

Integrating Poisson regression into the undergraduate curriculum

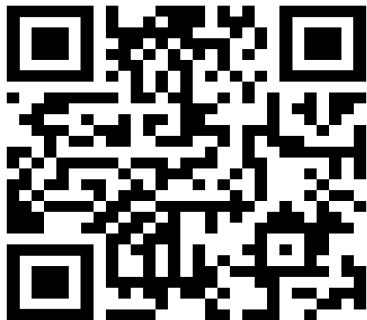
USCOTS25 Breakout Session B3H

Laura Boehm Vock and Paul Roback



A quick initial survey!

Please [click here](#) or use the following QR code:



Poisson regression at St. Olaf

- ▶ 2002-04: Not taught. Statistics concentration required Prob Theory and Math Stat plus 2 electives.
- ▶ 2004-18: Taught as part of Advanced Statistical Modeling (Stat 316). Concentration required Statistical Modeling (Stat 272) and 316 plus 2 electives.
- ▶ 2018-24: Still taught in Stat 316. Concentration renamed “Statistics and Data Science” and required 272 and Intro to Data Science plus 2 electives. Stat 316 now counts as an upper level elective.
- ▶ 2024-now: Still taught in Stat 316. Concentration became a major. Stat 316 counts as a “Level 3 Stats Depth” elective course.

Advanced Statistical Modeling at St. Olaf

- ▶ Covers generalized linear models (Poisson regr, binomial regr, negative binomial regr, zero-inflated models, hurdle models, etc.) and multilevel modeling
- ▶ Prerequisites: Intro Stats and Stat Modeling (nothing else – calculus, linear algebra, computing, ...)
- ▶ Applied focus using R
- ▶ Uses [Beyond Multiple Linear Regression: Applied Generalized Linear Models and Multilevel Models in R](#) by Roback and Legler. Second edition by Roback, Boehm Vock, and Legler expected by Fall 2026.

Initial survey results

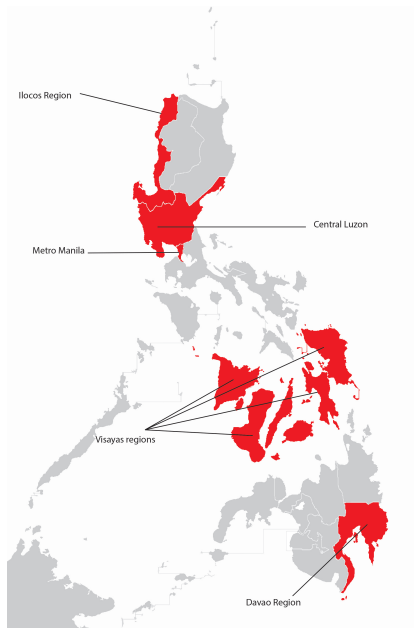
First case study: Philippine households

- ▶ International agencies often use household size to determine the magnitude of the household needs
- ▶ Want to discern factors associated with larger households
- ▶ We will model both the total household size and number of children under 5
- ▶ Data is subset from 2015 Philippine Statistics Authority's Family Income and Expenditure Survey (FIES)
- ▶ Primary response is a count, which can make linear regression problematic

Philippine household data

Key variables:

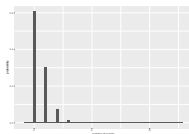
- ▶ `location` = region (Central Luzon, Davao, Ilocos, Metro Manila, or Visayas)
- ▶ `age` = the age of the head of household
- ▶ `total` = the number of people in the household other than the head
- ▶ `numLT5` = the number in the household under 5 years of age
- ▶ `roof` = the type of roof (stronger material can be used as a proxy for greater wealth)



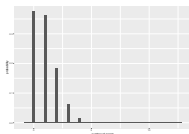
What is the Poisson *distribution*?

$$P(Y_i = y_i) = \frac{e^{-\lambda} \lambda^{y_i}}{y_i!} \quad \text{for } y_i = 0, 1, \dots, \infty,$$

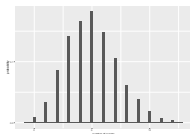
Note that both $E(Y_i) = \lambda_i$ and $Var(Y_i) = \lambda_i$.



(a) $\lambda = 0.5$



(b) $\lambda = 1$



(c) $\lambda = 5$

Figure 1: Poisson distributions with $\lambda = 0.5$, 1, and 5.

