

Lecture 17

Inference with quasi-Poisson models

$$(n-p) \frac{\hat{\emptyset}}{\emptyset} \approx \chi^2_{n-p}$$

$$F_{d_1, d_2} = \frac{V_1 / d_1}{V_2 / d_2} \quad \begin{matrix} V_1 \sim \chi^2_{d_1} \\ V_2 \sim \chi^2_{d_2} \end{matrix}$$

$$\hat{\text{var}}(\hat{\beta})_{QP} = \hat{\emptyset} \text{var}(\hat{\beta})_{\text{poisson}} \quad T \sim t_{n-p} \Rightarrow T^2 \sim F_{1, n-p} \quad V_1 \perp V_2$$

$$\text{Test: } H_0: \beta_j = 0 \quad \text{vs.} \quad H_A: \beta_j \neq 0$$

$$\text{Test stat: } \frac{\hat{\beta}_j - 0}{SE(\hat{\beta}_j)} \sim N(0, 1) \quad \text{or} \quad \frac{(\hat{\beta}_j - 0)^2}{\text{var}(\hat{\beta}_j)} \sim \chi^2_1$$

Test for QP:

$$\frac{(\hat{\beta}_j - 0)^2}{\hat{\emptyset} \hat{\text{var}}(\hat{\beta}_j)_{\text{poisson}}} = \frac{(\hat{\beta}_j - 0)^2}{\underbrace{\emptyset \hat{\text{var}}(\hat{\beta}_j)_{\text{poisson}}}_{\chi^2_1}} \cdot \underbrace{\frac{\emptyset}{\hat{\emptyset}}}_{\frac{1}{(\chi^2_{n-p})/(n-p)}} \approx F_{1, n-p}$$

$$\Rightarrow \frac{\hat{\beta}_j - 0}{\sqrt{\hat{\emptyset}} SE(\hat{\beta}_j)_{\text{poisson}}} \approx t_{n-p}$$

$$\text{LRT: } \frac{D(Y, \hat{\mu}_{\text{red}}) - D(Y, \hat{\mu}_{\text{full}})}{\phi}$$

$$\approx \chi^2_q$$

$q = \# \text{ parameters tested}$

$$= df_{\text{red}} - df_{\text{full}}$$

$$\frac{(D(Y, \hat{\mu}_{\text{red}}) - D(Y, \hat{\mu}_{\text{full}})) / q}{\hat{\phi}_{\text{full}}}$$

$$= \frac{(D(Y, \hat{\mu}_{\text{red}}) - D(Y, \hat{\mu}_{\text{full}})) / q}{\underbrace{\phi}_{\chi^2_q / q}}$$

$$\approx F_{q, n-p_{\text{full}}}$$

(intuition: nested F test for linear regression)

$$F = \frac{(SSE_{\text{red}} - SSE_{\text{full}}) / q}{(SSE_{\text{full}}) / (n - p_{\text{full}})} \approx F_{q, n-p_{\text{full}}}$$

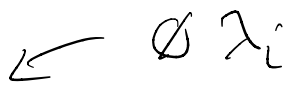
$$\frac{\frac{\phi}{\hat{\phi}}}{\frac{1}{\chi^2_{n-p_{\text{full}}} / (n-p_{\text{full}})}}$$

An alternative to quasi-Poisson

Poisson:

- Mean = λ_i
- Variance = λ_i

quasi-Poisson:

- Mean = λ_i
- Variance = $\varphi \lambda_i$ 
- Variance is a linear function of the mean

Question: What if we want variance to depend on the mean in a different way?

The negative binomial distribution

If $Y_i \sim \text{NB}(r, p)$, then Y_i takes values $y = 0, 1, 2, 3, \dots$ with probabilities

$$P(Y_i = y) = \frac{\Gamma(y + r)}{\Gamma(y + 1)\Gamma(r)} (1 - p)^r p^y$$

- $r > 0, \quad p \in [0, 1]$

- $\mathbb{E}[Y_i] = \frac{pr}{1 - p} = \mu$

- $\text{Var}(Y_i) = \frac{pr}{(1 - p)^2} = \mu + \frac{\mu^2}{r}$

For a fixed μ ,
 $\text{Var}(Y_i) \rightarrow \mu$ as
 $r \rightarrow \infty$

- Variance is a *quadratic* function of the mean

Key:

