

Inference and overdispersion

Last time: modeling article publication

We are interested in analyzing the number of articles published by biochemistry PhD students. The data contains the following variables:

- + art: articles published in last three years of Ph.D.
- + fem: gender (recorded as male or female)
- + mar: marital status (recorded as married or single)
- + kid5: number of children under age six
- + phd: prestige of Ph.D. program
- + ment: articles published by their mentor in last three years

$$Articles_i \sim Poisson(\lambda_i)$$

$$\log(\lambda_i) = \beta_0 + \beta_1 Female_i + \beta_2 Married_i + \beta_3 Kids_i + \beta_4 Prestige_i + \beta_5 Mentor_i$$

Checking mean vs. variance

```
mean(articles$art)
```

```
## [1] 1.692896
```

```
var(articles$art)
```

```
## [1] 3.709742
```

Overdispersion

Overdispersion occurs when the response Y has higher variance than we would expect if Y followed a Poisson distribution.

What problems do you think it causes to assume the mean and variance are the same, when they are not?

Formal checks for overdispersion

First, we need a formal measure of dispersion (relation between mean and variance):

$$\phi = \frac{\text{Variance}}{\text{Mean}}$$

What should ϕ be if there is no overdispersion?

Hypothesis test for overdispersion

$$\phi = \frac{\text{Variance}}{\text{Mean}}$$

$H_0 : \phi = 1$ (no overdispersion)

$H_A : \phi > 1$ (overdispersion)

Now we need to estimate ϕ ...

Pearson residuals and estimated dispersion

The *Pearson residual* for observation i is defined as

$$e_{(P)i} = \frac{Y_i - \hat{\lambda}_i}{\sqrt{\hat{\lambda}_i}}$$

$$\hat{\phi} = \frac{\sum_{i=1}^n e_{(P)i}^2}{n - p}$$

+ p = number of parameters in model

Example: estimating dispersion parameter in R

```
# fit the model
m1 <- glm(art ~ ., data = articles,
          family = poisson)

# get Pearson residuals
pearson_resids <- resid(m1, "pearson")

# estimate dispersion parameter
phihat <- sum(pearson_resids^2)/(915 - 6)
phihat
```

```
## [1] 1.828984
```


Back to the hypothesis test

$$\phi = \frac{\text{Variance}}{\text{Mean}}$$

$H_0 : \phi = 1$ (no overdispersion)

$H_A : \phi > 1$ (overdispersion)

$$\hat{\phi} = 1.829$$

Now what?

Calculating a p-value

```
library(AER)  
dispersiontest(m1)
```

```
##  
##      Overdispersion test  
##  
## data:  m1  
## z = 5.7825, p-value = 3.681e-09  
## alternative hypothesis: true dispersion is greater than  
## sample estimates:  
## dispersion  
##      1.82454
```

So there is strong evidence for overdispersion in the data.

Handling overdispersion

Overdispersion is a problem because our standard errors (for confidence intervals and hypothesis tests) are too low.

If we think there is overdispersion, what should we do?

Adjusting the standard error

- + In our data, $\hat{\phi} = 1.829$
- + This means our variance is 1.829 times bigger than it should be
- + So our standard error is $\sqrt{1.829} = 1.352$ times bigger than it should be

New confidence interval for β_4 :

$$0.0128 \pm 1.96 \cdot \sqrt{1.829} \cdot 0.0264 = (-0.0572, 0.0828)$$

Adjusting the standard error in R

```
m2 <- glm(art ~ ., data = articles,  
          family = quasipoisson)
```

```
...  
##               Estimate Std. Error t value Pr(>|t|)  
## (Intercept)  0.304617   0.139273   2.187 0.028983 *  
## femWomen    -0.224594   0.073860  -3.041 0.002427 **  
## marMarried   0.155243   0.083003   1.870 0.061759 .  
## kid5        -0.184883   0.054268  -3.407 0.000686 ***  
## phd          0.012823   0.035700   0.359 0.719544  
## ment        0.025543   0.002713   9.415 < 2e-16 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## (Dispersion parameter for quasipoisson family taken to be 1.829006)  
...
```

- ✚ Allowing ϕ to be different from 1 means we are using a *quasi-likelihood* (in this case, a *quasi-Poisson*)

Adjusting the standard error in R

Poisson:

```
...  
##           Estimate Std. Error z value Pr(>|z|)  
## (Intercept)  0.304617   0.102981   2.958   0.0031 **  
## femWomen    -0.224594   0.054613  -4.112  3.92e-05 ***  
## marMarried   0.155243   0.061374   2.529   0.0114 *  
## kid5        -0.184883   0.040127  -4.607  4.08e-06 ***  
...
```

Quasi-Poisson:

```
...  
##           Estimate Std. Error t value Pr(>|t|)  
## (Intercept)  0.304617   0.139273   2.187  0.028983 *  
## femWomen    -0.224594   0.073860  -3.041  0.002427 **  
## marMarried   0.155243   0.083003   1.870  0.061759 .  
## kid5        -0.184883   0.054268  -3.407  0.000686 ***
```

Class activity

...

