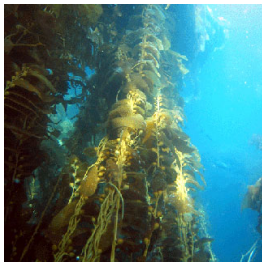


After the ANOVA

1/50

Group Properties: Kelp

- ▶ Kelp sampled at multiple sites annually
- ▶ At each transect, holdfast diameter and # of fronds counted



2/50

ANOVA

- ▶ Comparison of grouped means
- ▶ Really a special case of a linear model
- ▶ F-Tests to evaluate if a normal error generating process
- ▶ Can use ANODEV & LR Ratios for non-normal error generating process

3/50

Today

- ▶ plyr and visualizing things by groups
- ▶ Evaluating treatment means after ANOVA
- ▶ Contrasting treatment means
- ▶ Unplanned post-hoc testing

4/50

How can we get quick summaries by site?, year, or both?

#	YEAR	MONTH	DATE	SITE	TRANSECT	QUAD	SIDE	FRONDS
# 2	2000	9	2000-09-28	BULL	1	20		4
# 8	2000	9	2000-09-28	BULL	2	20		11
# 9	2000	9	2000-09-28	BULL	2	20		16
# 10	2000	9	2000-09-28	BULL	2	20		34
# 16	2000	9	2000-09-28	BULL	3	20		27
# 17	2000	9	2000-09-28	BULL	3	20		38
#	HLD_DIAM							
# 2	7							
# 8	65							
# 9	55							
# 10	55							
# 16	65							
# 17	60							

5/50

For loops for Summarization by Site

```
# number of groups
k <- length(levels(kelp$SITE))

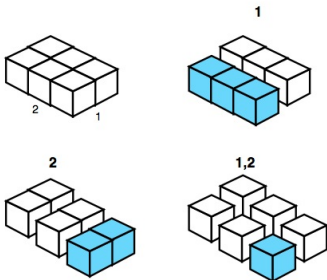
#blank means vector
means <- rep(NA, k)

#the loop
for(i in 1:k) {
  #split the data first
  subdata <- subset(kelp, kelp$SITE == levels(kelp$SITE)[i])

  #apply the means function,
  #combine with previous means
  means[i] <- mean(subdata$FRONDS, na.rm=T)
}
```

6/50

The Split, Apply, Combine Strategy



Wickham 2011

7/50

ddply from Hadley Wickham's plyr library

```
library(plyr)
#
kelpMeans <- ddply(kelp, .(SITE), summarise,
  mean.FRONDS = mean(FRONDS, na.rm=T))
```

8/50

ddply from Hadley Wickham's plyr library

```
kelpMeans
```

```
#   SITE mean.FRONDS
# 1 ABUR  28.810000
# 2 AHND  17.633508
# 3 AQUE  21.029720
# 4 BULL  27.272152
# 5 CARP  13.110985
# 6 GOLB  42.164319
# 7 IVEE  25.777251
# 8 MOHK  20.041916
# 9 NAPL  13.159147
#10 SCDI   1.058824
#11 SCTW  14.492063
```

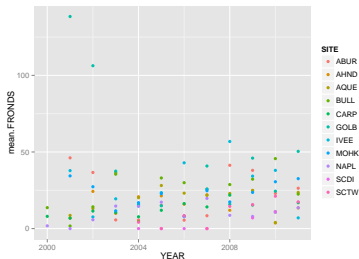
9/50

Multiple Groups & ddply

```
kelpMeans2 <- ddply(kelp, .(YEAR, SITE), summarise,  
                    mean.FRONDS = mean(FRONDS, na.rm=T))
```