What is a Poisson *regression model*?

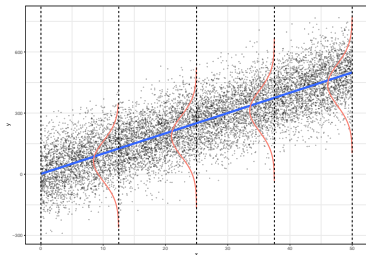
$$\log(\lambda_i) = \beta_0 + \beta_1 x_i$$

where the observed values $Y_i \sim \text{Poisson}$ with $\lambda = \lambda_i$ for a given x_i .

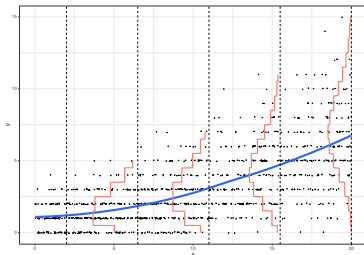
Poisson model conditions:

1. **Poisson Response** The response variable is a count per unit of time or space, described by a Poisson distribution.
2. **Independence** The observations must be independent of one another.
3. **Mean=Variance** By definition, the mean of a Poisson random variable must be equal to its variance.
4. **Linearity** The log of the mean rate, $\log(\lambda)$, must be a linear function of x .

Poisson regression conditions: A graphical look



(a) Linear Regression



(b) Poisson Regression

Figure 2: Comparison of regression models.

Pause to Ponder

With your neighbor(s), compare the Poisson regression conditions to the usual LINE conditions in linear regression. List similarities and differences. What implications might the differences have for modeling and checking conditions?

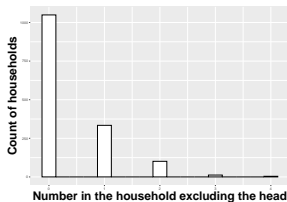
Differences with linear regression (LLSR)

1. For each level of X , the responses follow a Poisson distribution (Condition 1). For Poisson regression, small values of λ are associated with a distribution that is noticeably skewed with lots of small values and only a few larger ones. As λ increases the distribution of the responses begins to look more and more like a normal distribution.
2. In the LLSR model, the variation in Y at each level of X , σ^2 , is the same. For Poisson regression the responses at each level of X become more variable with increasing means, where variance=mean (Condition 3).
3. In the case of LLSR, the mean responses for each level of X , $\mu_{Y|X}$, fall on a line. In the case of the Poisson model, the mean values of Y at each level of X , $\lambda_{Y|X}$, fall on a curve, not a line, although the logs of the means should follow a line (Condition 4).

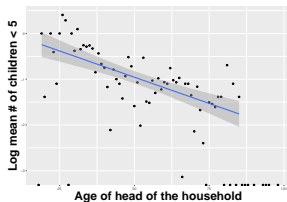
Side by side comparison

Linear Regression	Poisson Regression
$Y \sim N(\mu, \sigma)$	$Y \sim Pois(\lambda)$
$\mu = \beta_0 + X\beta_1$	$\log(\lambda) = \beta_0 + X\beta_1$
L inear relationship of X and Y	Log linear relationship of X and E(Y)
I ndependent observations	Independent observations
N ormally distributed residuals	Poisson distributed Y
E qual variance of all residuals	Variance increases with E(Y) ($\text{Var}(Y) = E(Y)$)

Exploratory data analysis 1: Number under 5 in household



(a) Distribution of number of children under 5 in households across all 5 Philippine regions.



(b) The log of the mean number of children under 5 by age of the head of household, with loess smoother.

Figure 3: EDA: selected plots

Initial model

$$\log(\hat{\lambda}) = 0.525 - 0.029\text{age}$$

##	Estimate	Std. Error	z value	Pr(> z)
## (Intercept)	0.52550316	0.149021042	3.526369	4.212997e-04
## age	-0.02938696	0.003052136	-9.628326	6.071065e-22

Pause to Ponder: How would you want your students to interpret coefficient estimates above?

Interpreting model coefficients

If your students have interpreted coefficients with a log-transformed response in linear regression, or log odds in logistic regression, this is similar.

