

EDMs and goodness of fit

Recap: EDMs and GLMs

Why is the canonical link function nice?

Data

A concerned parent asks us to investigate crime rates on college campuses. We have access to data on 81 different colleges and universities in the US, including the following variables:

- + type: college (C) or university (U)
- + nv: the number of violent crimes for that institution in the given year
- + enroll1000: the number of enrolled students, in thousands
- + region: region of the US C = Central, MW = Midwest, NE = Northeast, SE = Southeast, SW = Southwest, and W = West)

Model

$$Crimes_i \sim Poisson(\lambda_i)$$

$$\log(\lambda_i) = \beta_0 + \beta_1 MW_i + \beta_2 NE_i + \beta_3 SE_i + \beta_4 SW_i + \beta_5 W_i$$

Fitted model:

$$\log(\hat{\lambda}_i) = 1.34 + 0.48 MW_i + 0.44 NE_i + 0.77 SE_i + \\ 0.33 SW_i + 0.53 W_i$$

How would I interpret the intercept 1.34?

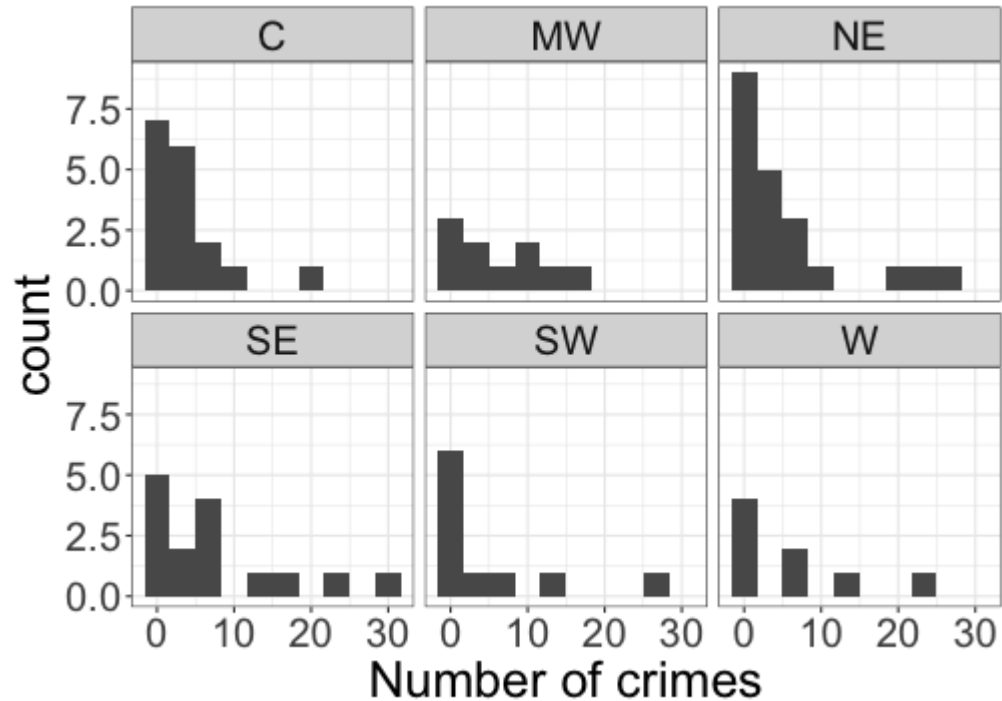
Model

$$Crimes_i \sim Poisson(\lambda_i)$$

$$\log(\lambda_i) = \beta_0 + \beta_1 MW_i + \beta_2 NE_i + \beta_3 SE_i + \beta_4 SW_i + \beta_5 W_i$$

What assumptions is this model making?

Exploratory data analysis



Exploratory data analysis

Mean and variance for number of crimes by region:

region	mean	variance
C	3.82	24.28
MW	6.20	37.07
NE	5.95	59.05
SE	8.27	84.35
SW	5.30	75.34
W	6.50	65.71

Goodness of fit

Goodness of fit test: If the model is a good fit for the data, then the residual deviance follows a χ^2 distribution with the same degrees of freedom as the residual deviance

...

Null deviance: 649.34 on 80 degrees of freedom

Residual deviance: 621.24 on 75 degrees of freedom

...

Residual deviance = 621.24, df = 75

How likely is a residual deviance of 621.24 if our model is correct?

Goodness of fit

Goodness of fit test: If the model is a good fit for the data, then the residual deviance follows a χ^2 distribution with the same degrees of freedom as the residual deviance

Residual deviance = 621.24, df = 75

```
pchisq(621.24, df=75, lower.tail=F)
```

```
## [1] 5.844298e-87
```

So our model might not be a very good fit to the data.

EDMs and deviance

Saddlepoint approximation

Goodness of fit

Goodness of fit test: If the model is a good fit for the data, then the residual deviance follows a χ^2 distribution with the same degrees of freedom as the residual deviance

Residual deviance = 621.24, df = 75

```
pchisq(621.24, df=75, lower.tail=F)
```

```
## [1] 5.844298e-87
```

So our model might not be a very good fit to the data.

Why might our model not be a good fit?

Offsets

We will account for school size by including an **offset** in the model:

$$\log(\lambda_i) = \beta_0 + \beta_1 MW_i + \beta_2 NE_i + \beta_3 SE_i + \beta_4 SW_i + \beta_5 W_i \\ + \log(Enrollment_i)$$

Motivation for offsets

We can rewrite our regression model with the offset:

$$\log(\lambda_i) = \beta_0 + \beta_1 MW_i + \beta_2 NE_i + \beta_3 SE_i + \beta_4 SW_i + \beta_5 W_i \\ + \log(Enrollment_i)$$

Fitting a model with an offset

```
m2 <- glm(nv ~ region, offset = log(enroll1000),
          data = crimes, family = poisson)
summary(m2)
```

```
...
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.30445    0.12403  -10.517  < 2e-16 ***
## regionMW     0.09754    0.17752   0.549   0.58270
## regionNE     0.76268    0.15292   4.987   6.12e-07 ***
## regionSE     0.87237    0.15313   5.697   1.22e-08 ***
## regionSW     0.50708    0.18507   2.740   0.00615 **
## regionW      0.20934    0.18605   1.125   0.26053
...
```

- ✚ The offset doesn't show up in the output (because we're not estimating a coefficient for it)

Fitting a model with an offset

$$\begin{aligned}\log(\hat{\lambda}_i) = & -1.30 + 0.10MW_i + 0.76NE_i + \\ & 0.87SE_i + 0.51SW_i + 0.21W_i \\ & + \log(Enrollment_i)\end{aligned}$$

How would I interpret the intercept -1.30?

When to use offsets

Offsets are useful in Poisson regression when our counts come from groups of very different sizes (e.g., different numbers of students on a college campus). The offset lets us interpret model coefficients in terms of rates instead of raw counts.

With your neighbor, brainstorm some other data scenarios where our response is a count variable, and an offset would be useful. What would our offset be?