

# Lecture 20

# Data from last time

Survey data from 77 college students on a dry campus (i.e., alcohol is prohibited) in the US. Survey asks students “How many alcoholic drinks did you consume last weekend?”

- `drinks`: number of drinks the student reports consuming
- `sex`: whether the student identifies as male
- `OffCampus`: whether the student lives off campus
- `FirstYear`: whether the student is a first-year student

Our goal: model the number of drinks students report consuming.

# Recap: Poisson hurdle model

$y_i = \# \text{ drinks student } i \text{ reported consuming}$

$$P(Y_i > 0) = p_i \quad \log\left(\frac{p_i}{1-p_i}\right) = \gamma^T X_i$$

$$Y_i | (Y_i > 0) \sim \text{PosPoisson}(\lambda_i) \quad \log(\lambda_i) = \beta^T X_i$$

(aka ZTP( $\lambda_i$ ))

$$P(Y_i = y) = \begin{cases} 1 - p_i & y=0 \\ p_i \frac{e^{-\lambda_i} \lambda_i^y}{y! (1-e^{-\lambda_i})} & y > 0 \end{cases} \quad \left. \begin{array}{l} \text{involves both} \\ \gamma \text{ and } \beta \end{array} \right.$$

$$\sum_{y=0}^{\infty} P(Y_i = y) = P(Y_i = 0) + \sum_{y=1}^{\infty} P(Y_i = y)$$

$$= 1 - p_i + p_i \sum_{y=1}^{\infty} \underbrace{\frac{e^{-\lambda_i} \lambda_i^y}{y! (1-e^{-\lambda_i})}}_{= 1 \text{ (PosPoisson)}}$$

# Fitting Poisson hurdle models

$$P(Y_i=y) = \begin{cases} 1 - \rho_i & y=0 \\ \rho_i \frac{e^{-\lambda_i} \lambda_i^y}{y! (1-e^{-\lambda_i})} & y > 0 \end{cases}$$

$$Z_i = \mathbb{1}\{\gamma_i > 0\}$$

$$P(Y_i=y) = (1-\rho_i)^{1-Z_i} \rho_i^{Z_i} \left( \frac{e^{-\lambda_i} \lambda_i^y}{y! (1-e^{-\lambda_i})} \right)^{Z_i}$$

$$\ell(\beta, \gamma) = \sum_{i=1}^n \left\{ Z_i \log \rho_i + (1-Z_i) \log (1-\rho_i) \right\} + \sum_{i=1}^n \left\{ Z_i (\gamma_i \log \lambda_i - \log(e^{\lambda_i} - 1)) - Z_i \log \gamma_i \right\}$$

$$= \sum_{i=1}^n \left\{ Z_i \log \left( \frac{\rho_i}{1-\rho_i} \right) + \log (1-\rho_i) \right\}$$

$$+ \sum_{i=1}^n \left\{ Z_i (\gamma_i \log \lambda_i - \log(e^{\lambda_i} - 1)) + Z_i \log \gamma_i \right\}$$

$$\frac{\partial L}{\partial \gamma} = X^T(Z - P)$$

$$X = \text{design matrix}$$

$$Z = \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix}$$

$$P = \begin{bmatrix} p_1 \\ \vdots \\ p_n \end{bmatrix}$$

$$\frac{\partial L}{\partial \beta} = X^T \text{diag}(z_i)(Y - \mu)$$

$$\mu_i = \mathbb{E}[Y_i | (Y_i > 0)]$$

$$= \frac{\lambda_i}{1 - e^{-\lambda_i}}$$

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

$$\text{set } = 0$$

$$\Rightarrow U(\beta, \gamma) = \begin{bmatrix} X^T \text{diag}(z_i)(Y - \mu) \\ X^T(Z - P) \end{bmatrix}$$

can solve for  $\beta, \gamma$  separately

$$P_i = P(Y_i \geq 0)$$

# Fitting the model in R

```

1 library(pscl)
2
3 m1 <- hurdle(drinks ~ sex + FirstYear + OffCampus,
4                 dist = "poisson", zero.dist = "binomial",
5                 data = wdrinks)
6
7 m1$coefficients

