

# Handling Categorical Predictors: ANOVA

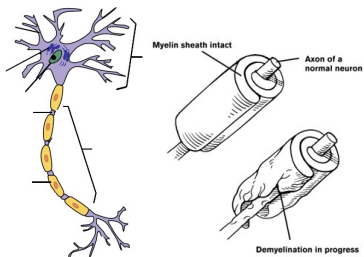
1/33

## I Hate Lines!

- ▶ When we think of experiments, we think of manipulating categories
- ▶ Control, Treatment 1, Treatment 2
- ▶ Models with Categorical Predictors still reflect an underlying data generating process and error generating process
- ▶ In many ways, it's like having many processed generating data
- ▶ Linear (or Generalized Linear) Models still underlie ANOVA or other analyses with categorical predictors

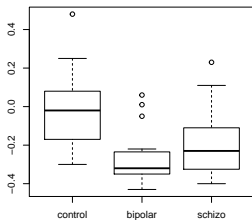
2/33

## Categorical Predictors: Gene Expression and Mental Disorders



3/33

## Categorical Predictors

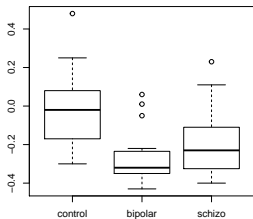


How do we determine the importance of categorical predictors?

4/33

## Aside: Reordering Factors

```
brainGene$group <- factor(brainGene$group,  
                           levels=c("control", "bipolar", "schizo"))
```



How do we determine the importance of categorical predictors?

5/33

## Categorical Predictors Ubiquitous

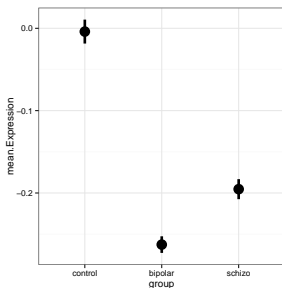
- ▶ Treatments in an Experiment
- ▶ Spatial groups - plots, Sites, States, etc.
- ▶ Individual sampling units
- ▶ Temporal groups - years, seasons, months

6/33

# Modeling Categorical Predictors

7/33

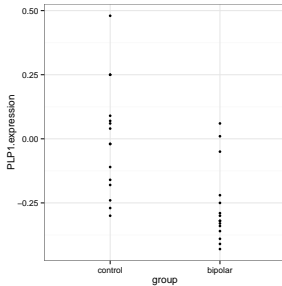
## Traditional Way to Think About Categories



What is the variance between groups v. within groups?

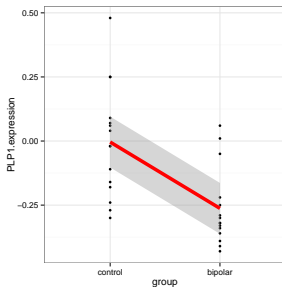
8/33

## But How is the Model Fit?



9/33

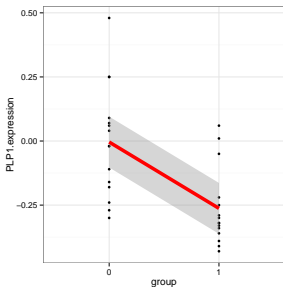
## But How is the Model Fit?



Underlying linear model with control = intercept, dummy variable for bipolar

10/33

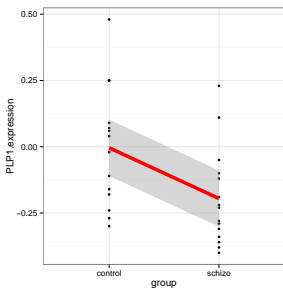
## But How is the Model Fit?



Underlying linear model with control = intercept, dummy variable for bipolar

11/33

## But How is the Model Fit?



Underlying linear model with control = intercept, dummy variable for schizo

12/33

## Different Ways to Write a Categorical Model

$$y_{ij} = \bar{y} + (\bar{y}_i - \bar{y}) + (y_{ij} - \bar{y}_i)$$

$$y_{ij} = \mu + \alpha_i + \epsilon_{ij}, \quad \epsilon_{ij} \sim N(0, \sigma^2)$$

$$y_j = \beta_0 + \sum \beta_i x_i + \epsilon_j, \quad x_i = 0, 1$$

$x_i$  indicates presence/absence of a category

Traditional ANOVA special case where all  $x_i$  are orthogonal

Often one category set to  $\beta_0$  for ease of fitting

13/33

## This is a Linear Model

```
bg.sub.lm <- lm(PLP1.expression ~ group, data=brainGene)
```

14/33

# Evaluating Categorical Predictors

15/33

## Hypothesis Testing with a Categorical Model: ANOVA

$$H_0 = \mu_1 = \mu_2 = \mu_3 = \dots$$

OR

$$\beta_0 = \mu, \quad \beta_i = 0$$

16/33



## F-Test to Compare

Regression:  $SS_{Total} = SS_{Model} + SS_{Error}$

$SS_{Total} = SS_{Between} + SS_{Within}$

$$SS_{Between} = \sum_i \sum_j (\bar{Y}_i - \bar{Y})^2, \text{ df} = k - 1$$

$$SS_{Within} = \sum_i \sum_j (Y_{ij} - \bar{Y}_i)^2, \text{ df} = n - k$$

To compare them, we need to correct for different DF. This is the Mean Square.

