

# Research Cycle 08: General Linear Model

Dale J. Barr

University of Glasgow

# What is the “General Linear Model” (GLM)?

## Definition (General Linear Model or GLM)

A general mathematical framework for expressing relationships among variables

- Differs from the “cookbook” approach to statistics
  - ▶ *t*-test, ANOVA, ANCOVA,  $\chi^2$  test, regression, correlation, etc.
- Can express/test linear relationships between a numerical dependent variable and any combination of independent variables (categorical or continuous)
- Can even be generalized to categorical dependent variables (through “Generalized Linear Models”; **NB**: advanced)

# ANOVA, Regression, ANCOVA

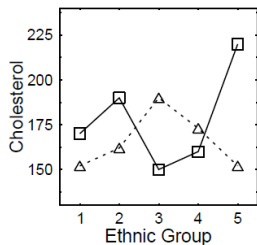


Fig 1a. Cholesterol levels by ethnic group and gender (male=sqr, female=tri).

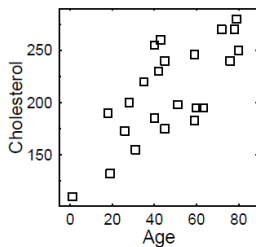


Fig 1a. Cholesterol levels by age.

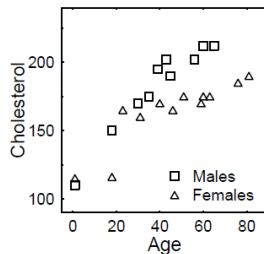


Fig 1a. Cholesterol levels by age and gender.

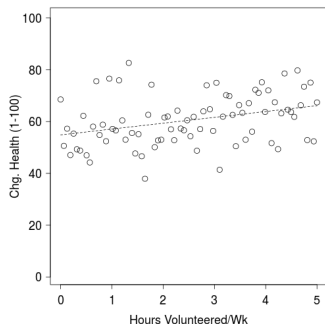
# How the GLM represents relationships

Component of GLM	Notation
DV	$Y$
Grand Average	$\mu$ "mu"
Main Effects	$A, B, C, \dots$
Interactions	$AB, AC, BC, ABC, \dots$
Random Error	$S(\text{Group})$

$$\begin{array}{r} \text{Score} \\ Y \end{array} = \begin{array}{r} \text{Grand Avg.} \\ \mu \end{array} + \begin{array}{r} \text{Main Effects} \\ A + B + C + \dots \end{array} + \begin{array}{r} \text{Interactions} \\ AB + AC + BC + ABC + \dots \end{array} + \begin{array}{r} \text{Error} \\ S(\text{Group}) \end{array}$$

- Components of the model are estimated from the observed data
- Tests are performed (  $F$  ) to see whether its variability is too large to be introduced by chance

# An example: Simple Linear Regression



$$\begin{aligned} Y_i &= \mu + b \times X_i + e_i \\ \text{Score}_i &= \text{Baseline} + \text{Slope} \times \text{Hours}_i + \text{Error}_i \\ Y_i &= 50 + 3 \times X_i + e_i \\ e_i &\sim N(\mu = 0, \sigma^2 = 10) \end{aligned}$$

# Making comparisons across groups

## Example (Spelling)

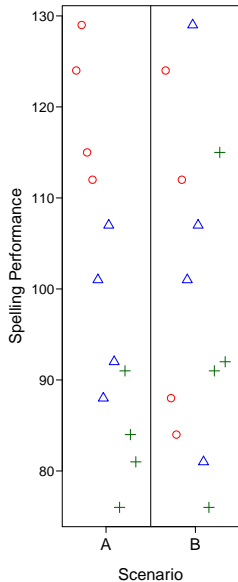
You wish to compare the benefits of three different spelling programs. Do these programs yield differences in spelling performance?

