

STAT 5200 Handout #5.a CRD Tests & ANOVA Table (Ch. 3)

Recall beet lice example of Handout #4: four treatments (A, B, C, D)

→ more than just two-sample t-test; this example is actually a Completely Randomized Design

CRD characteristics:

1. a single experimental factor with $g \geq 2$ levels

↳ #levels ($g=4$ here)

2. experimental units are randomly assigned to factor levels

? all arrangements are equally likely

3. one measurement of response variable is made on each experimental unit

- if multiple response variables are of interest,
we analyze each one separately

4. if # of experimental units assigned to each factor level is the same (say n), then we say this design is balanced

↳ in this course, we're focusing on balanced designs because they tend to have more statistical power

factor has multiple levels (1, 2, ..., g)
- each level corresponds to a specific "treatment"

↓
here, $n=25$

Two main ways to analyze a CRD:

- I. Means model
- II. Effects model

I. Means model

"Y sub i j"

Y_{ij}
 ↓
 Value of response variable for j^{th} experimental unit in factor level i

$i = 1, \dots, g$

=

example of a "parameter" - a fixed, unknown value

μ_i
 ↓
 mean for factor level i (in population)

$j = 1, \dots, n$

↳ or n_i if unbalanced

↳ residual error or random error "difference"

↳ overall sample size $N = g \cdot n$

Test of primary interest:

$H_0 : \mu_1 = \dots = \mu_g$

$H_a : \mu_i$'s not all equal

Can estimate pop. means μ_i with: sample means

$\hat{\mu}_i = \bar{Y}_{i\cdot} = \frac{1}{n} \sum_{j=1}^n Y_{ij} = \text{sample mean in factor level } i$
 ↳ "Y i dot bar"

Construct test statistic to measure size of differences among $\hat{\mu}_i$'s

↳ will be a standardized differences

Means model discussion points:

- very simple & intuitive
- easy to formulate sensible hypotheses
- obvious parameter estimates for μ_i 's, with nice mathematical properties
- BUT - hard to generalize to more than one factor



- especially when combined effect of multiple factors is of interest

II. Effects model

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij}$$

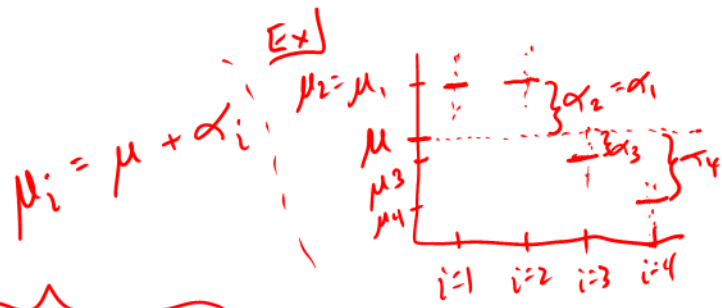
\downarrow
 value of response variable for j^{th} exp. unit in factor level i

\downarrow
 overall or grand mean

\downarrow
 effect of factor level i

\downarrow
 residual error

$i = 1, \dots, g$ $j = 1, \dots, n$ $N = g \cdot n$



Look at # of parameters:

- Here, there are $\mu, \alpha_1, \alpha_2, \dots, \alpha_g \rightarrow g+1$

- But how many groups of data?

g (# factor levels)

Need to impose some "identifiability" constraint to ensure unique estimates

- Default in SAS is to constrain: $\alpha_g \equiv 0$, so really, only have g parameters

– Then $\hat{\mu} = \bar{Y}_g$, and for each $i = 1, \dots, g-1$, $\hat{\alpha}_i = \bar{Y}_i - \bar{Y}_g$.

– So for every i , $\hat{Y}_{ij} = \hat{\mu} + \hat{\alpha}_i = \bar{Y}_i$.

- Another option (appropriate for balanced designs, and often useful for pedagogical purposes) is the "sum-to-zero" constraint:

$$\sum_{i=1}^g \alpha_i = 0$$

$$\hat{\mu} = \bar{Y}_{..} = \frac{1}{N} \sum_{ij} Y_{ij} = \frac{1}{g \cdot n} \sum_{i=1}^g \sum_{j=1}^n Y_{ij} = \frac{1}{g} \sum_{i=1}^g \bar{Y}_i$$

$$\hat{\alpha}_i = \bar{Y}_i - \bar{Y}_{..}$$

$$\sum_{i=1}^g \hat{\alpha}_i = \sum_{i=1}^g (\bar{Y}_i - \bar{Y}_{..}) = \sum_{i=1}^g \bar{Y}_i - g \bar{Y}_{..} = 0$$