- Count variable ($Y_i = 0, 1, 2, \dots$)
- A lot more flexibility in the variance

Negative binomial regression

Poisson:

$$Y_i \sim \text{Poisson}(\mu_i)$$

$$\log(\mu_i) = \beta^T X_i$$

↑ canonical link function

$$Y_i \sim \text{NB}(r, p_i)$$

$$\log(\mu_i) = \beta^T X_i$$

↑
not

the canonical link function

- $\mu_i = \frac{p_i r}{1 - p_i}$
- Note that r is the same for all i
- Note that just like in Poisson regression, we model the average count
 - Interpretation of β s is the same as in Poisson regression

In R

If r is known, then the NB is an EDM

If r is unknown, then the NB is not an EDM

```
1 library(MASS)
2 m2 <- glm.nb(cigsPerDay ~ male + age + education +
3             diabetes + BMI, data = smokers)
```

Don't need to specify family)

```
...
(Intercept)  2.877771    0.123477    23.306    < 2e-16 ***
male          0.459148    0.027641    16.611    < 2e-16 ***
age          -0.007010    0.001731    -4.050    5.12e-05 ***
education2    0.024518    0.032534     0.754     0.451
education3    0.009252    0.040802     0.227     0.821
education4   -0.027732    0.044825    -0.619     0.536
diabetes     -0.010124    0.099126    -0.102     0.919
BMI           0.003693    0.003573     1.033     0.301
```

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(Dispersion parameter for Negative Binomial(3.2981) family taken to be 1)

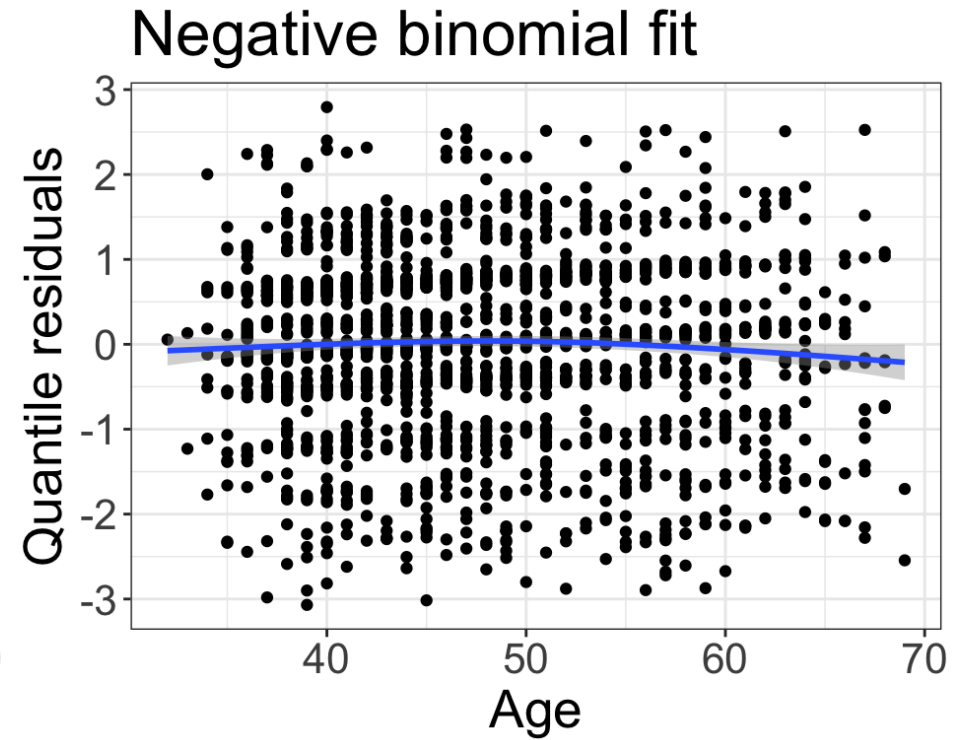
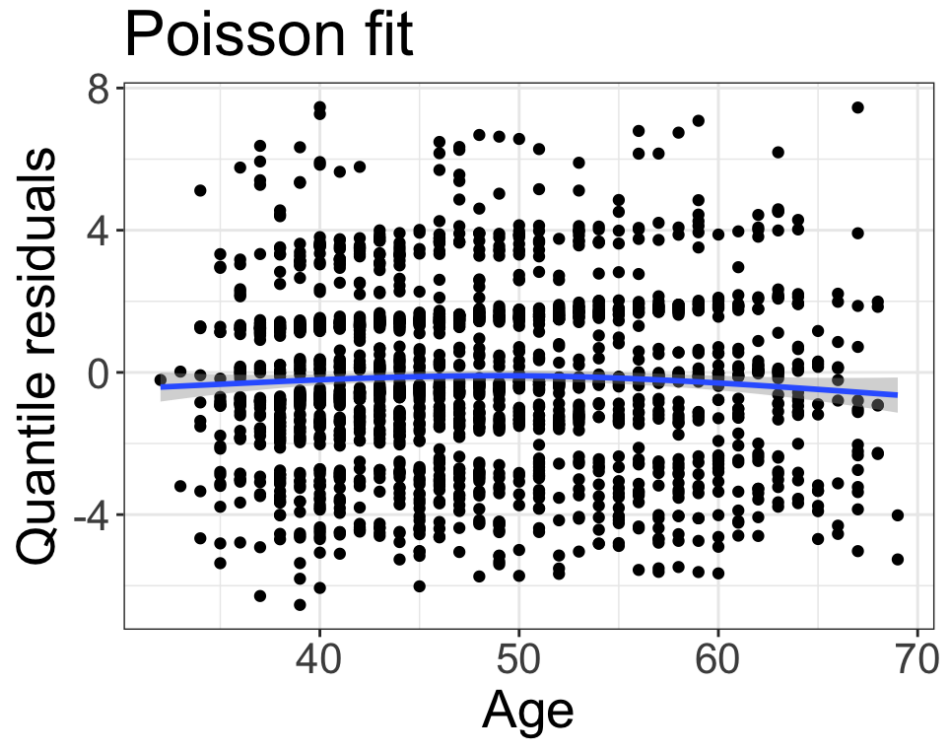
$$\hat{r} = 3.3$$