Consider how the estimated mean number in the house, λ , changes as the age of the household head increases by an additional year.

$$\begin{aligned} \log(\lambda_X) &= \beta_0 + \beta_1 X \\ \log(\lambda_{X+1}) &= \beta_0 + \beta_1(X + 1) \\ \log(\lambda_{X+1}) - \log(\lambda_X) &= \beta_1 \\ \log\left(\frac{\lambda_{X+1}}{\lambda_X}\right) &= \beta_1 \\ \frac{\lambda_{X+1}}{\lambda_X} &= e^{\beta_1} \end{aligned} \tag{1}$$

These results suggest that by exponentiating the coefficient on age we obtain the *multiplicative* factor by which the mean count changes.

Interpreting model coefficients (continued)

In this case, the mean number of children under 5 changes by a factor of $e^{-0.029} = 0.971$ or decreases by 2.9% (since $1 - .971 = .029$) with each additional year older the household head is;

Or, we predict a 3.0% *increase* in mean number of children less than 5 for a 1-year *decrease* in age of the household head (since $1/.971 = 1.0298$).

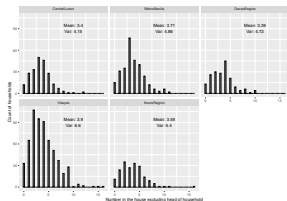
The quantity on the left-hand side of Equation 1 is referred to as a **rate ratio** or **relative risk**, and it represents a percent change in the response for a unit change in X.

Multiple Poisson regression

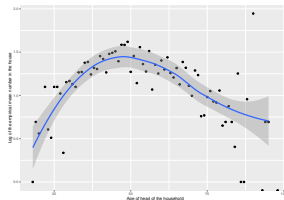
Most significant ideas from multiple linear regression are reinforced by Poisson regression:

- ▶ hypothesis testing
- ▶ confidence intervals
- ▶ indicator variables
- ▶ categorical variables (reference level / Tukey HSD)
- ▶ squared terms
- ▶ interaction terms
- ▶ control for / adjusting for / holding constant
- ▶ checking violations of model conditions
- ▶ maximum likelihood estimates (equals least squares in linear regr)

Exploratory data analysis 2: Total in household



(a) Distribution of household size by Philippine regions.



(b) The log of the mean household sizes by age of the head of household, with loess smoother.

Figure 4: EDA: selected plots

What does the graph on the right tell us about the relationship of λ and age?

Potential final model

```
modela2L <- glm(total ~ age + age2 + location, family = poisson, data = fHH1)
```

##	Estimate	Std. Error	z value	Pr(> z)
## (Intercept)	-0.3843337714	1.820919e-01	-2.1106581	3.480171e-02
## age	0.0703628330	6.905067e-03	10.1900292	2.196983e-24
## age2	-0.0007025856	6.420019e-05	-10.9436677	7.125764e-28
## locationMetroManila	0.0544800704	4.720116e-02	1.1542104	2.484139e-01
## locationDavaoRegion	-0.0193872310	5.378273e-02	-0.3604732	7.184933e-01
## locationVisayas	0.1121091959	4.174960e-02	2.6852758	7.246998e-03
## locationIlocosRegion	0.0609819668	5.265981e-02	1.1580362	2.468493e-01

For example, $\hat{\beta}_6 = -0.0194$ indicates that, after controlling for the age of the head of household, the log mean household size is 0.0194 lower for households in the Davao Region than for households in the reference location of Central Luzon.

In more interpretable terms, mean household size is $e^{-0.0194} = 0.98$ times “higher” (i.e., 2% lower) in the Davao Region than in Central Luzon, when holding age constant.

Maximum estimated additional number in the house occurs when the head of the household is around 50 years old, after adjusting for location.

Potential final model (continued)

To test for the effect of location, use a drop-in-deviance test (analogous to an extra-sum-of-squares F test in linear regression):

##	ResidDF	ResidDev	Deviance	Df	pval
## 1	1497	2200.944	NA	NA	NA
## 2	1493	2187.800	13.14369	4	0.01059463

Adding the four terms corresponding to location to the quadratic model with age produces a statistically significant improvement ($\chi^2 = 13.144$, $df = 4$, $p = 0.0106$), so there is significant evidence that mean household size differs by location, after controlling for age of the head of household.

Comparing non-nested models

Which is the single best predictor? quadratic age, roof type, or location/region?