```

\$count

	(Intercept)	sexm	FirstYearTRUE	OffCampusTRUE
	0.8132113	0.9706640	-0.2181068	0.3762608

Pos Poisson  
part

\$zero

	(Intercept)	sexm	FirstYearTRUE	OffCampusTRUE
	0.1230510	0.3377969	-0.8554289	1.5803472

logistic  
regression  
part

Among more students who drink, male students report drinking  
(holding First Year  $\setminus$  off campus fixed)

Male students are more likely to report drinking  
(holding First Year  $\setminus$  off campus fixed)

The odds of any drinking are  $e^{0.338}$  times higher for  
a male student (holding ... fixed)

# Model assumptions

```
1 m1 <- hurdle(drinks ~ sex + FirstYear + OffCampus,  
2                  dist = "poisson", zero.dist = "binomial",  
3                  data = wdrinks)  
4  
5 m1$coefficients
```

\$count

	(Intercept)	sexm	FirstYearTRUE	OffCampusTRUE
	0.8132113	0.9706640	-0.2181068	0.3762608

\$zero

	(Intercept)	sexm	FirstYearTRUE	OffCampusTRUE
	0.1230510	0.3377969	-0.8554289	1.5803472

What assumptions does this model make?

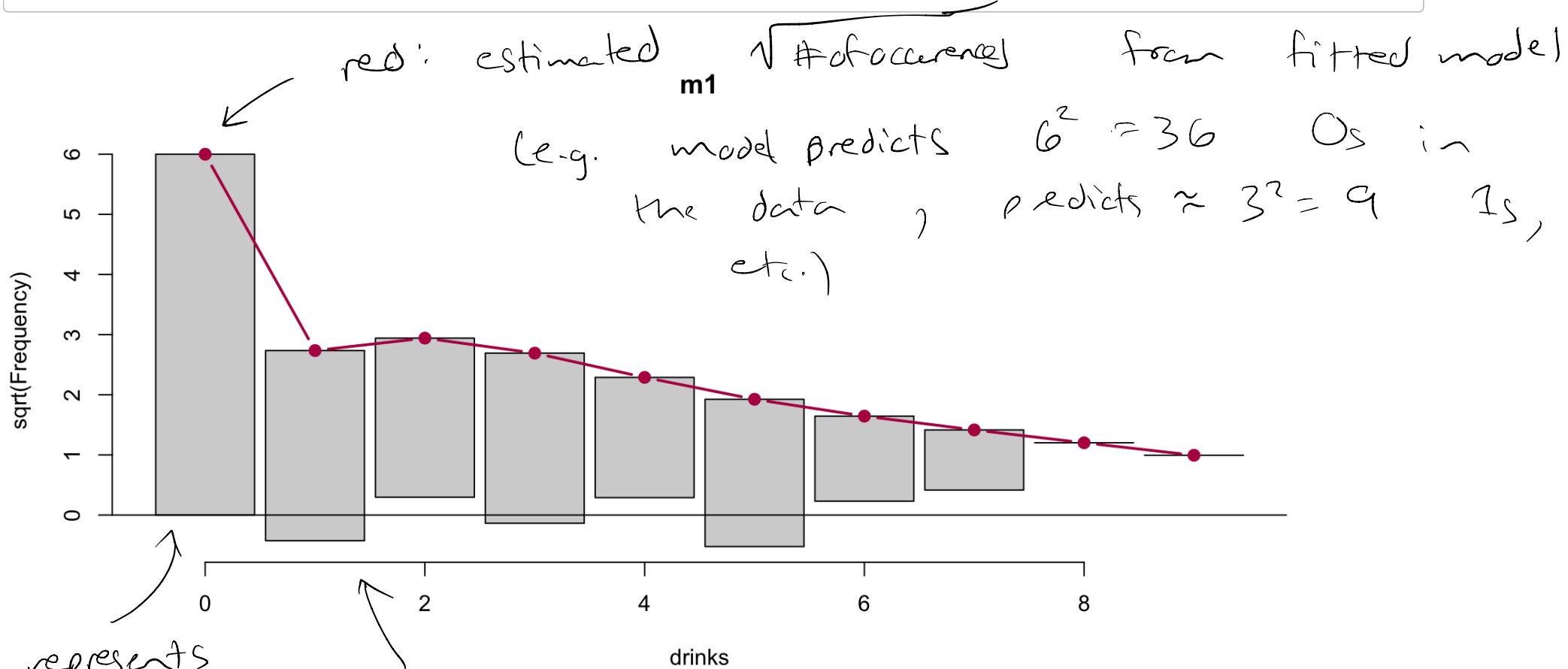
- Distribution:  $y_i | (y_i > 0) \sim \text{PosPoisson}(\lambda_i)$
- Independence (needed for ML)
- Shape:  $\log(\lambda_i) = \beta^T x_i$   
 $\log\left(\frac{\rho_i}{1-\rho_i}\right) = \gamma^T x_i$