10/50

Multiple Groups & ddply



11/50

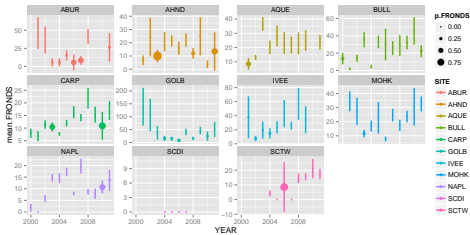
Complex Functions & ddply

```
kelpMeans3 <- ddply(kelp, .(YEAR, SITE), function(aFrame){
  #calculate metrics for a 1-sample T test comparison against
  #grand mean of 10 fronds/m^2
  m <- mean(aFrame$FRONDS, na.rm=T)
  n<-length(na.omit(aFrame$FRONDS))
  se <- sd(aFrame$FRONDS, na.rm=T)/sqrt(n)
  t <- (m-10)/se
  p <- 2*pt(abs(t), df=n-1, lower.tail=F)

  # return everything
  return(c(mean.FRONDS=m, n.FRONDS=n,
           se.FRONDS=se, t.FRONDS=t,
           p.FRONDS = p))
})
```

12/50

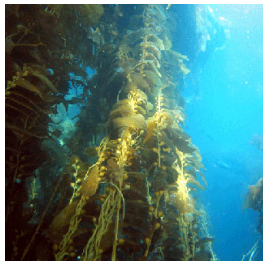
Complex Functions & dply



13/50

Exercise: Correlation!

- ▶ Evaluate the correlation between fronds and holdfasts by site and year
- ▶ Plot it
- ▶ Extra: include the SE of the correlation visually



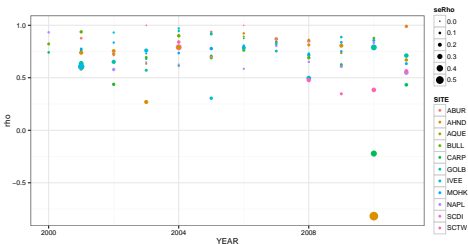
14/50

Exercise: Correlation!

```
kelpCor <- ddply(kelp, .(YEAR, SITE), function(adf){  
  #first get the correlation  
  cors <- cor(adf$FROND, adf$HLD_DIAM)  
  
  #use this to calculate it's SE  
  seCor <- sqrt((1-cors^2) / (nrow(adf)-2))  
  
  #return both  
  return(c(rho = cors, seRho = seCor))  
})
```

15/50

Exercise: Correlation!



16/50

Many plyr Functions

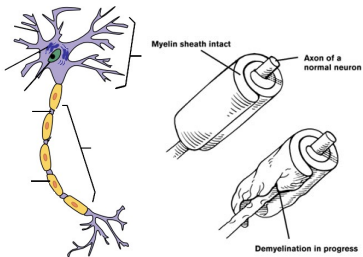
<i>Input</i> \ <i>Output</i>	Array	Data frame	List	Discarded
Array	aapply	adply	alply	a_ply
Data frame	dapply	ddply	dlply	d_ply
List	lapply	ldply	llply	l_ply

Also `r*ply` to replicate an action and return an object. Great for simulation.

See also `colwise` and `each` for everyday use!

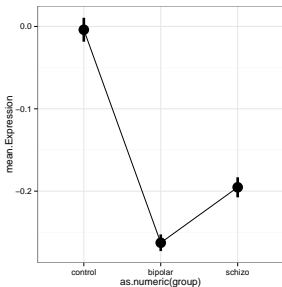
Looking at Groups After an ANOVA

Categorical Predictors: Gene Expression and Mental Disorders



19/50

The Data



20/50

Fit the Data with a Linear Model

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

21/50

F-Test to Compare Variation Within versus Between Groups

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

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

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

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

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

22/50

ANOVA

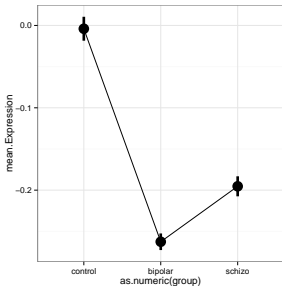
```
anova(bg.sub.lm)

# 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
```

Which groups are different from one another?

23/50

The Data



24/50

How would you have made that graph?

