

case1final

Nathaniel Brown, Annie Tang, William Yang

September 13, 2017

Introduction

UNFINISHED!!!

The purpose of this project is to analyze rat bioassays. is there variation in results between labs? there shouldn't be.

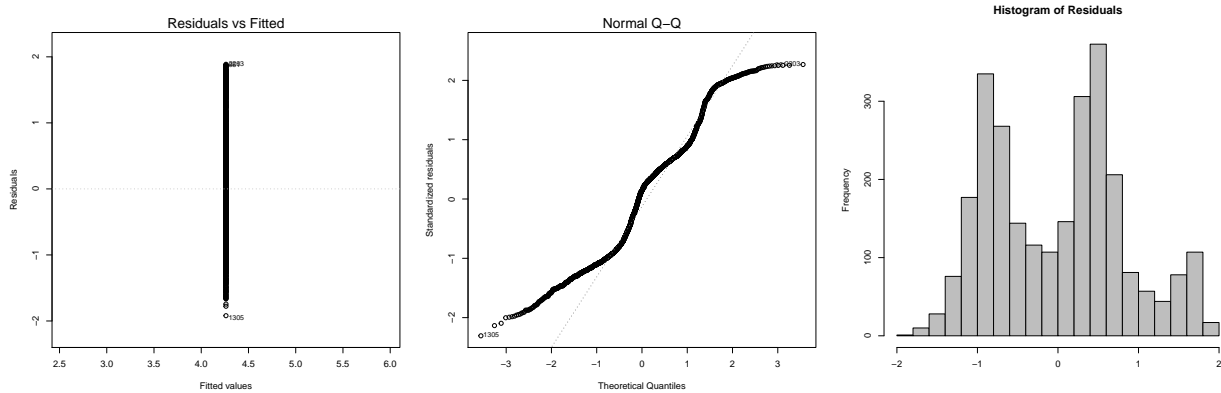
Model-Fitting

UNFINISHED!!!

We fit models to try to understand the variation between labs. Our approaches were univariate normal, multivariate normal, and mixed effects.

First Approach: Univariate Normal

$$\log(y) \sim \text{Norm}(\mu, \sigma^2)$$



This naive approach assumes the blotted weight follows a normal distribution with a mean centered around the mean of the log blotted weight (4.2619) and constant variance. The diagnostic plots below show that the same value is predicted for each input, the quantiles of the residuals do not follow a normal distribution, and that the residuals are bimodal. We believe bimodality may be caused by the separation of rats into juvenile and adult protocols.

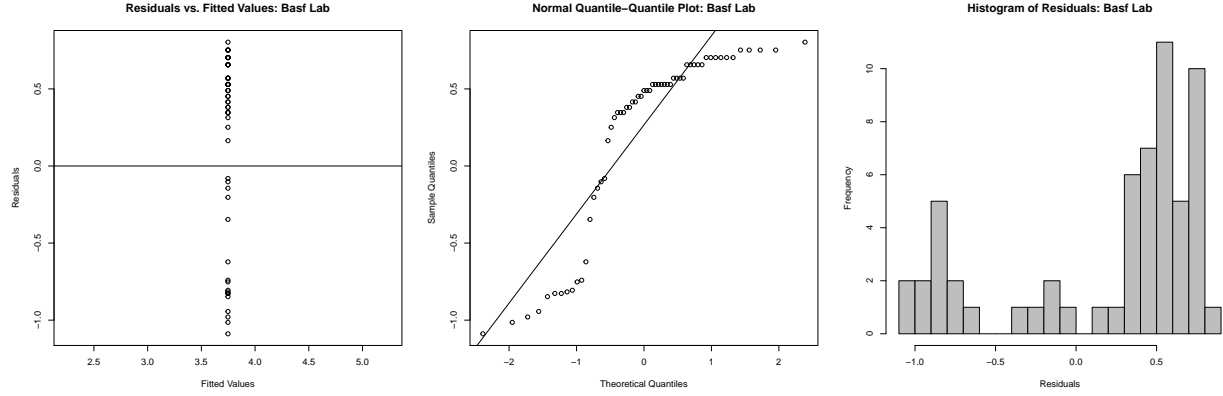
Second Approach: Multivariate Normal

$$y_i, \mu \in M_{n_i \times 1}(\mathbb{R}) \Sigma \in M_{n \times n}(\mathbb{R}) \log(y_i) \sim \text{MVNorm}(\mu, \Sigma)$$

In this model, y_i is a vector of n_i observations from lab i . It assumes that the log blotted weights of each lab follow approximately normal distributions with a their own means.

