

Generalized Linear Models

Poisson Regression and Multinomial (Logistic) Regression

University of California, San Diego
Instructor: Armin Schwartzman

1 / 16

Aircraft Damage dataset

- Consider the Aircraft Damage dataset taken from *Applied Linear Regression* (4th Edition) by Weisberg. This is a dataset on the result of strike missions during the Vietnam War with A-4 or A-6 aircraft.
- The variables are:
 - y : is the number of locations where the aircraft was damaged
 - x_1 : indicates the type of plane (0 for A-4; 1 for A-6)
 - x_2 : is the bomb load in tons
 - x_3 : is the total months of aircrew experience

2 / 16

Dealing with count data: standard model

- The response represents **counts**.

Here the number of different values it takes is not large compared to the sample size. It could be considered numerical, but we have another option.
- The standard linear model
$$y|\mathbf{x} \sim \mathcal{N}(\mu(\mathbf{x}), \sigma^2), \quad \mu(\mathbf{x}) = \mathbb{E}(y|\mathbf{x}) = \boldsymbol{\beta}^\top \mathbf{x}$$
is not appropriate because
 1. y is an integer
 2. $\mu(\mathbf{x})$ will be negative for some \mathbf{x} 'sThis model is relevant and may hold approximately if y takes a large number of values. This is because the Poisson distribution looks normal if its mean is large.

3 / 16

Dealing with count data: Poisson model

- A more appropriate is the **Poisson regression** model:

$$y|\mathbf{x} \sim \text{Poisson}(\mu(\mathbf{x})), \quad \log(\mu(\mathbf{x})) = \boldsymbol{\beta}^\top \mathbf{x}$$

- The logarithm could be replaced by any other (**link**) function $g : (0, \infty) \rightarrow (-\infty, \infty)$ monotone.
- Note that, by design, the variance is a function of the mean:

$$\sigma(\mathbf{x})^2 = \text{Var}(y|\mathbf{x}) = \mu(\mathbf{x})$$

4 / 16

MLE for Poisson regression

A Poisson model is usually fitted by maximum likelihood. The log-likelihood is:

$$\begin{aligned} \ell(\mu_1, \dots, \mu_n) &= \sum_{i=1}^n [y_i \log(\mu_i) - \mu_i - \log(y_i!)] \\ &= \sum_{i=1}^n [y_i \mathbf{b}^\top \mathbf{x}_i - \exp(\mathbf{b}^\top \mathbf{x}_i) - \log(y_i!)] \end{aligned}$$

since $\mu_i = \mu(\mathbf{x}_i) = \exp(\mathbf{b}^\top \mathbf{x}_i)$.

We want to maximize this function of \mathbf{b} . No closed form expression exists in general, but the problem is convex (maximize a concave function).

5 / 16

Deviance

The **deviance** is defined as twice the log-likelihood ratio:

$$\text{DEV} = 2 \log \frac{\mathcal{L}(y_1, \dots, y_n)}{\mathcal{L}(\hat{\mu}_1, \dots, \hat{\mu}_n)} = 2 [\ell(y_1, \dots, y_n) - \ell(\hat{\mu}_1, \dots, \hat{\mu}_n)]$$

For linear regression:

$$\ell(\mu_1, \dots, \mu_n) = -\log(\sqrt{2\pi}\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu_i)^2 \quad \Rightarrow \quad \text{DEV} = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \hat{\mu}_i)^2$$

For Poisson regression:

$$\begin{aligned} \text{DEV} &= 2 \left(\sum_{i=1}^n [y_i \log(y_i) - y_i - \log(y_i!)] - \sum_{i=1}^n [y_i \log(\hat{\mu}_i) - \hat{\mu}_i - \log(y_i!)] \right) \\ &= 2 \sum_{i=1}^n [y_i \log(y_i/\hat{\mu}_i) - y_i + \hat{\mu}_i] \end{aligned}$$

The deviance plays the role of the residual sum of squares.

6 / 16

Education by Age dataset

- Consider the Education by Age data taken from <http://lib.stat.cmu.edu/DASL/Datafiles/Educationbyage.html>.