\hat{r}

$$r \neq 0$$

For NB, $\phi = 1$

Poisson vs. negative binomial fits



Inference with negative binomial models

```

...
                                z          Pr(>|z|)
(Intercept)    2.877771    0.123477    23.306    < 2e-16 ***
male           0.459148    0.027641    16.611    < 2e-16 ***
age            -0.007010    0.001731    -4.050    5.12e-05 ***
education2     0.024518    0.032534     0.754     0.451
education3     0.009252    0.040802     0.227     0.821
education4    -0.027732    0.044825    -0.619     0.536
diabetes       -0.010124    0.099126    -0.102     0.919
BMI            0.003693    0.003573     1.033     0.301
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

How would I test whether there is a relationship between age and the number of cigarettes smoked, after accounting for other variables?

$$H_0: \beta_2 = 0$$

↑
age

$$H_A: \beta_2 \neq 0$$

$$p\text{-value} \approx 5 \times 10^{-5}$$

Inference with negative binomial models

```
...
(Intercept)  2.877771    0.123477    23.306    < 2e-16 ***
male         0.459148    0.027641    16.611    < 2e-16 ***
age        -0.007010    0.001731    -4.050    5.12e-05 ***
education2   0.024518    0.032534     0.754     0.451
education3   0.009252    0.040802     0.227     0.821
education4  -0.027732    0.044825    -0.619     0.536
diabetes    -0.010124    0.099126    -0.102     0.919
BMI          0.003693    0.003573     1.033     0.301
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

...

How would I test whether there is a relationship between education and the number of cigarettes smoked, after accounting for other variables?

LRT : reduced model has sex, age, diabetes, BMI
(no education)

Likelihood ratio test

```
1 m2 <- glm.nb(cigsPerDay ~ male + age + education +  
2             diabetes + BMI, data = smokers)  
3 m3 <- glm.nb(cigsPerDay ~ male + age +  
4             diabetes + BMI, data = smokers)  
5 m2$twologlik - m3$twologlik
```

```
[1] 1.423055
```

```
1 pchisq(1.423, df=3, lower.tail=F)
```

```
[1] 0.7001524
```

$$\text{LRT: } 2(\log L_{\text{full}} - \log L_{\text{reduced}}) \approx \chi^2_q$$

Class activity

https://sta712-f23.github.io/class_activities/ca_lecture_17.html

