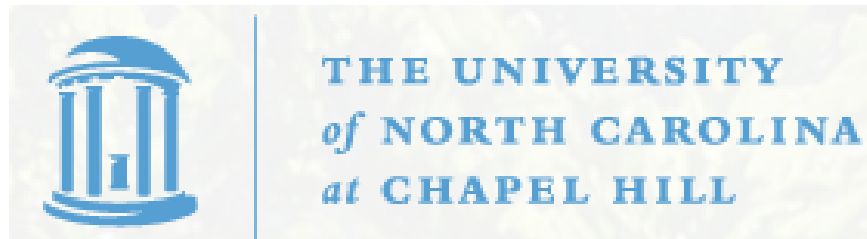


STOR 590:
ADVANCED LINEAR MODELS
Instructor: Richard L. Smith

Class Notes:
September 16, 2020



CLASS ANNOUNCEMENTS

- HW4: New deadline is Monday September 21 (late deadline Wednesday)
- From next week, we will revert to regular office hours, but look out for announced changes
- Take-home Midterm: Posted September 26, 6pm, to be returned by September 28, 6pm
- Spring 2020 midterm and final exams have been posted
- Final exam — still planning take-home exam, will update plans after the Midterm

CHAPTER 5: REGRESSION FOR COUNT DATA

1. Poisson Regression

Basics of Poisson model

- $\Pr\{Y = y\} = \frac{\mu^y e^{-\mu}}{y!}, \quad y = 0, 1, 2, \dots$
- Data: y_1, \dots, y_n Poisson with mean μ_1, \dots, μ_n
- Log link: $\log \mu_i = \eta_i = \sum_j x_{ij} \beta_j$
- Log likelihood $\ell(\mu_1, \dots, \mu_n) = \sum (y_i \log \mu_i - \mu_i - \log y_i!)$
- Unrestricted μ_i : maximized when $\mu_i = y_i$. Call this ℓ_1 .
- With log link and regressors:

$$\ell(\beta) = \sum_i \left\{ y_i \sum_j x_{ij} \beta_j - \exp \left(\sum_j x_{ij} \beta_j \right) - \log(y_i!) \right\},$$
$$\frac{\partial \ell(\beta)}{\partial \beta_k} = \sum_i \left\{ y_i x_{ik} - x_{ik} \exp \left(\sum_j x_{ij} \beta_j \right) \right\}.$$

Maximum Likelihood Estimators

- Write the *likelihood equations* as

$$\frac{\partial \ell(\hat{\beta})}{\partial \beta_k} = \sum_i \left\{ y_i x_{ik} - x_{ik} \exp \left(\sum_j x_{ij} \hat{\beta}_j \right) \right\} = 0.$$

- If we write $\exp \left(\sum_j x_{ij} \hat{\beta}_j \right) = \hat{\mu}_i$ we get

$$\sum_i (y_i - \hat{\mu}_i) x_{ik} = 0$$

which leads to the *normal equations*

$$X^T y = X^T \hat{\mu}.$$

- Note however we must still use numerical approximation to find $\hat{\mu}$.

Alternatives to Poisson Regression

- We can also try a standard linear regression, ignoring the fact that y is a count. The text starts out this way with the Species dataset
 - Simple linear regression did not give a good fit — variance increased with fitted value
 - Box-Cox transformation suggested $\lambda = 0.3$ but $\lambda = 0.5$ was almost as good on the plot
 - In fact taking $\lambda = 0.5$ is a standard trick for count data — the reason is given on the next slide
 - This improves on the untransformed linear regression but it still isn't perfect
 - Another problem with the square root transformation is difficulty of interpreting the resulting model — Poisson regression with log link is much easier to understand

Rationale for Square Root Transformation

- Suppose Y is Poisson with mean μ moderately large (say $\mu \geq 10$)
- The mean and variance of Y are both μ
- Write $Y = \mu(1 + \mu^{-1/2}\epsilon)$ where ϵ has mean 0 and variance 1
- Then $Y^{1/2} = \mu^{1/2}(1 + \mu^{-1/2}\epsilon)^{1/2} \approx \mu^{1/2} \left(1 + \frac{1}{2}\mu^{-1/2}\epsilon\right)$.
- $Y^{1/2}$ has mean approximately $\mu^{1/2}$ and variance approximately $\frac{1}{4}$ — *independent of μ*
- Therefore, a regression with $Y^{1/2}$ as the response should have approximately constant variance (standard deviation ≈ 0.5)
- However in the Species example, the residual standard error is 2.77, so this doesn't seem to work well either
- Further evidence of overdispersion

Deviance and Pearson X^2

- As for binary case, compare log likelihood for a saturated model (μ_i unrestricted) with the linear model being fitted,
- $\ell_1 = \sum_i (y_i \log y_i - y_i - \log y_i!)$
- $\ell_0 = \sum_i (y_i \log \hat{\mu}_i - \hat{\mu}_i - \log y_i!)$
- Deviance is

$$D = 2(\ell_1 - \ell_0) = 2 \sum_i \left(y_i \log \frac{y_i}{\hat{\mu}_i} - (y_i - \hat{\mu}_i) \right).$$

- We can also calculate the Pearson X^2 statistic

$$X^2 = \sum_i \frac{(y_i - \hat{\mu}_i)^2}{\hat{\mu}_i}.$$

Overdispersion

- Sometimes a more reasonable model may be $E(y_i) = \mu_i$, $\text{Var}(y_i) = \phi\mu_i$ where ϕ is a constant known as the *overdispersion* (usually but not necessarily $\phi > 1$)
- How to spot?
 - Plots of squared residuals against fitted values as in Fig. 5.3 (right — note that the plot is on a log scale here!)
 - Formal test of fit based on deviance or Pearson residuals (here leads to decisive rejection of the null hypothesis)
- Remedy — use `family=quasipoisson`
- For the species example we get a huge value $\phi = 31.7$
- There are still some observations with large Cook statistic but not nearly so bad as with the regular Poisson model

Use of Offset in R

“dicentric” example

```
rmod=glm(ca~offset(log(cells))+log(doserate)*dosef,  
family=poisson,dicentric)
```

Comment on HW problems

- Chapter 8, Question 4
- Chapter 8, Question 6