There are two categorical variables (factors): age group and highest degree.

- The main question is whether the two factors are **independent**.
- In general, suppose we have two paired categorical variables $\{(U_i, V_i) : i = 1, \dots, n\}$, with

$$U_i \in \{u_a : a = 1, \dots, A\}, \quad V_i \in \{v_b : b = 1, \dots, B\}$$

If the observations are independent, then the cell counts

$$y_{ab} = \#\{i : (U_i, V_i) = (u_a, v_b)\}$$

are **sufficient statistics**.

These counts are organized in a (two-way) **contingency table** with A rows and B columns, which is the analog of a two-way table for numerical data.

- Note that $y = (y_{ab} : a = 1, \dots, A; b = 1, \dots, B)$ is **multinomial** with sample size n and probabilities $p_{ab} = \mathbb{P}(U = u_a, V = v_b)$.

7 / 16

Pearson's χ^2 test

- Testing for independence means testing $H_0 : p_{ab} = p_{a.}p_{.b}$, where

$$p_{a.} = \mathbb{P}(U = u_a), \quad p_{.b} = \mathbb{P}(V = v_b)$$

- The most popular method is the **chi-square test of independence**. It rejects for large values of

$$\mathbb{X} = \sum_{a=1}^A \sum_{b=1}^B \frac{(y_{ab} - \hat{y}_{ab})^2}{\hat{y}_{ab}} \quad \text{where} \quad \hat{y}_{ab} = \frac{y_{a.} y_{.b}}{y_{..}}$$

- ▷ y_{ab} is the **observed** count for cell (a, b) , and

$$y_{a.} = \sum_b y_{ab}, \quad y_{.b} = \sum_a y_{ab}, \quad y_{..} = \sum_a \sum_b y_{ab} = n$$

are the sum for row a , the sum for column b , and the total sum (equal to the sample size).

- ▷ \hat{y}_{ab} is the **predicted** count for cell (a, b) under independence.
- ▷ Under the null, as $n \rightarrow \infty$, \mathbb{X} has the limiting distribution $\chi_{AB-A-B+1}^2 = \chi_{(A-1)(B-1)}^2$.

8 / 16

Poisson model for contingency tables

- Count data is, strictly speaking, multinomial data (assuming the observations were independently sampled from a homogeneous population).
- As an approximation, we model the count data as Poisson distributed:

$$y_{ab} \sim \text{Poisson}(\mu_{ab}), \quad \mu_{ab} = np_{ab}$$

This approximation is accurate if the sample is large enough.

- Then testing for independence of the two factors is formalized as testing

$$H_0 : \mu_{ab} = \frac{\mu_{a.} \mu_{.b}}{n} \quad \forall a, b$$

- From a Poisson regression point of view, testing for H_0 corresponds to testing for the restricted model with no interaction term.

9 / 16

Cleveland Clinic Foundation heart disease study

- Consider the cleveland dataset taken from <https://www.kaggle.com/datasets/chenngs/heart-disease-cleveland-uci>
8 variables are categorical, and 6 variables are numerical.
We first focus on predicting cond based on the other (14) characteristics.
- The response cond is categorical (binary), therefore this is a **classification** task.
- A standard linear model is not that relevant here.

10 / 16

Logistic regression

- Assume the response y is binary and “coded” as $y \in \{0, 1\}$.
- We want to fit the following model:

$$y|\mathbf{x} \sim \text{Bernoulli}(\mu(\mathbf{x})), \quad \mu(\mathbf{x}) = \mathbb{P}(y = 1|\mathbf{x}) = \mathbb{E}(y|\mathbf{x})$$

with

$$\mu(\mathbf{x}) = \frac{e^{\beta^\top \mathbf{x}}}{1 + e^{\beta^\top \mathbf{x}}}.$$

This relationship is defined through the logit link function, yielding the *log odds*

$$\text{logit}(\mu(\mathbf{x})) = \log\left(\frac{\mu(\mathbf{x})}{1 - \mu(\mathbf{x})}\right) = \beta^\top \mathbf{x}$$

