

- Final exam: Wednesday, May 3 2pm - 5pm  
in Manchester 121 (usual room)
  - cumulative
  - closed notes
  - bring a calculator
- Extra credit seminar (Extra credit on Final exam):  
Dr. Emily Griffith  
Manchester 121  
Thursday, 4/27 11am - 12pm

## Final Exam Review

### 1 Logistic Regression

#### 1.1 Cancer cells

In a study of patients with breast tumors, scientists were interesting in determining the relationship between the size of tumors in centimeters (X) found on lymph nodes and whether or not the tumor was cancerous (Y). Let  $Y_i = 1$  if patient  $i$  in the study has a tumor that is cancerous, and  $Y_i = 0$  if the tumor is not cancerous. Let  $Size_i$  be the size of the tumor of patient  $i$  in centimeters.

1. Write down the appropriate logistic regression model.

$$Y_i \sim \text{Bernoulli}(\pi_i)$$

$$\log\left(\frac{\pi_i}{1-\pi_i}\right) = \beta_0 + \beta_1 Size_i$$

2. The scientists fit the logistic regression model and obtain the following line:

$$\log\left(\frac{\hat{\pi}_i}{1-\hat{\pi}_i}\right) = -2.086 + 0.5117 Size_i.$$

Interpret the slope in terms of the log odds.

An increase of 1 cm in tumor size is associated with an increase of 0.5177 in the log odds of the tumor being cancerous.

$$\log\left(\frac{\hat{\pi}_i}{1-\hat{\pi}_i}\right) = -2.086 + 0.5117 \text{Size}_i.$$

3. Interpret the slope in terms of the odds.

An increase in size of 1 cm is associated with an increase in the odds by a factor of  $e^{0.5117} = 1.668$

4. What is predicted log odds that a tumor is cancerous for a patient with a tumor of size 5 cm?

$$\begin{aligned}\log\left(\frac{\hat{\pi}_i}{1-\hat{\pi}_i}\right) &= -2.086 + 0.5117(5) \\ &= 0.4725\end{aligned}$$

5. Based on your answer to Question 4, is the predicted probability that a tumor of size 5 cm is cancerous less than 50%, greater than 50%, or equal to 50%? Explain your reasoning. Note: You should perform no calculations.

Greater than 50%  
(because log odds > 0)

$$\log \left( \frac{\hat{\pi}_i}{1 - \hat{\pi}_i} \right) = -2.086 + 0.5117 \text{Size}_i.$$

6. What is the predicted probability that a tumor of size 7 cm is cancerous?

$$\begin{aligned}\hat{\pi}_i &= \frac{e^{-2.086 + 0.5117(7)}}{1 + e^{-2.086 + 0.5117(7)}} \\ &= 0.817\end{aligned}$$

7. What are the predicted odds that a tumor of size 7 cm is cancerous?

$$\begin{aligned}\frac{\hat{\pi}_i}{1 - \hat{\pi}_i} &= e^{-2.086 + 0.5117(7)} \\ &= 4.463\end{aligned}$$

## 1.2 Bird nests

A study was conducted to determine what factors contribute to a bird choosing to build a closed nest (a nest that is sealed except for a small opening) versus the traditional, bowl shaped open nest. We have information on  $n = 83$  bird species. Let  $Y_i = 1$  if a species builds a closed nest, and  $Y_i = 0$  otherwise.

We use the following predictors:

- **Length** : the mean body length of the species in cm.
- **Color**: takes the value 1 if the species lay colored eggs, and takes 0 if the species lay brown or white eggs.

We fit a logistic regression model (Model 1) and obtain the following output. You may assume the relationship between length and the log odds of making a closed nest is linear.