MS = SS/DF, e.g.  $MS_W = \frac{SS_W}{n - k}$

17/33

## F-Test to Compare

$$F = \frac{MS_B}{MS_W} \text{ with DF} = k - 1, n - k$$

(note similarities to  $SS_R$  and  $SS_E$  notation of regression)

18/33

## ANOVA

```
anova(bg.sub.lm)
```

```
# Analysis of Variance Table
#
# Response: PLP1.expression
#      Df Sum Sq Mean Sq F value Pr(>F)
# group   2  0.54025  0.270127   7.8231 0.001294
# Residuals 42  1.45023  0.034529
```

19/33

Is using ANOVA valid?

20/33

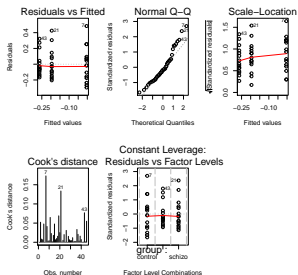
# Assumptions of Ordinary Least Squares Regression

- ▶ Independence of data points
- ▶ Normality within groups
- ▶ Homoscedasticity (homogeneity of variance)

21/33

## Inspecting Assumptions

```
par(mfrow=c(2,3))  
plot(bg.sub.lm, which=1:5 )  
par(mfrow=c(1,1))
```



22/33

## Levene's Test of Homogeneity of Variance

```
library(car)
leveneTest(PLP1.expression ~ group, data=brainGene)

# Levene's Test for Homogeneity of Variance (center = median)
#      Df F value Pr(>F)
# group  2  1.0067 0.3741
#      42
```

Levene's test robust to departures from normality

23/33

## What do I do if I Violate Assumptions?

- ▶ Nonparametric Kruskal-Wallis (uses ranks)
- ▶ Transform?
- ▶ GLM with ANODEV

24/33

## Kruskal Wallace Test

```
kruskal.test(PLP1.expression ~ group, data=brainGene)

#
# Kruskal-Wallis rank sum test
#
# data: PLP1.expression by group
# Kruskal-Wallis chi-squared = 13.1985, df = 2, p-value
# = 0.001361
```

25/33

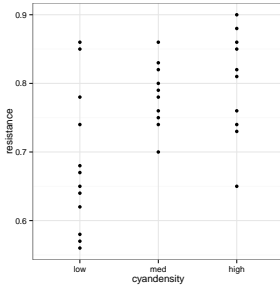
## Exercise: Daphnia Resistance

- ▶ Evaluate whether the data is appropriate for ANOVA
- ▶ Fit an ANOVA and check diagnostics
- ▶ Evaluate results & compare to Kruskal-Wallis and a glm with a Gamma distribution



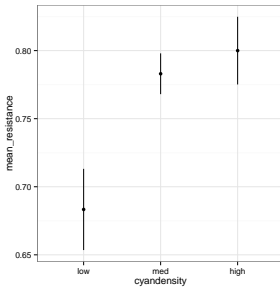
26/33

## Daphnia Data



27/33

## Daphnia Means



28/33

## How about HOV?

```
leveneTest(resistance ~ cyandensity, data=daphnia)

# Levene's Test for Homogeneity of Variance (center = median)
#      Df F value Pr(>F)
# group 2  2.0019 0.1533
#      29
```

29/33

## ANOVA shows an Effect

```
daphniaLM <- lm(resistance ~ cyandensity, data=daphnia)
anova(daphniaLM)

# Analysis of Variance Table
#
# Response: resistance
#           Df  Sum Sq Mean Sq F value  Pr(>F)
# cyandensity  2 0.089195  0.044598   6.6916 0.004078
# Residuals   29 0.193277  0.006665
```

30/33

## KW shows an Effect

```
#  
# Kruskal-Wallis rank sum test  
#  
# data: resistance by cyandensity  
# Kruskal-Wallis chi-squared = 8.1996, df = 2, p-value  
# = 0.01658
```

31/33

## Bad GLM Does Not

```
# Analysis of Deviance Table  
#  
# Model: Gamma, link: identity  
#  
# Response: resistance  
#  
# Terms added sequentially (first to last)  
#  
#  
#           Df Deviance Resid. Df Resid. Dev  
# NULL                31    0.52908  
# cyandensity  2  0.16216          29    0.36692
```

32/33



## Diagnostics Also Good