$$H_0 : \mu_1 = \mu_2 = \mu_3$$

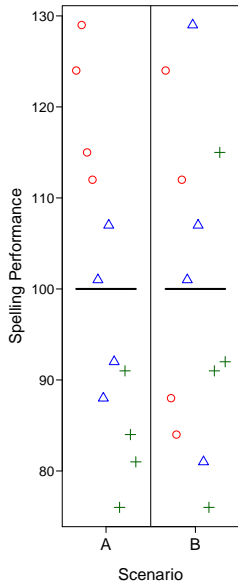
## Factors and Levels

Factor: a categorical variable that is used to divide subjects into groups, usually to draw some comparison. Factors are composed of different *levels*. **Do not confuse factors with levels!**

# Means, Variability, and Deviation Scores



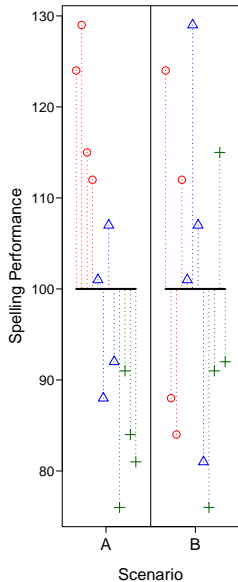
# Means, Variability, and Deviation Scores