- Note that, by design, the variance is a function of the mean:

$$\sigma(\mathbf{x})^2 = \text{Var}(y|\mathbf{x}) = \mu(\mathbf{x})(1 - \mu(\mathbf{x}))$$

11 / 16

Coefficient interpretation

- Let $\mathbf{e}_j = (0, \dots, 1, \dots, 0)^\top$ be a vector of zeros with a 1 in the j -th position. Suppose we increase variable x_j by 1 unit. Then the log odds ratio is:

$$\log\left(\frac{\mu(\mathbf{x} + \mathbf{e}_j)}{1 - \mu(\mathbf{x} + \mathbf{e}_j)}\right) - \log\left(\frac{\mu(\mathbf{x})}{1 - \mu(\mathbf{x})}\right) = \boldsymbol{\beta}^\top(\mathbf{x} + \mathbf{e}_j) - \boldsymbol{\beta}^\top\mathbf{x} = \beta_j$$

- The coefficient β_j is the log odds ratio when increasing x_j by unit while keeping all the other variables constant.

12 / 16

Classification boundary

- This model predicts (classifies) $y = 1$ at a new observation \mathbf{x} if $\mu(\mathbf{x}) > 1/2$, meaning that it predicts the class that is the most likely at \mathbf{x} .

As a consequence, the **boundary** b/w the two classes is the **hyperplane**:

$$\boldsymbol{\beta}^\top \mathbf{x} = 0$$

(If the first entry of \mathbf{x} is equal to 1 to represent the intercept, then this is an affine hyperplane.)

13 / 16

MLE and Deviance

- We again fit the model by maximum likelihood.

Let $g = \text{logit}$. The log-likelihood is:

$$\begin{aligned}\ell(\mu_1, \dots, \mu_n) &= \sum_{i=1}^n [y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)] \\ &= \sum_{i=1}^n [y_i \log(g^{-1}(\mathbf{b}^\top \mathbf{x}_i)) + (1 - y_i) \log(1 - g^{-1}(\mathbf{b}^\top \mathbf{x}_i))]\end{aligned}$$

Maximizing this concave function (of \mathbf{b}) is a convex optimization problem.

- The deviance has the following expression here:

$$\text{DEV} = -2 \sum_{i=1}^n [y_i \log(\hat{\mu}_i) + (1 - y_i) \log(1 - \hat{\mu}_i)]$$

where $\hat{\mu}_i = g^{-1}(\hat{\boldsymbol{\beta}}^\top \mathbf{x}_i)$.

14 / 16

Multinomial regression

- We turn to predicting `attplus` based on the individual characteristics. This is a categorical variable taking 5 distinct values.
- Assume the response y is categorical with K levels, e.g., $y \in \{1, \dots, K\}$.
- Let $\mu_k(\mathbf{x}) = \mathbb{P}(y = k|\mathbf{x})$. For $k = 1, \dots, K - 1$, we model these as

$$\log\left(\frac{\mu_k(\mathbf{x})}{\mu_K(\mathbf{x})}\right) = \beta_k^\top \mathbf{x}$$

same as

$$\mu_k(\mathbf{x}) = \frac{e^{\beta_k^\top \mathbf{x}}}{1 + \sum_{\ell=1}^{K-1} e^{\beta_\ell^\top \mathbf{x}}}, \quad k = 1, \dots, K - 1$$

$$\mu_K(\mathbf{x}) = \frac{1}{1 + \sum_{\ell=1}^{K-1} e^{\beta_\ell^\top \mathbf{x}}}$$

- In this model, the **boundary** b/w the classes k and ℓ is the **hyperplane**:

$$(\beta_k - \beta_\ell)^\top \mathbf{x} = 0$$

15 / 16

Overdispersion

- Assuming a one-parameter family as in the Poisson or logistic models implicitly ties the variance to the mean, in that $\sigma^2 = V(\mu)$. This may be found to be incongruent with the data.
- Introduce the **dispersion** parameter $\phi = \sigma^2/V(\mu)$. The one-parameter model is correct when $\phi = 1$. When $\phi > 1$, we have **overdispersion**.
- The function `glm` in R allows for an overdispersion parameter.

16 / 16