# Model diagnostics

- Quantile residual plots
  - For the whole model ( $\text{distribution of } Y_i$ )
  - For the zero component ( $\mathbb{1}\{Y_i \geq 0\}$ )
  - For the count component ( $Y_i | \{Y_i \geq 0\}$ )
- Rootogram (for assessing distributional assumption)

# Rootograms

```
1 library(countreg)  
2 rootogram(m1)
```

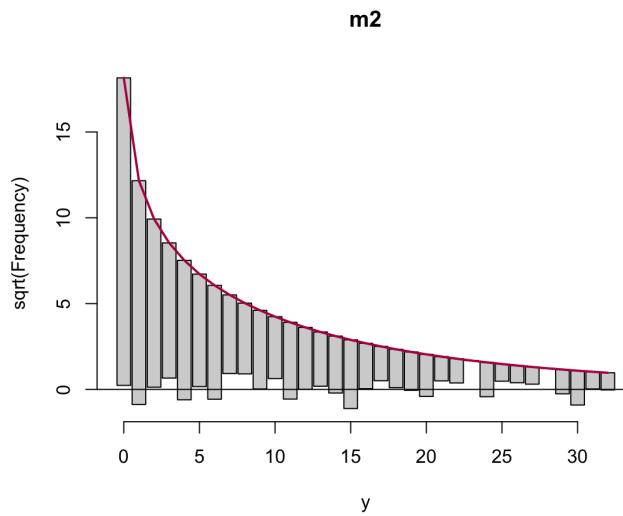
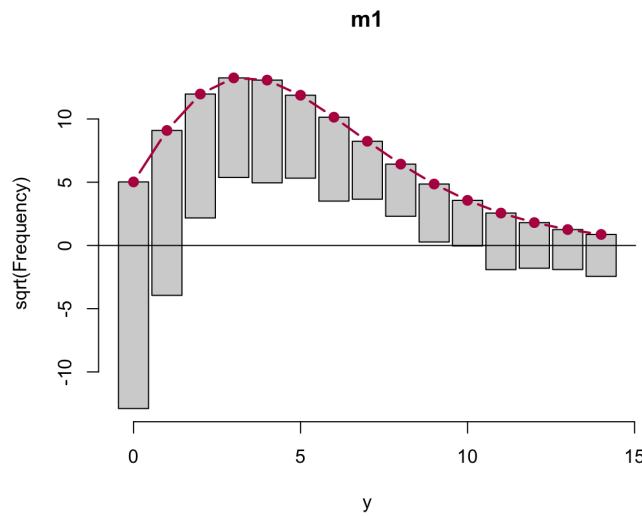


bar represents  
actual (observed)  
 $\sqrt{\# \text{occurrences}}$

bar below 0: ↑  
observed more occurrences than predicted on the horizontal

# Other examples with rootograms

```
1 par(mfrow=c(1, 2))
2
3 x <- rnorm(1000)
4 y <- rnbinom(1000, 0.5, mu=exp(1.5 + 0.2*x))
5 m1 <- glm(y ~ x, family = poisson)
6 m2 <- glm.nb(y ~ x)
7
8 rootogram(m1)
9 rootogram(m2)
```



suggests  
overdispersion

# Hurdle models for count data

$$P(Y_i > 0) = p_i \quad g_{\text{zero}}(p_i) = \gamma^T X_i$$

$$Y_i | (Y_i > 0) \sim ZT(\lambda_i) \quad g_{\text{count}}(\lambda_i) = \beta^T X_i$$