Compare the means model and effects model:

$$\begin{aligned}
 Y_{ij} &= \mu_i + \epsilon_{ij} && \text{means model} \\
 &= \mu + \alpha_i + \epsilon_{ij}, && \text{effects model} \\
 & && (i = 1, \dots, g ; j = 1, \dots, n)
 \end{aligned}$$

$$\hat{Y}_{ij} = \hat{\mu}_i = \hat{\mu} + \hat{\alpha}_i$$

In both models, parameter estimates are obtained by minimizing:

$$\sum_{ij} (Y_{ij} - \hat{Y}_{ij})^2$$

make predicted \hat{Y} close to observed Y

$$\hat{\epsilon}_{ij} = \epsilon_{ij} \left\{ \begin{array}{l} \text{residual for exp. unit } j \\ \text{in factor level } i \end{array} \right.$$

- this is the principle of least squares, and gives estimates nice properties

minimum bias
! minimum variance

To construct a test statistic for $H_0 : \mu_1 = \dots = \mu_g$ (or equivalently, $H_0 : \alpha_1 = \dots = \alpha_g$), we need to “decompose” observations into statistics. Here, we’ll focus on the effects model, with sum-to-zero constraint):

$$\begin{aligned}
 \frac{1}{n} \sum_{j=1}^n Y_{ij} &= \bar{Y}_i \\
 \frac{1}{g} \sum_{i=1}^g \bar{Y}_i &= \bar{Y}_{..}
 \end{aligned}$$

$$\begin{aligned}
 Y_{ij} &= \mu + \alpha_i + \epsilon_{ij} \\
 Y_{ij} &= \bar{Y}_{..} + (\bar{Y}_i - \bar{Y}_{..}) + (Y_{ij} - \bar{Y}_i) = a + b + c
 \end{aligned}$$

$\hat{Y}_{ij} = \bar{Y}_i$

Square both sides and sum over both subscripts:

$$(a + b + c)^2 = a^2 + b^2 + c^2 + \dots$$

$$\begin{aligned}
 \sum_{i=1}^g \sum_{j=1}^n Y_{ij}^2 &= \sum_{i=1}^g \sum_{j=1}^n (\bar{Y}_{..})^2 + \sum_{i=1}^g \sum_{j=1}^n (\bar{Y}_i - \bar{Y}_{..})^2 + \sum_{i=1}^g \sum_{j=1}^n (Y_{ij} - \bar{Y}_i)^2 \\
 &+ 2 \sum_{ij} [\bar{Y}_{..} (\bar{Y}_i - \bar{Y}_{..}) + \bar{Y}_{..} (Y_{ij} - \bar{Y}_i) + (\bar{Y}_i - \bar{Y}_{..}) (Y_{ij} - \bar{Y}_i)]
 \end{aligned}$$

cross product terms

these cross-product terms cancel out & sum to 0

$$= N\bar{Y}_{..}^2 + n \sum_{i=1}^g (\bar{Y}_i - \bar{Y}_{..})^2 + \sum_{i=1}^g \sum_{j=1}^n (Y_{ij} - \bar{Y}_i)^2$$

Now subtract $N\bar{Y}_{..}^2$ from both sides:

$$\sum_{i=1}^g \sum_{j=1}^n Y_{ij}^2 - N\bar{Y}_{..}^2 = n \sum_{i=1}^g (\bar{Y}_i - \bar{Y}_{..})^2 + \sum_{i=1}^g \sum_{j=1}^n (Y_{ij} - \bar{Y}_i)^2$$

$\epsilon_{ij} = Y_{ij} - \hat{Y}_{ij}$

Sum of Squares Total

$$\rightarrow SS_T \propto \text{Var}(Y) \quad SS_{T \text{ trt}} \quad SS_E$$

SS_T = corrected total sum of squares:

$$\begin{aligned}
 (N-1)Var(Y) &= \sum_{ij} (Y_{ij} - \bar{Y}_{..})^2 \\
 &= \sum_{ij} (Y_{ij}^2 - 2Y_{ij}\bar{Y}_{..} + \bar{Y}_{..}^2) \\
 &= \sum_{ij} Y_{ij}^2 - 2\bar{Y}_{..} \sum_{ij} Y_{ij} + N\bar{Y}_{..}^2 \\
 &= \sum_{ij} Y_{ij}^2 - 2\bar{Y}_{..} (N\bar{Y}_{..}) + N\bar{Y}_{..}^2 \\
 &= \sum_{i=1}^g \sum_{j=1}^n Y_{ij}^2 - N\bar{Y}_{..}^2 = SS_T \quad \& \quad Var(Y)
 \end{aligned}$$

overall sample mean (with arrow pointing to $\bar{Y}_{..}$)

