 2 4 1 2 4 1 ...
## $ vis : Factor w/ 2 levels "no","yes": 1 1 1 1 1 1 1 1 1 1 ...
## $ use : Factor w/ 2 levels "auto","noauto": 1 1 1 1 1 1 1 1 1 1 ...
```

```
fit1 <- glm(use ~ wind, family = binomial(link = logit), data = shuttle)
summary(fit1)$coefficients
```

```
##              Estimate Std. Error   z value Pr(>|z|)
## (Intercept) -0.25131443  0.1781742 -1.4104987 0.1583925
## windtail    -0.03181183  0.2522429 -0.1261159 0.8996402
```

```
shuttle$use.binary = as.integer(shuttle$use=="auto")
fit2 <- glm(use.binary ~ wind, family = binomial(link = logit), data = shuttle)
summary(fit2)$coefficients
```

```
##              Estimate Std. Error   z value Pr(>|z|)
## (Intercept)  0.25131443  0.1781742  1.4104987 0.1583925
## windtail     0.03181183  0.2522429  0.1261159 0.8996402
```

```
exp(summary(fit2)$coefficients[1, 1])/exp(sum(summary(fit2)$coefficients[, 1]))
```

```
## [1] 0.9686888
```

3.

```
fit3 <- glm(use.binary ~ wind + magn, family = binomial(link = logit), data = shuttle)
summary(fit3)$coefficients
```

```
##              Estimate Std. Error   z value Pr(>|z|)
## (Intercept)  3.635093e-01  0.2840608  1.279688e+00 0.2006547
## windtail     3.200873e-02  0.2530225  1.265055e-01 0.8993318
## magnMedium   -1.074656e-15  0.3599481 -2.985586e-15 1.0000000
## magnOut      -3.795136e-01  0.3567709 -1.063746e+00 0.2874438
## magnStrong   -6.441258e-02  0.3589560 -1.794442e-01 0.8575889
```

```
exp(summary(fit3)$coefficients[1, 1])/exp(sum(summary(fit3)$coefficients[, 1:2, 1]))
```

```
## [1] 0.9684981
```

4.

```
fit4 = glm(1 - use.binary ~ wind, family = binomial(link=logit), data = shuttle)
summary(fit4)$coefficients
```

```
##              Estimate Std. Error   z value Pr(>|z|)
## (Intercept) -0.25131443  0.1781742 -1.4104987 0.1583925
## windtail    -0.03181183  0.2522429 -0.1261159 0.8996402
```

```
summary(fit2)$coefficients
```

```
##              Estimate Std. Error   z value Pr(>|z|)
## (Intercept)  0.25131443  0.1781742  1.4104987 0.1583925
## windtail     0.03181183  0.2522429  0.1261159 0.8996402
```

Part 3: GLMs in R (Poisson regression)

Consider the insect spray data `InsectSprays`. Fit a Poisson model using spray as a factor level.

1. Report the estimated relative rate comparing spray A (numerator) to spray B (denominator).
2. Consider a Poisson glm with an offset, t . So, for example, a model of the form `glm(count ~ x + offset(t), family = poisson)` where x is a factor variable comparing a treatment (1) to a control (0) and t is the natural log of a monitoring time. What is impact of the coefficient for x if we fit the model `glm(count ~ x + offset(t2), family = poisson)` where $t2 < -\log(10) + t$? In other words, what happens to the coefficients if we change the units of the offset variable. (Note, adding $\log(10)$ on the log scale is multiplying by 10 on the original scale.)

```
data("InsectSprays")
head(InsectSprays)
```

```
##    count spray
## 1     10    A
## 2      7    A
## 3     20    A
## 4     14    A
## 5     14    A
## 6     12    A
```

Solution:

```
df <- InsectSprays
fit <- glm(count ~ spray - 1, family='poisson', df)
exp(fit$coef[1])/exp(fit$coef[2])
```

```
##    sprayA
## 0.9456522
```

As the Poisson regression fits log of expectation by the linear regression, when all the regressor data is multiplied by an offset, then the log of expectation is shifted, so the intercept is changes and the slope coefficients remains intact. So, the answer is: The coefficient estimate is unchanged.