We provide the diagnosis plots for the first lab, **Basf**, and we put the others in Appendix 1.1, since R does not support plotting for multivariate regression models.

Warning: package 'bindrcpp' was built under R version 3.3.3



The diagnostic plots for this model indicate that the model predicts the same fitted value for each lab, and that the residual quantiles do not follow a normal distribution. For the **Basf** lab, the histogram of residuals shows left skew. These same problems, along with multimodality of the residuals, are present for the diagnostic plots for the remaining labs in the Appendix.

A weakness of this model is that, as you add more structure (blocking factors, covariates for dosage, etc.) to the model, the covariance matrix becomes more complicated, and maximum likelihood estimation becomes unwieldy. The poor fit of this model indicates that we should control for more than just the lab effect, so we consider a mixed effects model.

Third Approach: Mixed Effects

$$y_{ij} \sim \beta_{0,i} + \beta_{1,i}x_{ij,d_1} + \beta_{2,i}x_{ij,d_2} + \beta_{3,i}x_{ij,p_B} + \beta_{4,i}x_{ij,p_C} + \beta_{5,i}x_{ij,p_D} + \beta_{6,i}x_{ij,\log(w)} + \epsilon \quad \beta_{0:6,i} \sim N(\mu_{0:6,i}, \sigma_{0:6,i}^2) \quad \epsilon \sim N(0, \sigma^2)$$

Let y_{ij} be the observation for log(blotted uterus weight) for subject x_{ij} , the j th individual in lab i . x_{ij,d_1} and x_{ij,d_2} are the values of dose1 and dose2 for subject x_{ij} . x_{ij,p_B} , x_{ij,p_C} , and x_{ij,p_D} are dummy variables for which protocol x_{ij} was subjected to. $x_{ij,\log(w)}$ is the log(body weight) for x_{ij} . We make the Gaussian assumption that the coefficients, β , are normally distributed according to some μ_i and σ_i . We add a random effect on all $\beta_{0:5,i}$ to account for lab-to-lab variability in the intercepts and slopes on the blotted weight for the different dosages and protocols.

We start at a reduced form of this model and augment it to its full form after initial analyses.

In our analysis, we choose to add a random effect for the lab variable, in order to account for lab-to-lab heterogeneity. Essentially, this allows us to avoid violating an independence assumption by assuming a different baseline response for each lab. First, we model these differences between individual labs by assuming random intercepts for each lab. We first include only dose1 and dose2 as fixed effects, then introduce proto as a fixed effect after confirming through anova ($p < 2.2e-16$, $\chi^2=1675.8776$) that it is a significant predictor. The summary statistics of this model are shown in the table below:

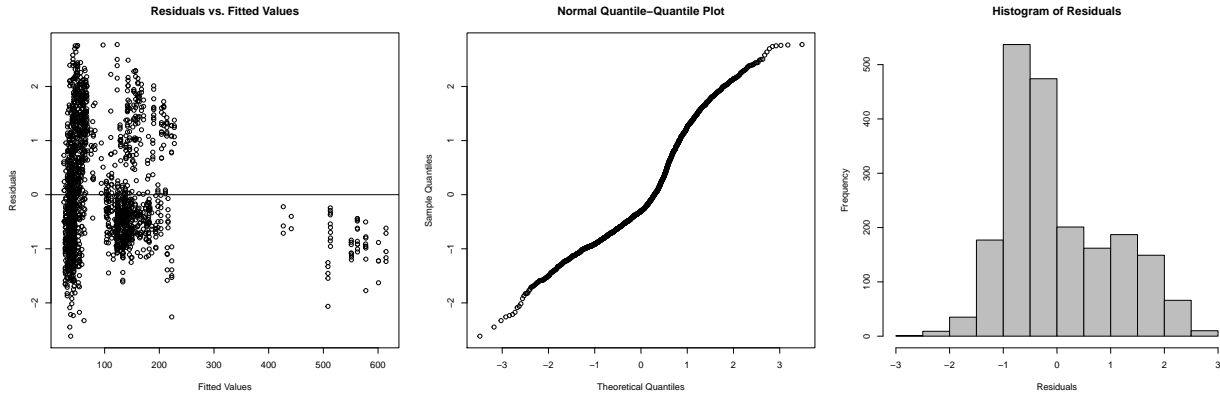
	Mean Fixed Effects	Variance Fixed Effects
(Intercept)	3.6311	0.0018