25/50

The Coefficients

```
summary(bg.sub.lm)

#
# Call:
# lm(formula = PLP1.expression ~ group, data = brainGene)
#
# Residuals:
#      Min       1Q   Median       3Q      Max
# -0.29600 -0.12733 -0.03467  0.07533  0.48400
#
# Coefficients:
#              Estimate Std. Error t value Pr(>|t|)
# (Intercept) -0.00400    0.04798  -0.083  0.933953
# groupbipolar -0.25867    0.06785  -3.812  0.000444
# groupschizo  -0.19133    0.06785  -2.820  0.007301
#
# Residual standard error: 0.1858 on 42 degrees of freedom
# Multiple R-squared:  0.2714, Adjusted R-squared:  0.2367
# F-statistic: 7.823 on 2 and 42 DF,  p-value: 0.001294
```

26/50

Default "Treatment" Contrasts

```
contrasts(brainGene$group)

#          bipolar schizo
# control      0      0
# bipolar      1      0
# schizo       0      1
```

27/50

The Coefficients

```
summary(lm(PLP1.expression ~ group -1, data=brainGene))

#
# Call:
# lm(formula = PLP1.expression ~ group - 1, data = brainGene)
#
# Residuals:
#      Min       1Q   Median       3Q      Max
# -0.29600 -0.12733 -0.03467  0.07533  0.48400
#
# Coefficients:
#              Estimate Std. Error t value Pr(>|t|)
# groupcontrol  -0.00400    0.04798  -0.083  0.933953
# groupbipolar  -0.26267    0.04798  -5.475  2.25e-06
# groupschizo  -0.19533    0.04798  -4.071  0.000202
#
# Residual standard error: 0.1858 on 42 degrees of freedom
# Multiple R-squared:  0.5257, Adjusted R-squared:  0.4918
# F-statistic: 15.52 on 3 and 42 DF,  p-value: 6.125e-07
```

28/50

OK, but WHICH GROUPS ARE DIFFERENT?

29/50

ANOVA is an Omnibus Test

Remember your Null:

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

This had nothing to do with specific comparisons of means.

30/50

A priori contrasts

Specific sets of a priori null hypotheses:

$$\mu_1 = \mu_2$$

$$\mu_1 = \mu_3 = \dots$$

Use t-tests.

31/50

A priori contrasts

```
library(contrast)

contrast(bg.sub.lm, list(group="control"),
         list(group="schizo"))

# lm model parameter contrast
#
# Contrast      S.E.      Lower      Upper      t df Pr(>|t|)
# 1 0.1913333 0.067852 0.05440245 0.3282642 2.82 42 0.0073
```

32/50

A priori contrasts

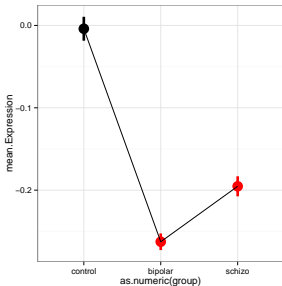
```
contrast(bg.sub.lm, list(group="control"),  
         list(group=c("schizo", "bipolar")))
```

```
# lm model parameter contrast  
#  
# Contrast      S.E.      Lower      Upper      t df Pr(>|t|)  
# 0.1913333 0.067852 0.05440245 0.3282642 2.82 42 0.0073  
# 0.2586667 0.067852 0.12173578 0.3955976 3.81 42 0.0004
```

Note: can only do k-1, as each takes 1df

33/50

The Data



34/50

Orthogonal A priori contrasts

Sometimes you want to test very specific hypotheses about the structure of your groups

```
#           control bipolar schizo
# Control v. Disorders      1   -0.5  -0.5
# Bipolar v. Schizo         0    1.0  -1.0
```

Note: can only do k-1, as each takes 1df

35/50

Orthogonal A priori contrasts with multcomp

```
library(multcomp)
#
bg_orthogonal <- glht(bg.sub.lm, linfct=contrast_mat,
                     test=adjusted("none"))
#
summary(bg_orthogonal)
```

Note adjusted p-value is set to none...

36/50

Orthogonal A priori contrasts

```
#  
# Simultaneous Tests for General Linear Hypotheses  
#  
# Fit: lm(formula = PLP1.expression ~ group, data = brainGene)  
#  
# Linear Hypotheses:  
#  
#           Estimate Std. Error t value  
# Control v. Disorders == 0  0.22100   0.10178   2.171  
# Bipolar v. Schizo == 0  -0.06733   0.06785  -0.992  
#  
#           Pr(>|t|)  
# Control v. Disorders == 0  0.0695  
# Bipolar v. Schizo == 0    0.5439  
# (Adjusted p values reported -- single-step method)
```

37/50

Post hoc contrasts

I want to test all possible comparisons!

