

# Lecture 28

# Recap: multinomial regression model

data  $(x_1, y_1), \dots, (x_n, y_n)$

$y_i \sim \text{Categorical}(\pi_{i1}, \dots, \pi_{iJ})$

$$\mu_i = (\pi_{i1}, \dots, \pi_{i,J-1})^T \in \mathbb{R}^{J-1}$$

$$g(\mu_i) = \begin{pmatrix} \log \left( \frac{\pi_{i1}}{1 - \sum_{j=1}^{J-1} \pi_{ij}} \right) \\ \vdots \\ \log \left( \frac{\pi_{i,J-1}}{1 - \sum_{j=1}^{J-1} \pi_{ij}} \right) \end{pmatrix} = \begin{pmatrix} \beta_1^T x_i \\ \vdots \\ \beta_{J-1}^T x_i \end{pmatrix}$$

$$\begin{bmatrix} x_i^T \\ x_i^T \\ \vdots \\ x_i^T \end{bmatrix} = X_i^* \beta \quad \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_{J-1} \end{bmatrix}$$

# Motivating example: earthquake data

We have data from the 2015 Gorkha earthquake in Nepal. After the earthquake, a large scale survey was conducted to determine the amount of damage the earthquake caused for homes, businesses and other structures. Variables

include:  $Damage_i \sim \text{Categorical}(\pi_{i(\text{none})}, \pi_{i(\text{mod})}, \pi_{i(\text{severe})})$

- **Damage:** the amount of damage suffered by the building (none, moderate, severe)
- **age:** the age of the building (in years)
- **condition:** a de-identified variable recording the condition of the land surrounding the building

$$\log\left(\frac{\pi_{i(\text{mod})}}{\pi_{i(\text{none})}}\right) = \beta_{(\text{mod})}^T X_i$$

$$\log\left(\frac{\pi_{i(\text{severe})}}{\pi_{i(\text{none})}}\right) = \beta_{(\text{severe})}^T X_i$$

# Exploratory data analysis

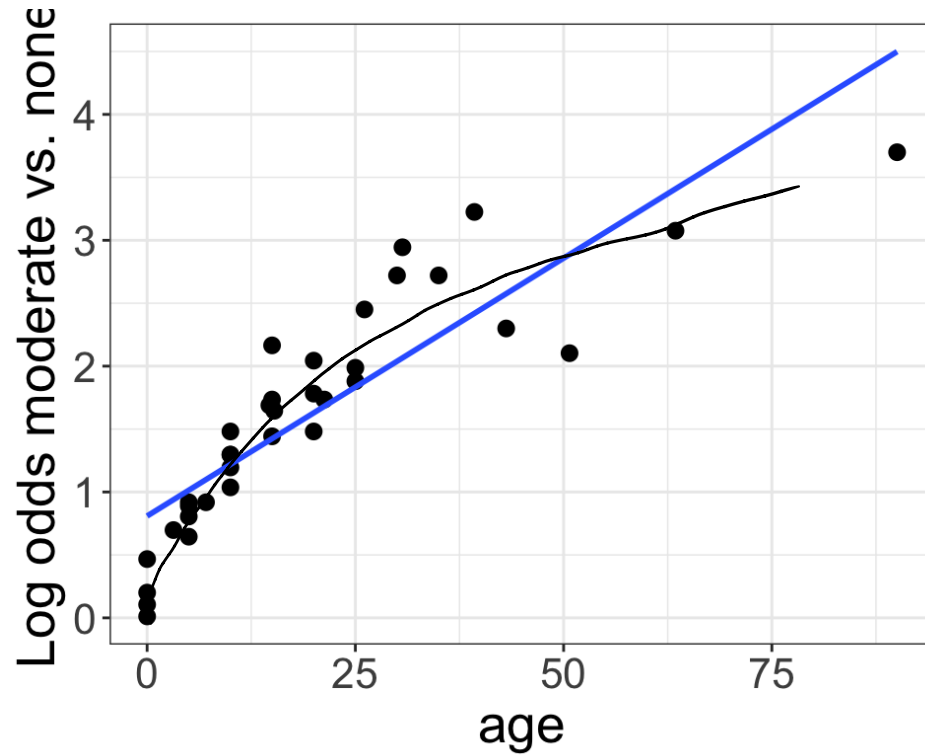
We want to model damage using age and land surface condition. What kind of EDA could I do?

Empirical      logit plots!