	Mean Fixed Effects	Variance Fixed Effects
dose1	0.1420	0.0000
dose2	-0.5195	0.0012
protoB	0.0379	0.0008
protoC	1.2579	0.0009
protoD	1.2574	0.0016

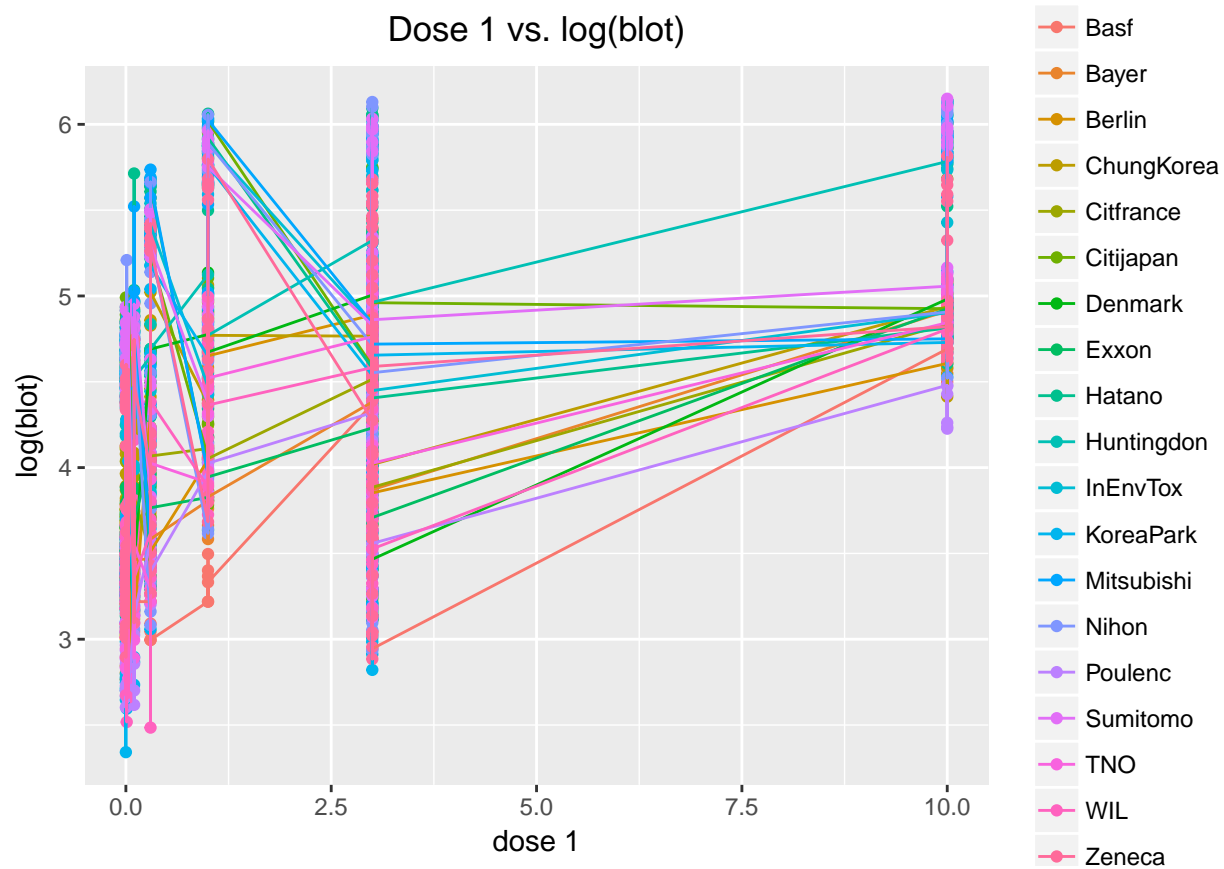
	Mean Random Effects	Variance Random Effects
(Intercept)	0.3834	0.0284