### Model 1

	Estimate	Std. Error	z value	Pr(>  z )
(Intercept)	2.0798	1.0468	1.99	0.0469
Length	-0.1709	0.0636	-2.69	0.0072

Null deviance: 103.199 on 82 degrees of freedom  
 Residual deviance: 93.591 on 81 degrees of freedom  
 AIC: 97.591

1. Do the data provide convincing evidence of a relationship between the length of a bird and the log odds of building a closed nest? Use a drop-in-deviance test to answer this question. Show all your steps. (The p-value is 0.001937)

Model:  $Y_i \sim \text{Bernoulli}(\pi_i)$   
 $\log\left(\frac{\pi_i}{1-\pi_i}\right) = \beta_0 + \beta_1 \text{Length}_i$

$$H_0: \beta_1 = 0 \quad H_A: \beta_1 \neq 0$$

Test statistic:  $G = 103.199 - 93.591 = 9.608$

$$\text{p-value} = P(\chi^2_1 > 9.608) = 0.001937$$

So we have strong evidence for a relationship between length & the log odds of building a closed nest

Now we are going to switch to a new predictor, color.

### Model 2

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-0.5878	0.5578	-1.05	0.2920
Color	-0.2389	0.6161	-0.39	0.6982

Null deviance: 103.199 on 82 degrees of freedom  
 Residual deviance: 103.05 on 81 degrees of freedom  
 AIC: 107.05

2. Build and interpret a 95% confidence interval for the slope.

$$\hat{\beta}_1 \pm z^* SE_{\hat{\beta}_1}$$

$$\hookrightarrow -0.2389 \pm 1.96(0.6161) = (-1.446, 0.969)$$

we are 95% confident that species laying colored eggs have a log odds of building a closed nest between 1.446 lower and 0.969 higher than species laying uncolored eggs

3. Is there convincing evidence of a relationship between the color of the eggs and the log odds of a bird species making a closed nest? Use a z-test to answer this question. Show your steps and clearly state your conclusion in context of the data.

Model:  $Y_i \sim \text{Bernoulli}(\pi_i)$   
 $\log\left(\frac{\pi_i}{1-\pi_i}\right) = \beta_0 + \beta_1 \text{Color}_i$

$$H_0: \beta_1 = 0 \quad H_A: \beta_1 \neq 0$$

$$z = \frac{-0.2389}{0.6161} = -0.39$$

$$p\text{-value} = 0.6982$$

So we have very weak evidence for a relationship between color & the log odds of building a closed nest

### Model 1

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	2.0798	1.0468	1.99	0.0469
Length	-0.1709	0.0636	-2.69	0.0072

Null deviance: 103.199 on 82 degrees of freedom  
Residual deviance: 93.591 on 81 degrees of freedom  
AIC: 97.591

### Model 2

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-0.5878	0.5578	-1.05	0.2920
Color	-0.2389	0.6161	-0.39	0.6982

Null deviance: 103.199 on 82 degrees of freedom  
Residual deviance: 103.05 on 81 degrees of freedom  
AIC: 107.05

4. If you could only choose one predictor, Length or Color, which would you choose and why?

Length. The model with length has a lower AIC than the model with color

## 2 Maximum likelihood estimation

If a Poisson distribution counts the number of events that occur in a fixed interval of time, then the length of time between each event follows what is called an *exponential* distribution. Suppose that  $Y \sim \text{Exponential}(\lambda)$  is an exponential random variable, with parameter  $\lambda$ . We observe  $n$  observations  $Y_1, \dots, Y_n$ , and we want to estimate  $\lambda$ . The likelihood of an estimate  $\hat{\lambda}$  is given by

$$L(\hat{\lambda}) = \prod_{i=1}^n \hat{\lambda} e^{-\hat{\lambda} Y_i}$$

Calculate the maximum likelihood estimate of  $\lambda$ . Show all steps.

$$L(\lambda) = \prod_{i=1}^n \lambda e^{-\lambda Y_i}$$

$$\textcircled{1} \text{ Take log: } \log L(\lambda) = \sum_{i=1}^n \log (\lambda e^{-\lambda Y_i}) = \sum_{i=1}^n (\log(\lambda) - \lambda Y_i)$$

$$\textcircled{2} \text{ Differentiate: } \frac{\partial}{\partial \lambda} \log L(\lambda) = \frac{\partial}{\partial \lambda} \sum_{i=1}^n (\log(\lambda) - \lambda Y_i) = \sum_{i=1}^n \frac{\partial}{\partial \lambda} (\log(\lambda) - \lambda Y_i)$$