Use Akaike Information Criterion (AIC) to compare

- ▶ lower values are better
- ▶ differences of 10 are “big”
- ▶ measures model fit while accounting for complexity, analog to adjusted R^2

	df	AIC
modela2	3	6579.823
model_roof	2	6739.269
model_location	5	6727.765

Lack of fit!

When a model is true, we can expect the residual deviance to be distributed as a χ^2 random variable with degrees of freedom equal to the model's residual degrees of freedom.

Our final model has a residual deviance of 2187.8 with 1493 df. The probability of observing a deviance this large if the model fits is essentially 0, saying that there is significant evidence of lack-of-fit.

There are several reasons why **lack-of-fit** may be observed:

- ▶ We may be missing important covariates or interactions.
- ▶ There may be extreme observations that cause the deviance to be larger than expected.
- ▶ There may be a problem with the Poisson model. In particular, the Poisson model has only a single parameter, λ , for each combination of the levels of the predictors which must describe both the mean and the variance.

Overdispersion

Often in Poisson models the variances in the response are larger than the corresponding means at different levels of the predictors. The response is then considered to be **overdispersed**.

Recall that we observed this in our EDA plot by region:

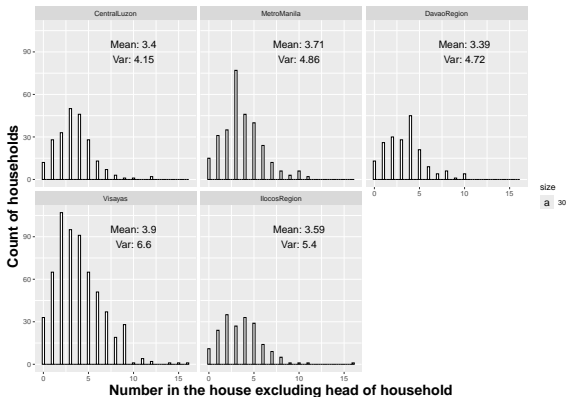


Figure 5: Distribution of household size by Philippine regions.

Quasi-Poisson models

Without adjusting for overdispersion, we use incorrect, artificially small standard errors leading to artificially small p-values for model coefficients.

The simplest way to take overdispersion into account is to use an estimated dispersion factor to inflate standard errors.

$\hat{\phi} = \frac{\sum(\text{Pearson residuals})^2}{n-p}$ where p is the number of model parameters.

$$SE_Q(\hat{\beta}) = \sqrt{\hat{\phi}} * SE(\hat{\beta})$$

Quasi-Poisson models (continued)

##	Estimate	Std. Error	t value	Pr(> t)
## (Intercept)	-0.3843337714	0.2166025174	-1.7743735	7.620514e-02
## age	0.0703628330	0.0082137357	8.5664837	2.616622e-17
## age2	-0.0007025856	0.0000763676	-9.2000473	1.168513e-19
## locationMetroManila	0.0544800704	0.0561468673	0.9703136	3.320474e-01
## locationDavaoRegion	-0.0193872310	0.0639757901	-0.3030401	7.619015e-01
## locationVisayas	0.1121091959	0.0496621109	2.2574392	2.412461e-02
## locationIlocosRegion	0.0609819668	0.0626400545	0.9735299	3.304477e-01
## Residual deviance =	2187.8	on	1493 df	
## Dispersion parameter =	1.414965			

In the absence of overdispersion, we expect the dispersion parameter estimate to be 1.0. The estimated dispersion parameter here is larger than 1.0 (1.415).

For example, the standard error for the Visayas region term from a likelihood based approach is 0.0417, whereas the quasi-likelihood standard error is $\sqrt{1.415} * 0.0417$ or 0.0497. This term is still statistically significant at the 0.05 level under the quasi-Poisson model, but the evidence is not as strong (quasi-Poisson p-value of .024 vs. Poisson p-value of .007).

Quasi-Poisson models (continued)

We can take another look at our final model:

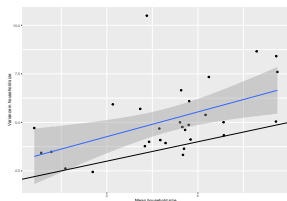
##	ResidDF	ResidDev	F	Df	pval
## 1	1497	2200.944	NA	NA	NA
## 2	1493	2187.800	2.322264	4	0.05477322

Here, after adjusting for overdispersion, we find that there is *not* statistically significant evidence at the 0.05 level ($F = 2.32, p = .055$) that mean household size differs among regions after adjusting for age.