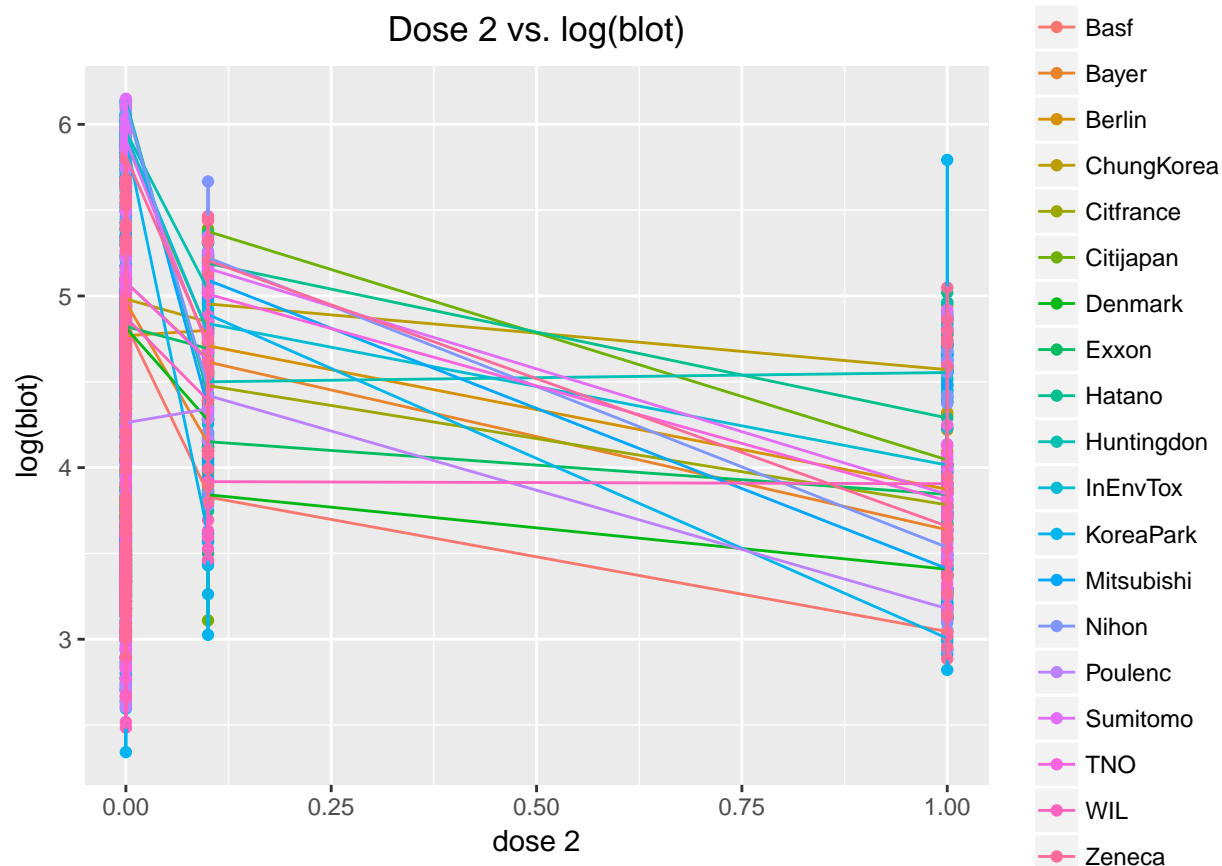
	Variance
Residual	0.1933

With this full model, we see that the variability due to lab is 0.028 (in terms of variance), and variability due to non-lab sources is 0.193. And when we take a look at the fixed effects, we see that an increase in the amount of dose1 corresponds to an increase (0.142 units) in the response variable, $\log(\text{blot})$. Furthermore, an increase in the amount of dose2 corresponds to a decrease (-0.519 units) in the response variable. Protocols B, C, and D all lead to an increase in the response variable, relative to protocol A. The diagnostic plots below show that this model fits better than the previous two. Although the residual variance decreases as the fitted values increase, the quantiles of the residuals are approximately normal.