$$= \sum_{i=1}^n \left( \frac{1}{\lambda} - Y_i \right)$$

$$\boxed{\frac{\partial}{\partial \lambda} \lambda x = x}$$

$$\textcircled{3} \quad \begin{array}{l} \text{set} \\ = 0 \end{array} \quad \text{and} \quad \begin{array}{l} \text{solve for } \lambda \\ \sum_{i=1}^n \left( \frac{1}{\lambda} - Y_i \right) = 0 \end{array}$$

$$\Rightarrow \sum_{i=1}^n \frac{1}{\lambda} - \sum_{i=1}^n Y_i = 0$$

$$\frac{n}{\lambda} - \sum_{i=1}^n Y_i = 0$$

$$\Rightarrow \frac{n}{\lambda} = \sum_{i=1}^n Y_i$$

$$\Rightarrow \boxed{\hat{\lambda} = \frac{n}{\sum_{i=1}^n Y_i}}$$

### 3 Poisson regression

#### 3.1 Model choice and offsets

1. Suppose we have data on elementary schools. We are interesting in modeling the number of children from each school who participate in a special summer program, with our explanatory variable as the average reading level of students at the school. Write down the appropriate model (taking care to include an offset if you need one).

$$\text{Participants}_i \sim \text{Poisson}(\lambda_i)$$
$$\log(\lambda_i) = \beta_0 + \beta_1 \text{ReadingLevel}_i + \log(\text{SchoolSize}_i)$$

2. Suppose we have data on a random sample of senate votes from a given political year. We are interesting in modeling the number of people voting yes on a motion is related to how politically charged the topic is. We have an explanatory variable "charged" that provides a numeric measure of how politically charged (contentious) a topic is. Write down the appropriate model (taking care to include an offset if you need one).

$$\text{Votes}_i \sim \text{Poisson}(\lambda_i)$$
$$\log(\lambda_i) = \beta_0 + \beta_1 \text{charged}_i$$

If we think there are different #s of senators participating in each vote:

$$\log(\lambda_i) = \beta_0 + \beta_1 \text{charged}_i + \log(\text{Voters}_i)$$

### 3.2 Knitting

A group of knitters are attempting to determine if Brand A or Brand B of yarn breaks less often. To test this, 54 individuals from their group are randomly selected. From those 54, 27 are randomly assigned to knit using Brand A and the rest are assigned knit using Brand B. Each individual recorded how many times the yarn broke during an hour of knitting time. The results of fitting the appropriate regression model are below.

#### Model 1:

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	3.43518	0.03454	99.443	< 2e-16
woolB	-0.20599	0.05157	-3.994	6.49e-05
---				

(Dispersion parameter for poisson family taken to be 1)

Null deviance: 297.37 on 53 degrees of freedom  
 Residual deviance: 281.33 on 52 degrees of freedom

- Our question of interest is: "Is there convincing evidence of a relationship between the wool type and yarn breaks?" What kind of test would you perform to respond to this question? You do not need to perform the test.

we could use either a wald test or likelihood ratio test

- Build and interpret a 95% Wald CI for the population slope.

$$-0.20599 \pm 1.96(0.05157)$$

$$= (-0.307, -0.105)$$

$$(e^{-0.307}, e^{-0.105}) = (0.736, 0.900)$$

we are 95% confident that the log mean # of breaks is between 0.307 and 0.105 lower for brand B.

we are 95% confident that the mean # of breaks for brand B is less than the mean # of breaks for brand A by a factor of between 0.736 and 0.900

## Model 2:

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	3.69196	0.04541	81.302	< 2e-16
woolB	-0.20599	0.05157	-3.994	6.49e-05
tensionM	-0.32132	0.06027	-5.332	9.73e-08
tensionH	-0.51849	0.06396	-8.107	5.21e-16
---				
(Dispersion parameter for poisson family taken to be 1)				

Null deviance: 297.37 on 53 degrees of freedom

Residual deviance: 210.39 on 50 degrees of freedom

3. Now we are considering a new model, Model2, that uses both wool type (A or B) and tension type (Low, Medium, High) as predictors. What test could we use to determine if there was convincing evidence that Model 2 explains more variability in yarn breaks than Model 1? You do not need to perform the test yet, just state the name.