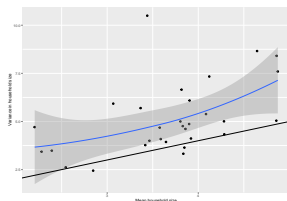
Alternative methods for modeling overdispersion

Diagnostic plot for overdispersion: plot mean vs. variance for groups of households based on predicted size

- ▶ linear with slope $> 1 \Rightarrow$ quasi-Poisson
- ▶ quadratic with incr. slope \Rightarrow negative binomial



(a) Linear fit



(b) Loess smoother

Figure 6: Mean and variance of predicted household sizes.

Negative Binomial models

A negative binomial model introduces another parameter in addition to λ , which gives the model more flexibility and, as opposed to the quasi-Poisson model, the negative binomial model assumes an explicit likelihood model.

Mathematically, you can think of the negative binomial model as a Poisson model where λ is also random, following a gamma distribution.

These results are very similar to the quasi-Poisson model in terms of estimated coefficients (which can change), standard errors, test statistics, and p-values.

Pause to Ponder

Check in with your neighbor(s). What questions do you have at this point? What can we clarify or discuss more fully?

Second case study: Bald eagles

Every year in late December, since 1921, birdwatchers in the Hamilton area of Ontario, Canada, have counted and recorded all the birds they see or hear in a day.

The data was made available by the [Bird Studies Canada website](#) and distributed through the R for Data Science TidyTuesday project.

We are particularly interested in how the Bald Eagle population has changed over time.

Bald eagle data

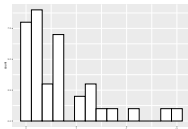
Each row of `bald_eagles.csv` contains information about bald eagles counts in Hamilton, Ontario, for one year. There are 37 rows covering 1981 through 2017. The variables include:

- ▶ `year` = year of data collection
- ▶ `count` = number of birds observed
- ▶ `hours` = total person-hours of observation period
- ▶ `count_per_hour` = count divided by hours
- ▶ `count_per_week` = `count_per_hour` multiplied by 168 hours per week

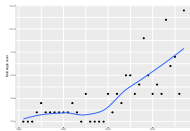


Credit: © Ron Niebrugge/wildnatureimages

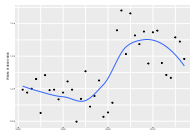
Exploratory data analysis



(a) Histogram of number of bald eagle sightings by year.



(b) Bald eagle counts vs. year



(c) Hours of observation vs. year

Figure 7: EDA: selected plots

Sampling effort

Poisson random variables are often used to represent counts (e.g., number of bald eagles) *per unit of time or space* (previously, number of people in one household).

But what if observation (sampling) effort (as measured by the number of weeks people observed birds) is changing over time?

We cannot directly compare the 2 eagles observed in 1985 to the 7 eagles observed in 2015 when there were only 143 person-hours (0.85 person-weeks) of observation in 1985 compared with 221 person-hours (1.32 person-weeks) in 2015.

We should examine time trends in the *rate* of bald eagles sightings; for example, we will calculate the bald eagle counts per week ($\frac{\text{number of bald eagles}}{\text{hours of observation}} \cdot (168 \text{ hours/week})$).

Offsets

We let λ_i be the expected number of eagles in year i per with weeks_i weeks observed in year i .

Then λ_i/weeks_i is the number of eagles expected in a week!

Adjusting the yearly count by observation time is equivalent to adding $\log(\text{weeks})$ to the right-hand side of the Poisson regression equation—essentially adding a predictor with a fixed coefficient of 1, called an **offset**:

$$\log\left(\frac{\lambda_i}{\text{weeks}_i}\right) = \beta_0 + \beta_1 x_i$$

$$\log(\lambda_i) - \log(\text{weeks}_i) = \beta_0 + \beta_1 x_i$$

$$\log(\lambda_i) = \beta_0 + \beta_1 x_i + \log(\text{weeks}_i)$$

Thus, modeling $\log(\lambda)$ and adding an offset is equivalent to modeling **rates**, and coefficients can be interpreted in terms of rates.