In this random intercept model, we account for baseline differences between labs, but we also assume that the effects of doses is the same for each lab. After plotting the $\log(\text{blot})$ versus dose by lab, we see that this is not a valid assumption to make since the lines are clearly not parallel. So we introduce a random slope model. Now, dose1 and dose2 can have varying slopes.



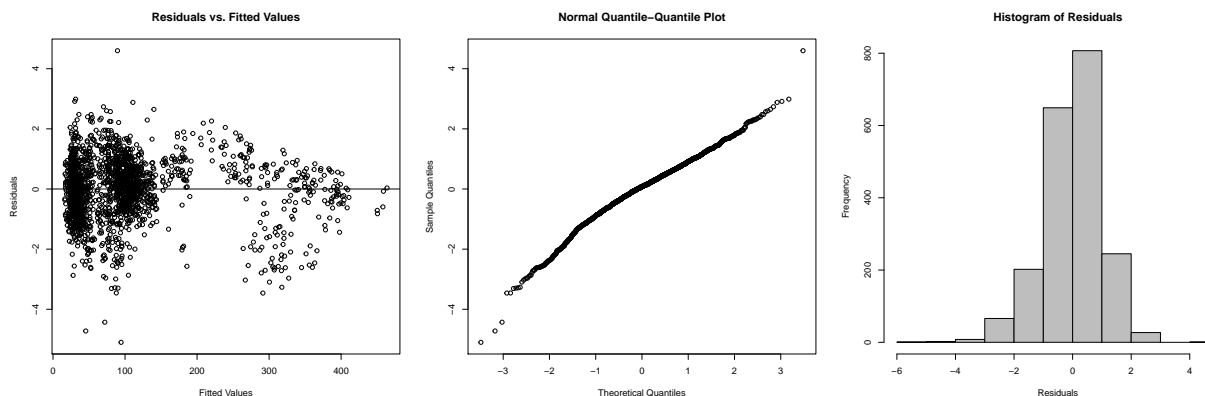


After adding the random slopes, we transformed dose 1 using a reciprocal transformation and added body weight as a predictor in order to have the best fitting model. We evaluate this model using summary statistics, diagnostic plots, and out-of-sample predictive accuracy, which we evaluate using Mean Absolute Error ($MAE = E[|y - \hat{y}|]$) and Root Mean Squared Error ($RMSE = \sqrt{E[(y - \hat{y})^2]}$). Root Mean Squared Error penalizes more for extreme errors, while Mean Absolute Error simply averages all of the errors.

	Mean Fixed Effects	Variance Fixed Effects
(Intercept)	2.6379	0.0728
protoB	0.0457	0.0003
protoC	0.8553	0.0094
protoD	0.8524	0.0100
log(body)	0.2946	0.0046

	Mean Random Effects	Variance Random Effects
(Intercept)	3.4691	0.9593
$I(1/(dose1 + 1/2))$	0.4969	0.5371
dose2	0.7532	1.0481

	Variance
Residual	0.0797



	MAE	RMSE
Random Dose + Transformations	21.7968	37.7587
Fixed Dose	35.2341	57.9896

We see that the variability due to sources other than the random effects is 0.0797 (decreased from 0.193). Also, the diagnostic plots improved, although it is notable that the residual variance decreases for larger fitted values. The fixed dose model is also much worse at predicting than the random dose model with transformations.

Lab Variation

INCOMPLETE

The plot below randomly samples dose1 effects and a lab effects (or intercepts), and plots the results. Since there is so much variance in the dose1 effects by lab, we say that the bioassay does depend on lab and thus this study fails miserably.

Conclusion

Because the dose 1 effect has such a high variance relative to the mean, our model shows that the bioassay being studied does not measure a consistent response to dosage across labs.

Contributions

Nathaniel Brown made the visualizations for this report. He also organized the relevant files in a Github repository for the group to access and edit. Annie Tang compiled the group work done on EDA into a .rmd and wrote the accompanying explanations for the EDA and approaches to analysis. William Yang helped pair on EDA analysis and identify approaches to handle the data. Approaches to analysis were a joint effort by all members of the group. Nathaniel implemented analysis for the univariate normal and multivariate normal approaches. Implementation and analysis of the mixed effects model was a joint effort by all members of the group.

Appendix

1.1: Multivariate Normal Diagnosis Plots