A likelihood ratio test

4. Write down the hypotheses for the test you suggested.

$$\text{Breaks}_i \sim \text{Poisson}(\lambda_i) \quad \log(\lambda_i) = \beta_0 + \beta_1 \text{WoolB}_i + \beta_2 \text{Medium}_i + \beta_3 \text{High}_i$$

$$H_0: \beta_2 = \beta_3 = 0$$

$$H_A: \text{at least one of } \beta_2, \beta_3 \neq 0$$

5. Based on the output from Model 1 and Model 2, test the hypothesis that Model 2 explains more variability in yarn breaks than Model 1. Show all of your steps. Explain how you would calculate the p-value (what is your test statistic, and what distribution would you compare with?)

$$\text{Test statistic: } G = 287.33 - 210.39$$

$$= 70.94$$

To calculate a p-value, calculate  $P(X^2 > 70.94)$

↑

2 df bc testing 2 parameters

### 3.3 Campus burglaries

We have data on 47 college campuses across the United States, and we are interested in determining what features of a university are related to the number of burglaries on campus. We have the following variables.

- `burg` = the number of burglaries on the campus in the past year.
- `campusName` = the name of the school.
- `tuition` = tuition, in thousands of dollars.
- `sat.tot` = the average total SAT score for admitted students.

**Model 1:**

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	4.134437	0.051403	80.43	< 2e-16
tuition	-0.027125	0.003799	-7.14	9.33e-13
---				

(Dispersion parameter for poisson family taken to be 1)

Null deviance: 1400.0 on 46 degrees of freedom  
 Residual deviance: 1345.9 on 45 degrees of freedom  
 AIC: 1595.4

1. Build and interpret a 95% Wald confidence interval for the slope of tuition in terms of the count.

$$-0.0271 \pm 1.96(0.0038) = (-0.0345, -0.0197)$$

$$(e^{-0.0345}, e^{-0.0197}) = (0.966, 0.980)$$

we are 95% confident that an increase in tuition by \$1000 is associated with a decrease in the average number of crimes by a factor of between 0.966 and 0.980

2. What does it mean that the dispersion parameter is “taken to be 1”?

If  $\text{crimes}_i \sim \text{Poisson}(\lambda_i)$ , then the mean # of crimes and the variance in the # of crimes are both  $\lambda_i$  (this is a feature of Poisson distributions)

$$\text{So, } \phi = \frac{\text{variance}}{\text{mean}} = 1$$

By using Poisson regression, we're implicitly assuming that the dispersion,  $\phi$ , is 1

## Model 2

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	1.5266833	0.2658280	5.743	9.29e-09
sat.tot	0.0046122	0.0004552	10.132	< 2e-16
tuition1000	-0.0275432	0.0037643	-7.317	2.54e-13
---				
(Dispersion parameter for poisson family taken to be 1)				

Null deviance: 1400.0 on 46 degrees of freedom  
Residual deviance: 1245.7 on 44 degrees of freedom  
AIC: 1497.1

3. What are the names of two possible tests we could use to compare Model 2 with Model 1?

Wald test

Likelihood ratio test

4. What is the name of, and conclusion of, the test shown below? Write the null and alternative hypothesis.

Analysis of Deviance Table

Model 1: burg09 ~ sat.tot		Model 2: burg09 ~ sat.tot + tuition1000			
		Resid.	Df	Resid.	Df
1	45	1302.8			
2	44	1245.7	1	57.146	4.047e-14

Likelihood ratio test

$$H_0: \beta_2 = 0 \quad H_A: \beta_2 \neq 0$$

We have strong evidence for a relationship between tuition & the average # of crimes, after accounting for SAT score

5. Suppose instead of Model 2, we fit a quasi-Poisson model. Our estimated dispersion parameter is  $\hat{\phi} = 27.7$ . Create a 95% Wald confidence interval for the change in the average number of burglaries associated with a one point increase in the average SAT score for admitted students, holding tuition constant.