Interpretation

Ordinary Poisson:

$$\log(\text{COUNT}) = \log(\lambda) = \beta_0 + \beta_1 X$$

The average *number of eagles* has increased by ...% per year.

Poisson with Offset:

$$\log(\text{RATE}) = \log(\lambda/t) = \beta_0 + \beta_1 X$$

The average *number of eagles per week* has increase by ...% per year.

Modeling results

We are interested primarily in trends over time in eagle sightings. We have no control variables other than sampling effort, so we simply fit a model with year (centered at 1981) and our offset.

```
model_eagles <- glm(count ~ year_1981, family = poisson,  
                    offset = log(weeks), data = eagles)
```

```
##              Estimate Std. Error   z value    Pr(>|z|)  
## (Intercept) -0.79971104 0.31969178 -2.501506 1.236662e-02  
## year_1981    0.07566483 0.01155782  6.546634 5.884818e-11  
  
## Residual deviance = 42.39031 on 35 df  
## Dispersion parameter = 1
```

Bald eagle counts are significantly increasing over time ($Z = 6.55$, $p < .001$), even after adjusting for observation time. The average eagle sighting rate per week has grown about 7.9% per year (since $e^{0.0757} = 1.0786$) in Hamilton, Ontario.

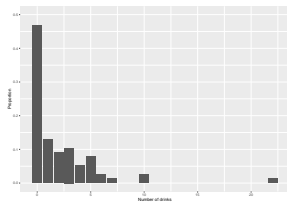
Adjustments for potential overdispersion using either quasi-Poisson or negative binomial regression provide minimal changes to model coefficients and tests.

Other extensions of the Poisson regression model

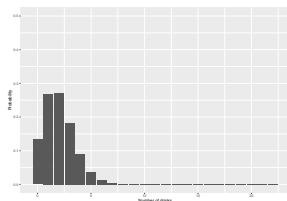
- ▶ Zero-inflated models
- ▶ Hurdle models
- ▶ Multilevel models

Zero-inflated Poisson models: Weekend drinking

An informal survey of students in an intro stats course included the question, “How many alcoholic drinks did you consume last weekend?”. We wish to identify factors associated with increased drinking.



(a) Actual distribution



(b) Poisson distribution with same mean

Figure 8: Count of drinks consumed last weekend

There are more zeros than expected under a Poisson model.

ZIP: Weekend drinking (continued)

Our zeros are a **mixture** of responses from non-drinkers (who would always report 0) and drinkers who abstained during the past weekend. Ideally, we'd like to sort out the non-drinkers and drinkers when performing our analysis.

Define λ to be the mean number of drinks *among those who drink*, and α to be the proportion of *non-drinkers* (“true zeros”).

Model λ and α (or functions of λ and α) simultaneously using covariates like sex, first-year status, and off-campus residence. For example:

$$\log(\lambda) = \beta_0 + \beta_1 \text{offcampus} + \beta_2 \text{sex}$$

$$\log(\alpha/(1 - \alpha)) = \beta_0 + \beta_1 \text{firstyear}$$

The first part of a ZIP model is a regular Poisson regression model, and the second part is a logistic regression model.

Hurdle models: Going vague

In a 2018 study, Chapp et al. scraped every issue statement from webpages of candidates for the U.S. House of Representatives, counting the number of issues candidates commented on and scoring the level of ambiguity of each statement.

Research questions:

- ▶ Which candidates for U.S. House are more likely to have at least one issue page and to offer statements on a greater number of issues?
- ▶ How are a candidate's political party, incumbency status, and political beliefs related to their willingness to post stands and ideas on issues?
- ▶ How do the demographics and political beliefs of voters in the candidate's district impact a candidate's willingness to engage?
- ▶ How does the interplay between candidate profile and voter profile affect a candidate's willingness to comment on issues?

Hurdle: Going vague (continued)

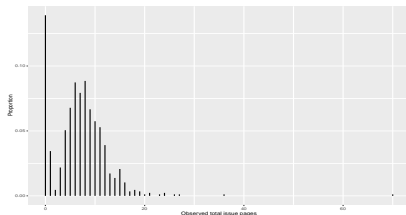


Figure 9: Total issue pages for 2014 US House candidates

Once again, there are more zeros than expected under a Poisson model. But *unlike* ZIP models, it is not natural to consider these zeros to be a mixture. Candidates decide either to post issue statements or not. Those who decide to not post any issue statements comprise our zeros, and for those who decide to post issue statements, we can model the number they choose to post.

