

STA305/1004 - Class 14

March 6, 2017

Today's class

- ▶ ANOVA demonstration
- ▶ Estimating Treatment Effects in ANOVA using Regression
- ▶ Coding Qualitative Predictors in Regression Models

ANOVA Demonstration



▶

Figure 1:

- ▶ Count the total number of each colour (e.g., yellow, purple, pink, green).
- ▶ Eat the Smarties.

ANOVA Data Setup

- ▶ How should the data be setup?

Box	Colour	Count
1	Green	
1	Pink	
1	Purple	
1	Yellow	
2	Green	
2	Pink	
2	Purple	
2	Yellow	
3	Green	
3	Pink	
3	Purple	
3	Yellow	
4	Green	
4	Pink	
4	Purple	
4	Yellow	
5	Green	
5	Pink	
5	Purple	
5	Yellow	

Figure 2:

Smarties Data from 3 boxes

$$\text{Colour 4} = \bar{y}_{4\cdot} - \bar{y}_{1\cdot} = 3.7 - 2.7 = 1$$

```
count <- c(4,3,4,3,1,4,2,5,1,1,2,4)
colour <- as.factor(c(rep("Yellow",3),rep("Purple",3),
                      rep("Green",3),rep("Pink",3)))
#Get means for each flavour
sapply(split(count, colour), mean)
```

	Green	Pink	Purple	Yellow
	2.666667	2.333333	2.666667	3.666667

$$\bar{y}_{..} = \text{ave}$$

$$\text{Intercept} = \bar{y}_{1\cdot} = 2.7 \text{ (Green)}$$

$$\text{Colour 2} = \bar{y}_{2\cdot} - \bar{y}_{1\cdot} = 2.3 - 2.7 = -0.4$$

$$\text{Colour 3} = \bar{y}_{3\cdot} - \bar{y}_{1\cdot} = 2.7 - 2.7 = 0$$

Estimating Treatment Effects in ANOVA using Regression

X_{ij} = not uniquely defined.

ϵ_{ij} = represents the error
 ϵ_{ij} = Variation of # Smarties within each Colour.

- ▶ y_{ij} is the j^{th} observation under the i^{th} treatment.
- ▶ The model for smarties $y_{ij} = \mu + \tau_i + \epsilon_{ij}$, $\epsilon_{ij} \sim N(0, \sigma^2)$ can be written in terms of the dummy variables X_1, X_2, X_3 as:

$$y_{ij} = \mu + \tau_1 X_{1j} + \tau_2 X_{2j} + \tau_3 X_{3j} + \epsilon_{ij}.$$

- What is $y_{ij}, \mu, \tau_i, X_{ij}, \epsilon_{ij}$?

Smarties : Unit of analysis : box

Treatment : Colour

y_{ij} = # Smarties in box j
of Colour i .

Green
Pink
Purple
Yellow

The ANOVA Table

```
#ANOVA table  
anova(lm(count~colour))
```

Analysis of Variance Table

Response: count

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
colour	3	3.000	1.0000	0.4286	0.7381
Residuals	8	18.667	2.3333		

Does not require prob.
Requires probability.

$H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$
 $H_1: \mu_i \neq \mu_j$

$$SS_{\text{Colour}} + SS_{\text{Resid}} = SS_{\text{Total}}$$

$$3.000 + 18.667 = 3 + 18.667$$

Dummy coding

- ▶ Dummy coding compares each level to the reference level. The intercept is the mean of the reference group.
- ▶ Dummy coding is the default in R and the most common coding scheme. It compares each level of the categorical variable to a fixed reference level.

```
contrasts(colour) <- contr.treatment(4) #Treatment contrast  
contrasts(colour) # print dummy coding
```

	X_{1i}	X_{2i}	X_{3i}	
	2	3	4	4
Green	0	0	0	0
Pink	1	0	0	0
Purple	0	1	0	0
Yellow	0	0	1	0

the columns define X_{1i}, X_{2i}, X_{3i}

\rightarrow this is $\underline{\text{O}}$ So that $(X^T X)^{-1}$ exists.

```
lm(count~colour)
```

Call:

```
lm(formula = count ~ colour)
```

Coefficients:
 $\bar{y}_{1.} = -3.3 \times 10^{-1}$

(Intercept) colour2 colour3 colour4
2.667e+00 -3.333e-01 4.710e-16 1.000e+00

$$X_{ij} = \begin{cases} 1 & \text{if Smartie in box}_j \text{ is } \text{pink} \\ 0 & \text{otherwise} \end{cases}$$

$$X_{ij} = \begin{cases} 1 & \text{if Smartie in box}_j \text{ is } \text{purple} \\ 0 & \text{otherwise} \end{cases}$$

$$X_{3j} = \begin{cases} 1 & \text{if Smartie is yellow} \\ 0 & \text{otherwise} \end{cases}$$

least squares (LS)
estimates.