$$\begin{aligned}
 & 0.00461 \pm 1.96(\sqrt{\hat{\phi}})(0.00046) \\
 & = 0.00461 \pm 1.96(\sqrt{27.77})(0.00046) \\
 & = (-0.00014, 0.00936) \\
 (e^{-0.00014}, e^{0.00936}) & = (0.99986, 1.0094) \\
 \text{we are 95\% confident that a one point increase in} \\
 \text{average SAT score is associated with a change in the} \\
 \text{average \# of burglaries by a factor of between} \\
 & 0.99986 \text{ and } 1.0094
 \end{aligned}$$

### 3.4 ZIP models: Brownies

Each year, a particular club sells brownies as a way of raising money for charity. This year, a new advertising campaign was used to try and increase brownie sales. To explore the effectiveness of this campaign, a survey was sent out to 300 individuals, asking how many brownies the individual purchased. Some individuals in the survey never purchase brownies, but some individuals have purchased in past years. The data is anonymous, so these distinctions are not known the individuals providing the data. Suppose you are tasked with analyzing this data. Explain why you might choose a zero inflated Poisson (ZIP) model to approach this task, and write down the model you would use and what the model parameters represent.

We are interested in the # of brownies purchased, which is a count variable, so Poisson regression may be useful. However, there are two groups of people who might record 0s here: those who never buy brownies, and those who sometimes buy brownies but didn't this year. So, a ZIP model may be useful to account for the excess 0s.

Model: Let  $y_i = \# \text{ brownies purchased}$

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

where  $\alpha_i = \text{probability an individual never purchases brownies}$ ,  
 $\lambda_i = \text{average \# of brownies purchased by individuals}$   
 $\text{who do sometimes buy brownies}$

$$\log\left(\frac{\alpha_i}{1-\alpha_i}\right) = \dots \quad (\text{we aren't told what predictors to consider})$$

$$\log(\lambda_i) = \dots \quad (\text{we aren't told what predictors to consider})$$

## 4 Case Study: Nurses

*This case study involves analyzing data and models that use several of the techniques we have learned in this course.*

Data from this study provided by Weiss (2005) includes 9573 observations on blood pressure measurements taken on nurses during a single day. In addition to physical measurements, the nurses also rate their mood on several dimensions, including how stressed they feel at the moment the blood pressure is taken. In addition, the activity of each nurse during the 10 minutes before each reading was measured using an actigraph worn on the waist. Each of the variables in is described below:

- SNUM: subject identification number
- SYS: systolic blood pressure (mmHg)
- DIA: diastolic blood pressure (mmHg)
- HRT: heart rate (beats per minute)
- MACT5: activity level (frequency of movements in 1-minute intervals, over a 10-minute period )
- DAY: workday or non-workday
- POSTURE: position during blood pressure measurement—either sitting, standing, or reclining
- STR, HAP, TIR: self-ratings by each nurse of their level of stress, happiness and tiredness at the time of each blood pressure measurement on a 5-point scale, with 5 being the strongest sensation of that feeling and 1 the weakest
- AGE: age in years
- FH123: coded as either NO (no family history of hypertension), YES (1 hypertensive parent), or YESYES (both parents hypertensive)
- time: in hours since the beginning of shift

## 4.1 Poisson regression

1. We are interested in modeling  $Y =$  the number of heart beats per minute, and choose to use  $X =$  happiness rating (HAP) as an explanatory variable. Though HAP is record in numbers from 1-5, we choose to treat it as numeric for this model. Assuming there is no over dispersion, write down the appropriate Poisson regression model.

$$Y_i \sim \text{Poisson}(\lambda_i)$$

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

### Model 1

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	4.411114	0.003408	1294.369	<2e-16
HAP	-0.009711	0.001035	-9.386	<2e-16
---				

(Dispersion parameter for poisson family taken to be 1)

Null deviance: 19130 on 8817 degrees of freedom  
 Residual deviance: 19042 on 8816 degrees of freedom  
 AIC: 73783