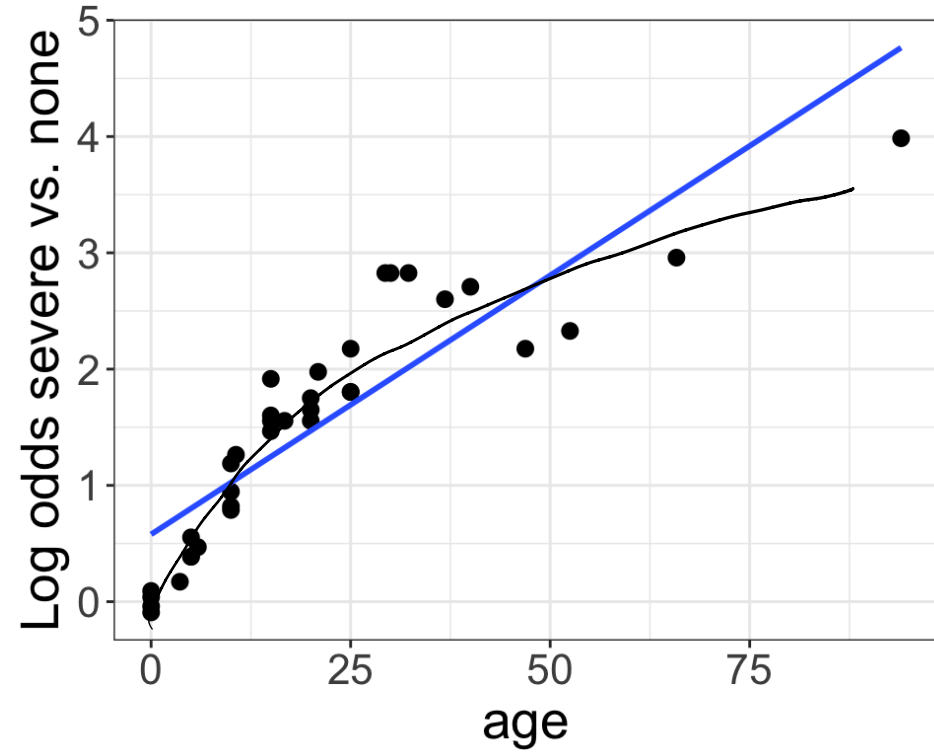
Compare              Moderate vs. None

Severe              vs. None

# Empirical logit plots

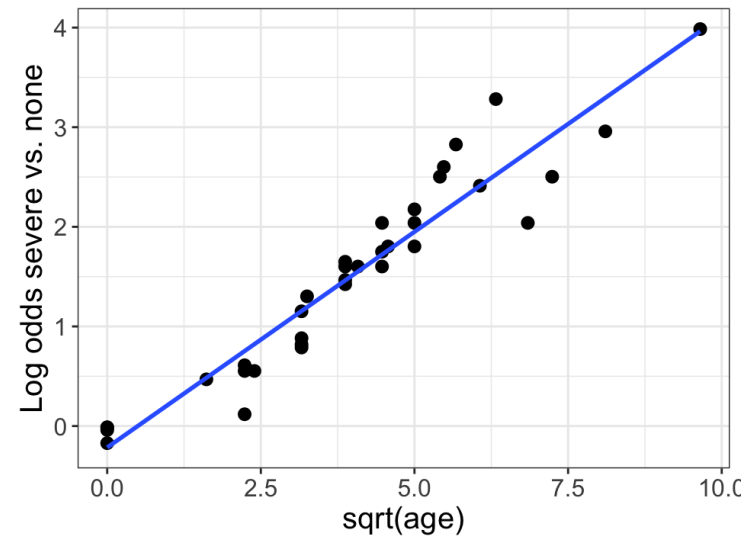
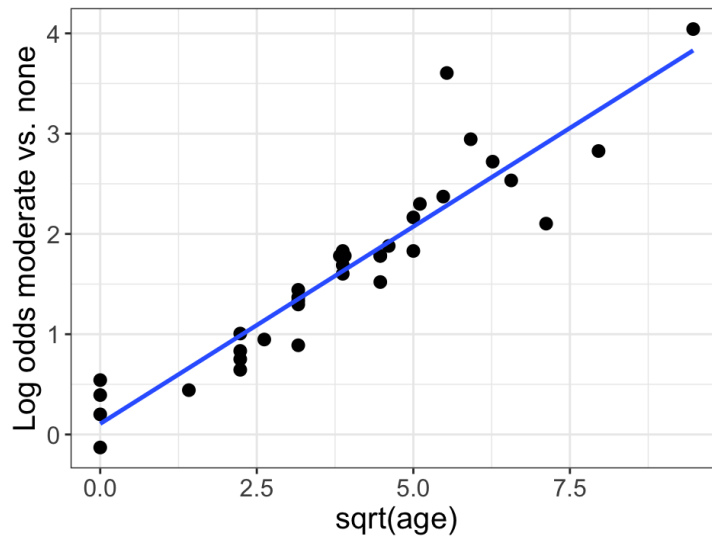
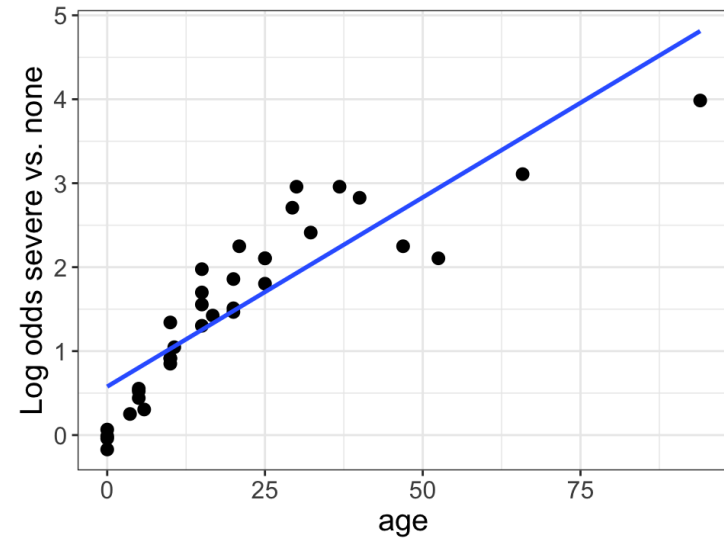
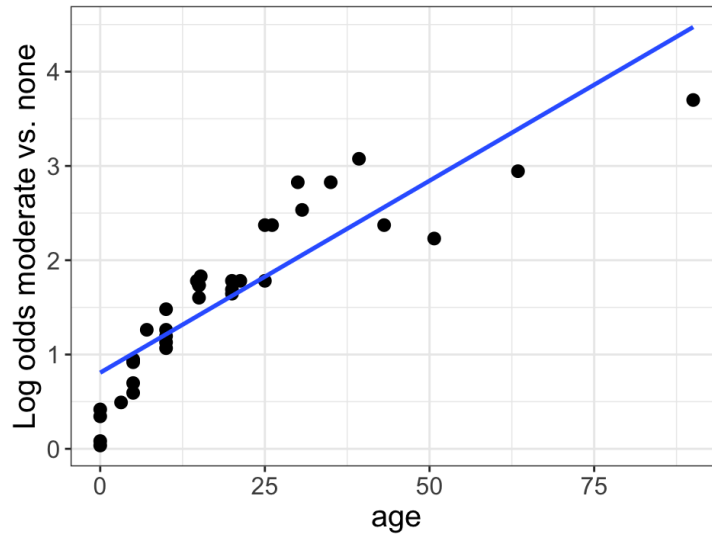