Since those who decide to post issue statements “leap over” the zero category, these models are referred to as **hurdle models**.

Hurdle: Going vague (continued)

Similar to ZIP models, we will define λ to be the mean number of drinks *among those who posted at least one issue page*, and α to be the proportion of *candidates who post at least one issue page* (“true non-zeros” or “true hurdlers”).

Then, model λ and α (or functions of λ and α) simultaneously using characteristics of the candidates and their districts:

$$\log(\lambda) = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

$$\log(\alpha/(1 - \alpha)) = \beta_0 + \beta_1 X_1 + \beta_2 X_3$$

Hurdle models focus on modeling *non-zeros*, whereas ZIP models focus on modeling *true zeros*. This really only affects the interpretation of logistic coefficients.

The count portion of the hurdle model is actually based on a **truncated Poisson** distribution (domain starting at 1) rather than a full Poisson distribution (domain starting at 0). This does not affect the interpretation of Poisson coefficients.

Multilevel modeling: Going vague

Problem: Even with a hurdle model for number of issue pages, a condition is still violated in the Going Vague example.

Observations are *not independent*! Some covariates are measured at the candidate level (incumbent, party, ideology), while others are measured at the district level (demographics, ideology of voters). Candidates from the same district will have the same values for any covariate at the district level.

Implications: overstate effective sample size, underestimate standard errors, and overstate significance of covariates

Solution: multilevel (hierarchical / mixed effects) modeling!

Idea: build a regression model for issue pages at the candidate level, and then build another regression model for coefficients from the first model using covariates at the district level. Combine into a single composite model.

Pause to Ponder

With your neighbor(s) discuss plans, ideas, questions, and concerns you have about teaching Poisson regression somewhere in your own world.

Preview of materials

GitHub repo for this session, containing:

- ▶ slides for this presentation (including R source code and data)
- ▶ St. Olaf SDS 316 class folder with class activities for the Poisson regression unit (keys available upon request).
- ▶ Draft of brand new Chapter 4 for BMLR2e (Roback, Boehm Vock, and Legler; expected Fall 2026)

Beyond Multiple Linear Regression: Applied Generalized Linear Models and Multilevel Models in R by Roback and Legler (2021).

Thanks!

Please be in touch with any questions, thoughts, feedback, etc.!

Laura Boehm Vock: (boehm@stolaf.edu)

Paul Roback: (roback@stolaf.edu)

Bonus material: Interpreting ZIP

\$count

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	0.7542747	0.1440035	5.237891	1.624221e-07
offcampus	0.4159420	0.2058602	2.020507	4.333078e-02
sexm	1.0209002	0.1751937	5.827264	5.634331e-09

\$zero

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.6036151	0.3114485	-1.938090	0.05261230
firstyear	1.1363880	0.6095155	1.864412	0.06226388

$e^{0.415} = 1.51$: Students who drink consume 51% more drinks if living off campus, after accounting for sex.

$e^{1.136} = 3.11$: First year students have 3.11 times greater odds of *not drinking* compared to older students.

Bonus material: Interpreting hurdle models

incumbent: $e^{-1.2085} = 0.299$ and $1/e^{-1.2085} = 3.35$. Thus, the odds a challenger posts at least one issue page are 3.35 times greater than the odds an incumbent posts at least one issue page, holding all else constant.

demHeterogeneity: $e^{-1.173} = 0.309$ and $1/e^{-1.173} = 3.23$. Thus, among candidates who choose to post at least one issue page, the mean number of issue pages posted are 3.23 times greater with each one unit decrease in demographic heterogeneity score, holding all else constant.

ideology:democrat: $1/e^{-0.490} = 1.632$ and $1/e^{-0.490+0.222} = 1.307$. Thus, for each 1 unit decrease in ideology (more liberal), the mean number of issue pages (for candidates with at least one issue page) increases by 63.2% for democrats but only 30.7% for republicans.

Overall candidates are disincentivized to take public stances on issues if they are an incumbent, if their constituents have a wide ranges of backgrounds, and if their beliefs are less aligned with their constituents.