$$Y_{ij} = \mu + \gamma_1 X_{1j} + \gamma_2 X_{2j} + \gamma_3 X_{3j} + \varepsilon_{ij}$$

$$\varepsilon_{ij} \sim N(0, \sigma^2)$$

If the Colour is green (this is the ref. Colour).

$$E(Y_{1j}) = \mu = \mu_1$$

My Colour is pink

$$E(Y_{2j}) = \mu + \gamma_2 = \mu_2 \Rightarrow \gamma_2 = \mu_2 - \mu_1$$

My Colour is purple

$$E(Y_{3j}) = \mu + \gamma_2 = \mu_3 \Rightarrow \gamma_2 = \mu_3 - \mu_1$$

$$E(Y_{4j}) = \mu + \gamma_3 = \mu_4 \Rightarrow \gamma_3 = \mu_4 - \mu_1$$

The LS Estimates :

$$\hat{\mu} = \bar{y}_{1.} = \hat{\mu}_1$$

$$y_{1.} = \sum_{j=1}^{n_1} y_{ij}$$

$$\hat{\gamma}_1 = \bar{y}_{2.} - \bar{y}_{1.}$$

$$\bar{y}_{2.} = \frac{\sum_{j=1}^{n_1} y_{ij}}{n_1}$$

$$\hat{\gamma}_2 = \bar{y}_{3.} - \bar{y}_{1.}$$

$$\hat{\gamma}_3 = \bar{y}_{4.} - \bar{y}_{1.}$$

Deviation coding

- This coding system compares the mean of the dependent variable for a given level to the overall mean of the dependent variable.

```
contrasts(colour) <- contr.sum(4) # Deviation contrast  
contrasts(colour) # print deviation coding
```

	[,1]	[,2]	[,3]
Green	1	0	0
Pink	0	1	0
Purple	0	0	1
Yellow	-1	-1	-1

```
lm(count~colour)
```

Call:

```
lm(formula = count ~ colour)
```

Coefficients:

(Intercept)	colour1	colour2	colour3
2.8333	-0.1667	-0.5000	-0.1667

Compares the mean dependent variable for given colour to overall mean.

$$Y_{ij} = \mu + \gamma_1 X_{1j} + \gamma_2 X_{2j} + \gamma_3 X_{3j} + \epsilon_{ij}$$

$$X_{1j} = \begin{cases} 1 & \text{green} \\ 0 & \text{ow} \\ -1 & \text{yellow} \end{cases}, X_{2j} = \begin{cases} 1 & \text{pink} \\ 0 & \text{ow} \\ -1 & \text{yellow} \end{cases}$$

$$X_{3j} = \begin{cases} 1 & \text{purple} \\ 0 & \text{ow} \\ -1 & \text{yellow} \end{cases}$$

$$E(y_{1j}) = \mu_1 = \mu + \gamma_1$$

$$E(y_{2j}) = \mu_2 = \mu + \gamma_2$$

$$E(y_{3j}) = \mu_3 = \mu + \gamma_3$$

$$E(y_{4j}) = \mu_4 = \mu - \gamma_1 - \gamma_2 - \gamma_3$$

$$\mu = \frac{\mu_1 + \mu_2 + \mu_3 + \mu_4}{4} = \underbrace{\mu + \gamma_1 + \mu + \gamma_2 + \mu + \gamma_3 + \mu - \gamma_1 - \gamma_2 - \gamma_3}_{4}$$

$$\mu_1 - \left(\frac{\mu_1 + \mu_2 + \mu_3 + \mu_4}{4} \right) = \mu + \gamma_1 - \mu = \gamma_1$$

$$\mu_2 - \left(\frac{\mu_1 + \mu_2 + \mu_3 + \mu_4}{4} \right) = \mu + \gamma_2 - \mu = \gamma_2$$

$$\mu_3 \Rightarrow \left(\frac{\mu_1 + \mu_2 + \mu_3 + \mu_4}{4} \right) = \gamma_3$$

LS Estimates :

$$\bar{y}_{..} = \frac{2.7 + 2.3 + 2.7 + 3.7}{4} \\ = 2.85$$

$$\bar{y}_{1.} - \bar{y}_{..} = \hat{\gamma}_1 = 2.7 - 2.85 = 0.15$$

$$y_{2.} - \bar{y}_{..} = \hat{\gamma}_2 = 2.3 - 2.85 = -0.55$$

$$\bar{y}_{3.} - \bar{y}_{..} = \hat{\gamma}_3 = 2.7 - 2.85 = -0.15$$