Try a transformation?



Square root?

# Trying a transformation



# Fitting the model in R

```
1 library(nnet)
2 m1 <- multinom(Damage ~ sqrt(age) +
3               condition,
4               data = earthquake)
```

```
1 summary(m1)
```

...

Coefficients:

	(Intercept)	sqrt(age)	conditiono	conditiont
moderate	0.6581163	0.3747641	<del>-0.45376940</del>	<del>-0.5803708</del>
severe	0.1881145	0.4251732	0.04706934	-0.4623774

} coefficients for logits

Std. Errors:

	(Intercept)	sqrt(age)	conditiono	conditiont
moderate	0.1208913	0.01684468	0.2305975	0.1155475
severe	0.1243799	0.01725782	0.2292533	0.1180182

} standard errors for coefficients

...

$$\log \left( \frac{\hat{\pi}_i(\text{mod})}{\hat{\pi}_i(\text{None})} \right) = 0.658 + \underbrace{0.375}_{\text{sqrt(age)}} \sqrt{\text{Age}_i} - 0.454 O_i - 0.580 T_i$$

A one-unit increase in  $\sqrt{\text{Age}_i}$  is associated w/ an increase in the odds of moderate vs. no damage by a factor of  $e^{0.375} \approx 1.45$ , holding surface condition fixed.

# Class activity

<https://sta712->

[f23.github.io/class\\_activities/ca\\_lecture\\_28.html](https://sta712-f23.github.io/class_activities/ca_lecture_28.html)



# Class activity

1) odds:  $\exp \{ 0.658 + 0.375 \sqrt{25} - 0.454 \}$

$\approx 8$

(moderate is 8 times as likely as no damage)

2)  $\hat{\pi}_{i(\text{moderate})}$

In general:  $\hat{\pi}_{ij} = \frac{\text{odds}(j \text{ vs. } \bar{J})}{1 + \sum_{k=1}^{J-1} \text{odds}(k \text{ vs. } \bar{J})}$

$$\frac{p_m}{p_n} = 7.996$$

by similar logic:

$$\frac{p_s}{p_n} = 10.591$$

$$\frac{p_n}{p_n} = 1$$

implies  $7.996 p_n = p_m$

$$10.591 p_n = p_s$$

$$p_m + p_n + p_s = 1$$

$$p_n (7.996 + 10.591 + 1) = 1 \Rightarrow p_n = .051$$

$$7.996 p_n = .408 = p_m$$

# Fisher scoring for multinomial regression