2. Based on Model 1, build and interpret a 95% confidence interval for average happiness score in terms of the count.

$$\begin{aligned} -0.00971 &\pm 1.96(0.001035) = (-0.0117, -0.0077) \\ (e^{-0.0117}, e^{-0.0077}) &= (0.988, 0.992) \end{aligned}$$

we are 95% confident that a unit increase in happiness score is associated with a decrease in the average number of heart beats per minute by a factor of between 0.988 and 0.992

## Model 1

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	4.411114	0.003408	1294.369	<2e-16
HAP	-0.009711	0.001035	-9.386	<2e-16
---				

(Dispersion parameter for poisson family taken to be 1)

Null deviance: 19130 on 8817 degrees of freedom  
Residual deviance: 19042 on 8816 degrees of freedom  
AIC: 73783

## Model 2

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	4.411245	0.005019	878.960	< 2e-16
HAP	-0.009754	0.001521	-6.414	1.42e-10
---				

(Dispersion parameter for Negative Binomial(68.9626) family taken to be 1)

Null deviance: 8881.3 on 8817 degrees of freedom  
Residual deviance: 8840.3 on 8816 degrees of freedom  
AIC: 70346  
Theta: 68.96  
Std. Err.: 1.94

3. What is the difference between Model 1 and Model 2?

Model 1 uses a Poisson distribution for the response  
Model 2 uses a negative binomial distribution for the response

4. Let  $\mu_i$  be the mean of the response variable. For each model (Model 1 and Model 2), what is a reasonable estimate of the standard deviation of the response variable? Hint: This will not be a number, it will involve  $\mu_i$ .

$$\text{Model 1: Standard deviation} = \sqrt{\mu_i}$$

$$\text{Model 2: Standard deviation} = \sqrt{\mu_i + \frac{\mu_i^2}{\theta}}$$

$$\text{and } \hat{\theta} = 68.963$$

5. What does overdispersion mean? Explain in 1-2 sentences.

Overdispersion means that there is more variability in our response than is assumed by our model

## 4.2 Zero inflated Poisson (ZIP)

Now we are modeling  $Y$  = the amount of coffee, in cups, that a nurse consumes on a given day. During the study, some of the coffee machines on the 3rd floor of the hospital were not working, meaning that some of the nurses were not able to get coffee when they worked on the third floor, even though they usually drink coffee. We do not have information on which floor the nurses were working on during the study.

We now fit a zero inflated Poisson (ZIP) model, and get the following fitted model:

### Model 4

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

where  $\alpha_i$  is the probability a nurse was not able to get coffee, and  $\lambda_i$  is the average number of cups consumed by a nurse able to get coffee. Our estimates are

$$\log\left(\frac{\hat{\alpha}_i}{1 - \hat{\alpha}_i}\right) = 0.40 + 0.20 \text{ DayNW}_i$$

$$\log(\hat{\lambda}_i) = 0.65 + 0.141 \text{ Time}_i$$

What is the probability that a nurse who is 7 hours into their shift, on a work day, drinks 3 cups of coffee?

$$\hat{\alpha}_i = \frac{e^{0.4}}{1 + e^{0.4}} = 0.6$$

$$\hat{\lambda}_i = e^{0.65 + 0.141(7)} = 5.14$$

$$\hat{P}(Y_i = 3) = \frac{e^{-5.14} 5.14^3}{3!} (1 - 0.6)$$

$$= 0.053$$

### 4.3 Linear mixed effect models

Now we are provided new information about our data. The data include 9573 rows, but this is made up of observations taken on only a random sample of 203 nurses over the course of a single day. This means 40-60 measurements were taken per nurse. The first blood pressure measurement was taken half an hour before the nurse's normal start of work, and was measured approximately every 20 minutes for the rest of the day.

1. We are now interested in modeling  $Y$  = the systolic blood pressure (SYS), using posture as our explanatory variable. Write an appropriate model.