Usually:

$$g_{\text{zero}}(p_i) = \log\left(\frac{p_i}{1-p_i}\right)$$

Usually:

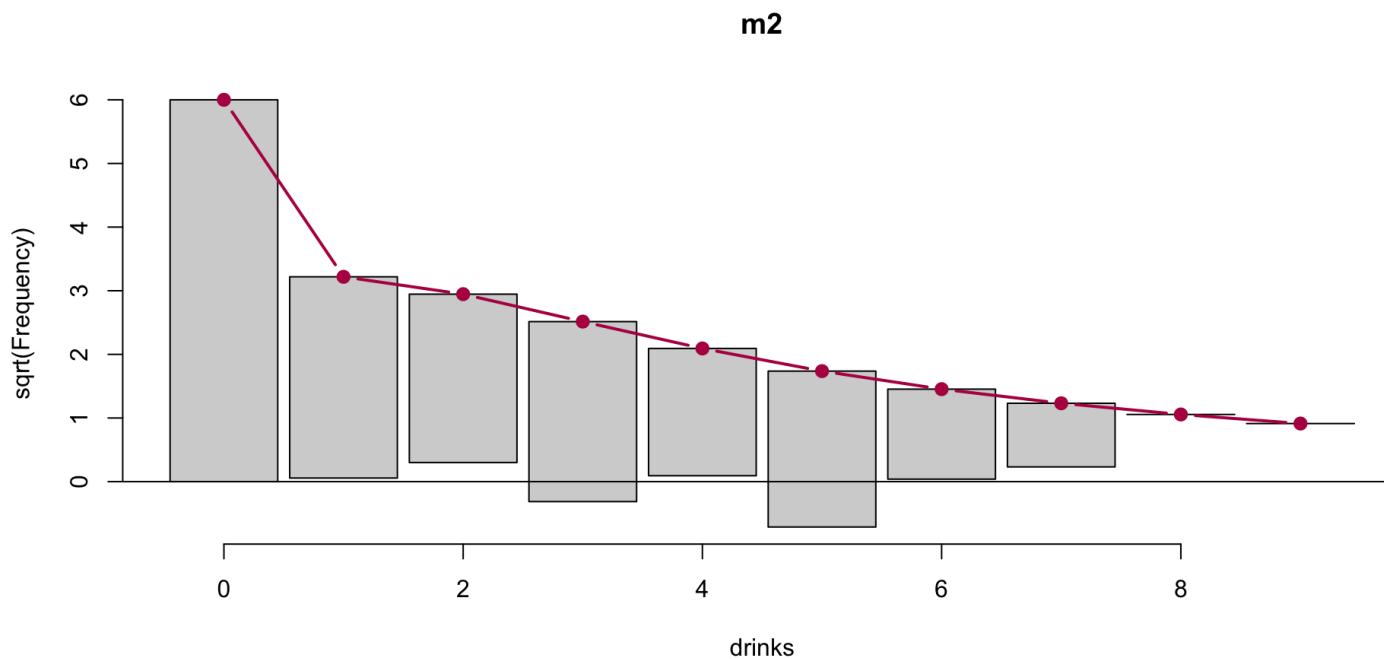
$$g_{\text{count}}(\lambda_i) = \log(\lambda_i)$$

$ZT(\lambda_i)$  : zero-truncated distribution

$\lambda_i$  = mean of non-truncated dist.

# Negative binomial hurdle model

```
1 m2 <- hurdle(drinks ~ sex + FirstYear + OffCampus,  
2                  dist = "negbin", zero.dist = "binomial",  
3                  data = wdrinks)  
4  
5 rootogram(m2)
```



# Class activity

[https://sta712-f23.github.io/class\\_activities/ca\\_lecture\\_20.html](https://sta712-f23.github.io/class_activities/ca_lecture_20.html)