$$Y_{..} = \frac{\sum_{ij} Y_{ij}}{N}$$



# Means, Variability, and Deviation Scores

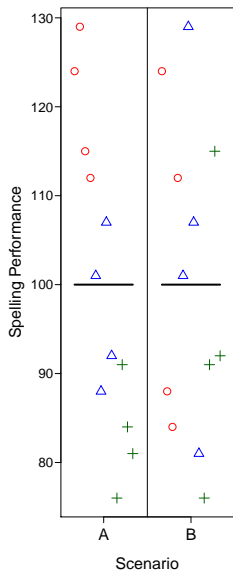


grand mean  $Y_{..} = \frac{\sum_{ij} Y_{ij}}{N}$

$$SD_Y = \sqrt{\frac{\sum_{ij} (Y_{ij} - Y_{..})^2}{N}}$$

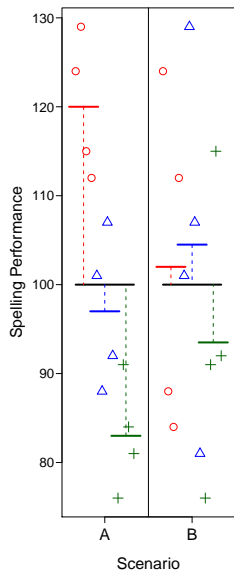
deviation score:  $Y_{ij} - Y_{..}$

# GLM for One-Factor ANOVA



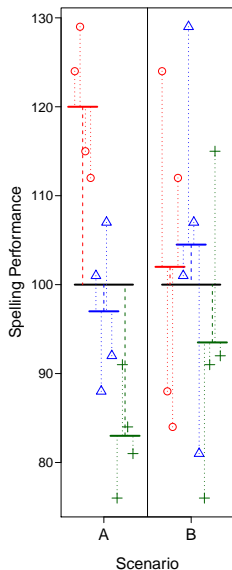
$$Y_{ij} = \mu$$

# GLM for One-Factor ANOVA



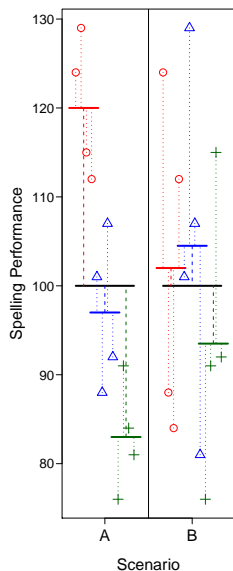
$$Y_{ij} = \mu + A_i$$

# GLM for One-Factor ANOVA



$$Y_{ij} = \mu + A_i + S(A)_{ij}$$

# GLM for One-Factor ANOVA



$$Y_{ij} = \mu + A_i + S(A)_{ij}$$

## Estimation Equations

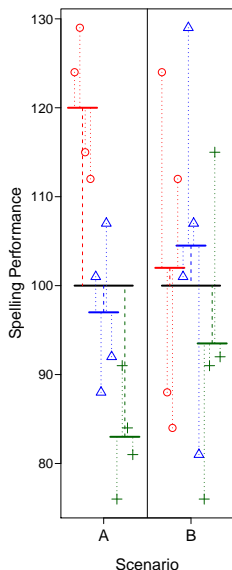
$$\hat{\mu} = Y_{..}$$

$$\hat{A}_i = Y_{i.} - \hat{\mu}$$

$$\widehat{S(A)}_{ij} = Y_{ij} - \hat{\mu} - \hat{A}_i$$

Note that  $\sum_i \hat{A}_i = 0$  and  $\sum_{ij} \widehat{S(A)}_{ij} = 0$

# Sources of Variance



$$Y_{ij} = \mu + A_i + S(A)_{ij}$$

$$Y_{ij} - \mu = A_i + S(A)_{ij}$$

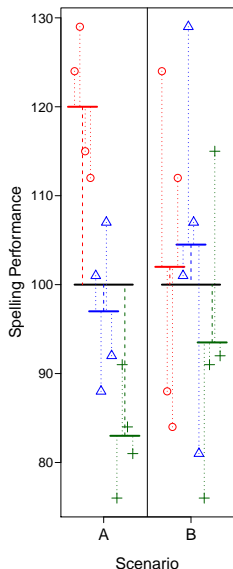
*individual* = *group* + *random*

## Sum of Squares (SS)

A measure of variability consisting of the sum of squared *deviation* scores, where a deviation score is a score minus a mean.

$$SS_A = \sum (Y_i - \mu)^2$$

# Decomposition Matrix



$$\hat{\mu} = 100$$

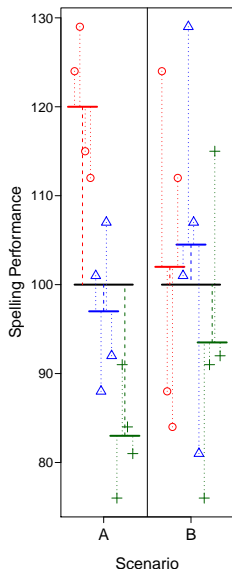
$$\hat{A}_1 = 120 - 100 = 20$$

$$\hat{A}_2 = 97 - 100 = -3$$

$$\hat{A}_3 = 83 - 100 = -17$$

$Y_{ij} =$	$\hat{\mu} +$	$\hat{A}_i +$	$\widehat{S(A)}_{ij}$
124 =	100 +	20 +	4
129 =	100 +	20 +	9
115 =	100 +	20 +	-5
112 =	100 +	20 +	-8
101 =	100 +	-3 +	4
88 =	100 +	-3 +	-9
107 =	100 +	-3 +	10
92 =	100 +	-3 +	-5
76 =	100 +	-17 +	-7
91 =	100 +	-17 +	8
84 =	100 +	-17 +	1
81 =	100 +	-17 +	-2
SS = 123318 =	120000 +	2792 +	526

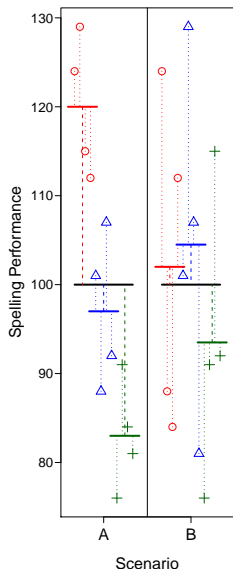
# Logic of ANOVA



- Compare two estimates of the variability, the *between-group* estimate ( $SS_{\text{between}}$ ) and the *within-group* estimate ( $SS_{\text{within}}$ )
- If  $H_0 : \mu_1 = \mu_2 = \mu_3$  is true, then these two measures estimate the same quantity.
- The extent to which the between-group variability exceeds the within-group variability gives evidence against  $H_0 : \mu_1 = \mu_2 = \mu_3$ .