$y_{ij}$  = systolic blood pressure for nurse  $i$  at measurement  $j$

$$y_{ij} = \beta_0 + \beta_1 \text{Sit}_{ij} + \beta_2 \text{Stand}_{ij} + u_i + \varepsilon_{ij}$$

↑  
random effect for  
nurse  $i$

$$u_i \stackrel{\text{iid}}{\sim} N(0, \sigma_u^2) \quad \varepsilon_{ij} \stackrel{\text{iid}}{\sim} N(0, \sigma_\varepsilon^2)$$

iid: independent & identically distributed

## Model 5

Random effects:

Groups	Name	Variance	Std.Dev.
SNUM	(Intercept)	70.54	8.399
Residual		166.22	12.893
Number of obs:	9573, groups:	ID, 203	

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	109.9183	0.7987	137.62
POSTURESIT	7.9044	0.5746	13.76
POSTURESTAND	9.8293	0.5806	16.93

Output

SNUM	Intercept	POSTURESIT	POSTURESTAND
1006	108.5899	7.904368	9.82931

2. For nurse 1006, what is the predicted average systolic blood pressure while standing?

Estimated blood pressure is  $\hat{\beta}_0 + u_i + \hat{\beta}_1 S_{i,j} + \hat{\beta}_2 S_{i,j}$   
 $108.59 = \hat{\beta}_0 + \hat{u}_{1006}$

For nurse 1006, standing:

$$\text{est. blood pressure} = 108.59 + 9.83 \\ = 118.42$$

3. Interpret the estimated random effect for nurse 1006.

$$\hat{u}_{1006} = 108.59 - 109.92 \\ = -1.33$$

On average, blood pressure for Nurse 1006 is about 1.33 points lower than the mean blood pressure for all nurses, holding posture fixed

4. Using the output from Model 5, does there appear to be systematic variation in systolic blood pressure between nurses? Calculate an appropriate statistic.

$$\hat{P}_{\text{group}} = \frac{\hat{\sigma}_u^2}{\hat{\sigma}_u^2 + \hat{\sigma}_e^2} = \frac{70.54}{70.54 + 166.62} = 0.297$$

between-group variance  
total within-group variance

$\approx 30\%$  of variance in blood pressure can be explained by systematic differences between nurses, after accounting for posture

5. Suppose you want to test whether there is systematic variation in systolic blood pressure between nurses. Write down the null and alternative hypotheses, and describe what your reduced model would be (you may treat Model 5 as your full model).

$$H_0: \sigma_u^2 = 0$$

$$H_A: \sigma_u^2 > 0$$

accounting  
for posture

reduced model

$$y_{ij} = \beta_0 + \beta_1 S_{itij} + \beta_2 S_{standij} + \varepsilon_{ij}$$

$$\varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma_e^2)$$

You fit your reduced model from the previous question, producing the following output:

Coefficients:

	Estimate	Std. Error	
(Intercept)	112.3121	0.8253	
POSTURESIT	6.8760	0.6412	
POSTURESTAND	8.5402	0.4937	

Fitted reduced model

Residual standard error: 14.175

6. Describe how you would use parametric bootstrapping to carry out the hypothesis test from the previous question. Provide as much detail as you can, so that someone could turn your description into R code if they wanted to (you do not need to write code, though you may choose to if it helps you explain your procedure). Your description should include details like values for the parameters of the model you will simulate from, how many simulations you will use, how you will calculate a test statistic for each simulation, and how you will calculate a p-value from your bootstrap results at the end.

- 1) Calculate a test statistic on observed data e.g.  $\hat{P}_{group} = 0.3$
- 2) Simulate data from reduced model  
 $\varepsilon_{ij}^* \sim N(0, 14.175^2)$   
 $y_{ij}^* = \text{fitted value from reduced model} + \varepsilon_{ij}^*$   
 $= 112.31 + 6.88 S_{itij} + 8.54 S_{standij} + \varepsilon_{ij}^*$
- 3) Calculate a test statistic on simulated data  
- e.g. fit full model on simulated data,  
calculate  $\hat{P}_{group}^*$
- 4) Repeat many times, store the resulting test statistics
- 5) Calculate bootstrap p-value:  

$$\frac{\#\{\hat{P}_{group}^* > 0.3\}}{\#\text{simulations}}$$