Coefficients:

##		Estimate	Std. Error	t value	Pr(> t)	
##	(Intercept)	-1.30445	0.34161	-3.818	0.000274	***
##	regionMW	0.09754	0.48893	0.199	0.842417	
##	regionNE	0.76268	0.42117	1.811	0.074167	.
##	regionSE	0.87237	0.42175	2.068	0.042044	*
##	regionSW	0.50708	0.50973	0.995	0.323027	
##	regionW	0.20934	0.51242	0.409	0.684055	

...

What confidence interval should I calculate to compare western and central schools?

Class activity

...

Coefficients:

##		Estimate	Std. Error	t value	Pr(> t)	
##	(Intercept)	-1.30445	0.34161	-3.818	0.000274	***
##	regionMW	0.09754	0.48893	0.199	0.842417	
##	regionNE	0.76268	0.42117	1.811	0.074167	.
##	regionSE	0.87237	0.42175	2.068	0.042044	*
##	regionSW	0.50708	0.50973	0.995	0.323027	
##	regionW	0.20934	0.51242	0.409	0.684055	

...

95% confidence interval for β_5 :

$$0.209 \pm 1.96 \cdot 0.512 = (-0.795, 1.213)$$

95% confidence interval for e^{β_5} :

$$(e^{-0.795}, e^{1.213}) = (0.452, 3.364)$$

Comparing Poisson and quasi-Poisson

Poisson:

- + Mean = λ_i
- + Variance = λ_i

quasi-Poisson:

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

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

Introducing the negative binomial

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

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

+ $\theta > 0, \quad p \in [0, 1]$

+ Mean = $\frac{p\theta}{1 - p} = \mu$

+ Variance = $\frac{p\theta}{(1 - p)^2} = \mu + \frac{\mu^2}{\theta}$

+ Variance is a *quadratic* function of the mean

Mean and variance for a negative binomial variable

If $Y_i \sim NB(\theta, p)$, then

+ Mean = $\frac{p\theta}{1-p} = \mu$

+ Variance = $\frac{p\theta}{(1-p)^2} = \mu + \frac{\mu^2}{\theta}$

How is θ related to overdispersion?

Negative binomial regression

$$Y_i \sim NB(\theta, p_i)$$

$$\log(\mu_i) = \beta_0 + \beta_1 X_i$$

- + $\mu_i = \frac{p_i \theta}{1 - p_i}$
- + Note that θ 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

Comparing Poisson, quasi-Poisson, negative binomial

Poisson:

- + Mean = λ_i
- + Variance = λ_i

quasi-Poisson:

- + Mean = λ_i
- + Variance = $\phi\lambda_i$

negative binomial:

- + Mean = μ_i
- + Variance = $\mu_i + \frac{\mu_i^2}{\theta}$

In R

```
m3 <- glm.nb(art ~ ., data = articles)
```

```
...  
##              Estimate Std. Error z value Pr(>|z|)  
## (Intercept)  0.256144   0.137348   1.865 0.062191 .  
## femWomen    -0.216418   0.072636  -2.979 0.002887 **  
## marMarried   0.150489   0.082097   1.833 0.066791 .  
## kid5        -0.176415   0.052813  -3.340 0.000837 ***  
## phd          0.015271   0.035873   0.426 0.670326  
## ment         0.029082   0.003214   9.048 < 2e-16 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1  
##  
## (Dispersion parameter for Negative Binomial(2.2644) fami  
...  
  
 $\hat{\theta} = 2.264$ 
```

In R

```
...  
##  
##           Estimate Std. Error z value Pr(>|z|)  
## (Intercept)  0.256144   0.137348   1.865 0.062191 .  
## femWomen    -0.216418   0.072636  -2.979 0.002887 **  
## marMarried   0.150489   0.082097   1.833 0.066791 .  
## kid5        -0.176415   0.052813  -3.340 0.000837 ***  
## phd          0.015271   0.035873   0.426 0.670326  
## ment         0.029082   0.003214   9.048 < 2e-16 ***  
...
```

How do I interpret the estimated coefficient -0.176?

quasi-Poisson vs. negative binomial

quasi-Poisson:

- + linear relationship between mean and variance
- + easy to interpret $\hat{\phi}$
- + same as Poisson regression when $\phi = 1$
- + simple adjustment to estimated standard errors
- + estimated coefficients same as in Poisson regression

negative binomial:

- + quadratic relationship between mean and variance
- + we get to use a likelihood, rather than a quasi-likelihood
- + Same as Poisson regression when θ is very large and p is very small