# Calculating $SS_{\text{between}}$ and $SS_{\text{within}}$

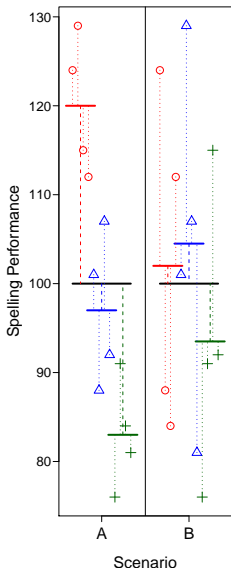


$Y_{ij}$	$=$	$\hat{\mu}$	$+$	$\hat{A}_i$	$+$	$\widehat{S(A)}_{ij}$		
124	=	100	+	20	+	4		
129	=	100	+	20	+	9		
115	=	100	+	20	+	-5		
112	=	100	+	20	+	-8		
101	=	100	+	-3	+	4		
88	=	100	+	-3	+	-9		
107	=	100	+	-3	+	10		
92	=	100	+	-3	+	-5		
76	=	100	+	-17	+	-7		
91	=	100	+	-17	+	8		
84	=	100	+	-17	+	1		
81	=	100	+	-17	+	-2		
$SS$	=	123318	=	120000	+	2792	+	526

check your math

$$SS_Y = SS_{\mu} + SS_A + SS_{S(A)}$$

# $H_0$ and Sums of Squares



$$Y_{ij} - \mu = A_i + S(A)_{ij}$$

## Scenario A

$$SS_A = 2792$$

$$SS_{S(A)} = 526$$

$$SS_A + SS_{S(A)} = 3318$$

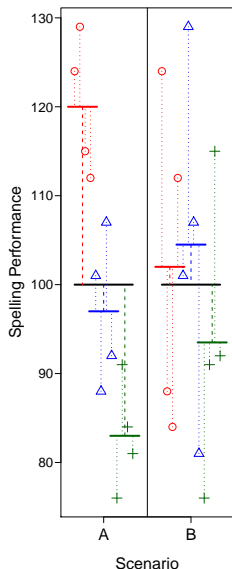
## Scenario B

$$SS_A = 266$$

$$SS_{S(A)} = 3052$$

$$SS_A + SS_{S(A)} = 3318$$

# Mean Square and Degrees of Freedom



## Degrees of Freedom (df)

The number of observations that are “free to vary”.

$$df_A = K - 1$$

$$df_{S(A)} = N - K$$

where  $N$  is the number of subjects and  $K$  is the number of groups.

## Mean Square (MS)

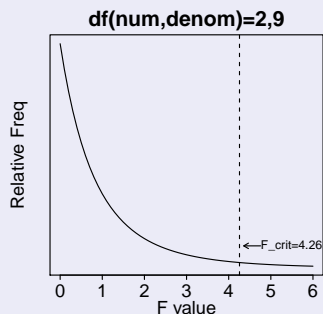
A sum of squares divided by its degrees of freedom.

$$MS_A = \frac{SS_A}{df_A} = \frac{2792}{2} = 1396$$

$$MS_{S(A)} = \frac{SS_{S(A)}}{df_{S(A)}} = \frac{526}{9} = 58.4$$

# The $F$ -ratio

## F density function



If  $F_{obs} > F_{crit}$ , then reject  $H_0$

## F ratio

A ratio of mean squares, with  $df_{\text{numerator}}$  and  $df_{\text{denominator}}$  degrees of freedom.