38/50

Post hoc contrasts

Only to be done if you reject H_0

- ▶ All possible comparisons via t-test
- ▶ But...with many comparisons, does type I error rate increase?
- ▶ Consider adjusted alpha
- ▶ But, adjusting alpha also may increase type II error rate!
- ▶ Additional multiple comparison methods calculate family-wise critical values of differences.

39/50

All Possible T-Tests

```
with( brainGene, pairwise.t.test(PLP1.expression, group,
                                p.adjust.method = "none" )

#
# Pairwise comparisons using t tests with pooled SD
#
# data:  PLP1.expression and group
#
#          control bipolar
# bipolar 0.00044 -
# schizo  0.00730 0.32671
#
# P value adjustment method: none
```

40/50

P-Value Adjustments

Bonferroni : $\alpha_{adj} = \frac{\alpha}{m}$ where $m = \#$ of tests

- VERY conservative

False Discovery Rate: $\alpha_{adj} = \frac{k\alpha}{m}$

- Order your p values from smallest to largest, rank = k,

- Adjusts for small v. large p values

- Less conservative

Other Methods: Sidak, Dunn, Holm, etc.

We're very focused on p here!

41/50

Bonferroni Correction

```
with( brainGene, pairwise.t.test(PLP1.expression, group,
                                p.adjust.method = "bonferroni") )

#
# Pairwise comparisons using t tests with pooled SD
#
# data:  PLP1.expression and group
#
#          control bipolar
# bipolar 0.0013  -
# schizo  0.0219  0.9801
#
# P value adjustment method: bonferroni
```

42/50

False Discovery Rate

```
with( brainGene, pairwise.t.test(PLP1.expression, group,
                                p.adjust.method = "fdr" )

#
# Pairwise comparisons using t tests with pooled SD
#
# data:  PLP1.expression and group
#
#          control bipolar
# bipolar 0.0013  -
# schizo  0.0110  0.3267
#
# P value adjustment method: fdr
```

43/50

Other Methods Use Critical Values

- ▶ Tukey's Honestly Significant Difference
- ▶ Dunnet's Test for Comparison to Controls
- ▶ Ryan's Q (sliding range)
- ▶ etc...

44/50

Tukey Test

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

# Tukey multiple comparisons of means
# 95% family-wise confidence level
#
# Fit: aov(formula = PLP1.expression ~ group, data = brainGene)
#
# $group
#           diff           lwr           upr
# bipolar-control -0.25866667 -0.42351268 -0.09382065
# schizo-control  -0.19133333 -0.35617935 -0.02648732
# schizo-bipolar   0.06733333 -0.09751268  0.23217935
#
#           p adj
# bipolar-control 0.0012670
# schizo-control  0.0195775
# schizo-bipolar  0.5857148
```

45/50

Final Notes of Caution

- ▶ Often you DO have a priori contrasts in mind
- ▶ If you reject H_0 with ANOVA, differences between groups exist
- ▶ Consider Type I v. Type II error before correcting

46/50

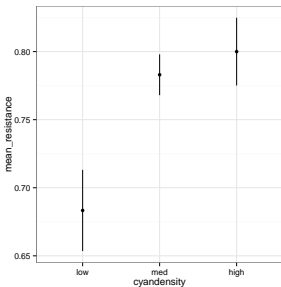
Exercise: Daphnia Resistance

- ▶ Fit an ANOVA
- ▶ Which groups are different?



47/50

Daphnia Data



48/50

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
```

49/50

High and Med Not Different

```
summary( glht(daphniaLM, linfct=mcp(cyandensity="Tukey")),
         test=adjusted("none"))

#
# Simultaneous Tests for General Linear Hypotheses
#
# Multiple Comparisons of Means: Tukey Contrasts
#
#
# Fit: lm(formula = resistance ~ cyandensity, data = daphnia)
#
# Linear Hypotheses:
#           Estimate Std. Error t value Pr(>|t|)
# med - low == 0  0.09967    0.03496   2.851  0.00794
# high - low == 0  0.11667    0.03496   3.338  0.00233
# high - med == 0  0.01700    0.03651   0.466  0.64496
# (Adjusted p values reported -- none method)
```

50/50