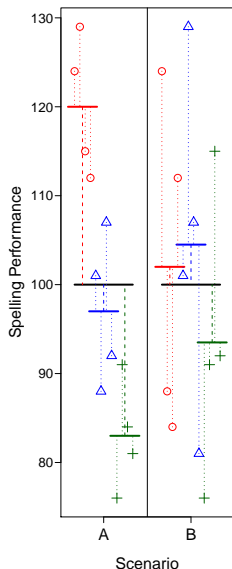
$$F_A = \frac{MS_A}{MS_{S(A)}} = \frac{1396}{58.4} = 23.886$$

df in denominator	df in numerator							
	1	2	3	4	5	6	7	8
1	161.40	199.50	215.70	224.60	230.20	234.00	236.80	238.90
2	18.51	19.00	19.16	19.25	19.30	19.33	19.35	19.37
3	10.13	9.55	9.28	9.12	9.01	8.94	8.89	8.85
4	7.71	6.94	6.59	6.39	6.26	6.16	6.09	6.04
5	6.61	5.79	5.41	5.19	5.05	4.95	4.88	4.82
6	5.99	5.14	4.76	4.53	4.39	4.28	4.21	4.15
7	5.59	4.74	4.35	4.12	3.97	3.87	3.79	3.73
8	5.32	4.46	4.07	3.84	3.69	3.58	3.50	3.44
9	5.12	4.26	3.86	3.63	3.48	3.37	3.29	3.23

# Density/Quantile functions for $F$ -distribution

name	function
<code>pf(x, df1, df2, lower.tail = FALSE)</code>	density (get $p$ given $F_{obs}$ )
<code>qf(p, df1, df2, lower.tail = FALSE)</code>	quantile (get $F_{crit}$ given $p$ )

# Summary Table



## Scenario A

Source	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>p</i>	Error
$\mu$	1	120000	120000.0	2053.232	<.001	S(A)
A	2	2792	1396.0	23.886	<.001	S(A)
S(A)	9	526	58.4			
Total	12	123318				

## Scenario B

Source	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>p</i>	Error
$\mu$	1	120000	120000.0	353.878	<.001	S(A)
A	2	266	133.0	.392	.687	S(A)
S(A)	9	3052	339.1			
Total	12	123318				

# Overview of One-Way ANOVA

- 1 Write the GLM:  $Y_{ij} = \mu + A_i + S(A)_{ij}$
- 2 Write down the estimating equations:
  - ▶  $\hat{\mu} = Y_{..}$
  - ▶  $\hat{A}_i = Y_{i.} - \hat{\mu}$
  - ▶  $\widehat{S(A)}_{ij} = Y_{ij} - \hat{\mu} - \hat{A}_i$
- 3 Compute estimates for all terms in model.
- 4 Create *decomposition matrix*.
- 5 Compute  $SS$ ,  $MS$ ,  $df$ .
  - ▶  $df_{\mu} = 1$
  - ▶  $df_A = K - 1$
  - ▶  $df_{S(A)} = N - K$
  - ▶  $MS = SS/df$
- 6 Construct a summary ANOVA table.
- 7 Compare  $F_{\{obs\}}$  with  $F_{\{crit\}}$ .

R

use the `aov()` function, e.g.:

```
1 spelling$A <- factor(spelling$A)
2 mod <- aov(Y ~ A, data = spelling)
3 summary(mod)
```

<http://talklab.psy.gla.ac.uk/stats/onefactoranova.html#sec-3-2>

- one-sample  $t$ -test  $Y_i - c = \beta_0 + e_i$
- two-sample  $t$ -test  $Y_i = \beta_0 + \beta_1 X_i + e_i$ 
  - ▶ where  $X_i \in (0, 1)$
- paired-samples  $t$ -test  $Y_{1i} - Y_{2i} = \mu + e_i$
- simple linear regression  $Y_i = \beta_0 + \beta_1 X_i + e_i$
- multiple regression  $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + e_i$
- ANCOVA  $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{1i} X_{2i} + e_i$ 
  - ▶ where  $X_{1i} \in (0, 1)$  and  $X_{2i} \in \mathbb{R}$