SS_{Trt} = sum of squares due to treatment (factors in model)

↳ like variance between factor levels

SS_E = sum of squares due to error; also "residual sum of squares":
 (This is the quantity that our parameter estimates are guaranteed to minimize.)

$$\sum_{ij} e_{ij}^2 = \sum_{ij} (Y_{ij} - \hat{Y}_{ij})^2$$

↑ $\bar{Y}_{.i}$

↳ like variance within factor levels

Note that $SS_{Trt} = n \sum_{i=1}^g \hat{\alpha}_i^2$, so we can use this to test

"no trt. effect" $\rightarrow H_0 : \alpha_1 = \dots = \alpha_g = 0$ (or $\mu_1 = \dots = \mu_g$)
 vs.
 $H_a : \alpha_i \neq 0$ for at least one i (or $\mu_i \neq \mu_{i'}$ for some i, i')

implied by constraint

- If H_0 is true, SS_{Trt} should be: *smaller*
- If H_a is true, SS_{Trt} should be: *larger*
- How large must SS_{Trt} be to say H_a is true? *Need to standardize first*

Test statistic for $H_0 : \alpha_1 = \dots = \alpha_g$:

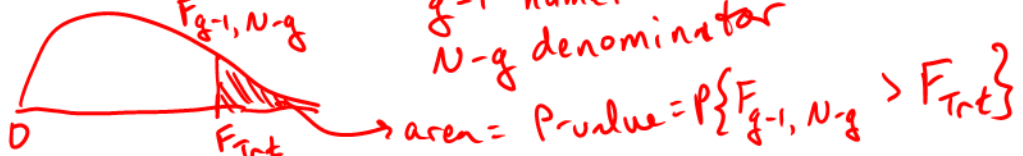
$$F_{Trt} = \frac{SS_{Trt} / (g-1)}{SS_E / (N-g)} = \frac{MS_{Trt}}{MS_E}$$

sampling dist'n

Statistical theory says under H_0 , $F_{Trt} \sim F_{g-1, N-g}$

if H_0 true and model assumptions are met ($H_0 \neq L_i, b_a$)

Sampling distribution is F with degrees of freedom: $g-1$ numerator, $N-g$ denominator



We summarize this decomposition and test statistic in the ANOVA table:

- “Analysis of Variance”: we partition (or decompose) variance (in response Y) into components, or sources (of variance)

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	$g - 1$	SS_{Trt}	$MS_{Trt} = \frac{SS_{Trt}}{g-1}$	$F_{Trt} = \frac{MS_{Trt}}{MSE}$	P-value = $P\{F_{g-1, N-g} > F_{Trt}\}$
Error	$N - g$	SS_E	$MSE = \frac{SS_E}{N-g}$		
Corrected Total	$N-1$	SS_T			

In beet lice data (Handout # 5), $F_{Trt} = 13.86$, with p-value $< .0001$

Conclude: At least two chemical treatments have different effects

Now – which two (or more) factor levels are significantly different from each other?

- This involves post hoc mean comparisons (H0 # 7 ! 7a)
- We’ll come back to this later, because before we can trust any inference, we need to assess model assumptions (H0 # 6 ! 6a)

NOTE: DF = degrees of freedom, or # of unconstrained data points

- Every time you estimate a parameter, you lose one DF
- With sum-to-zero constraint: $N = \text{total sample size} \rightarrow \sum_i \alpha_i = 0$
 - Estimate $\hat{\mu} = \bar{Y}.$ → lose 1 DF → DF for SS_T is $N-1$
(Knowing average $\bar{Y}.$ and Y_{ij} values for $N - 1$ observations would allow you to calculate the remaining observation)
 - Estimate $\hat{\alpha}_1, \dots, \hat{\alpha}_{g-1}$ → lose $g - 1$ DF → DF for SS_{Trt} is $g-1$
(Why not α_g ? → constraint)
 - This leaves $N - g$ DF for SS_E
↳ lost g DF to get there

Another useful statistic from ANOVA table:

- $SS_T = SS_{Trt} + SS_E$

$\rightarrow \propto$ variation in y ; unexplained by model

\propto total variation in y

$\rightarrow \propto$ variation in y "explained" by model

$$R^2 = \frac{SS_{Trt}}{SS_T} = 1 - \frac{SS_E}{SS_T}$$

= proportion of variation (or variability) in observed response values "explained" by differences in factor level effects (trt. effects)

- $0 \leq R^2 \leq 1$

- small R^2 : variability not really explained by factor levels

- large R^2 : factor levels cause most of the variability

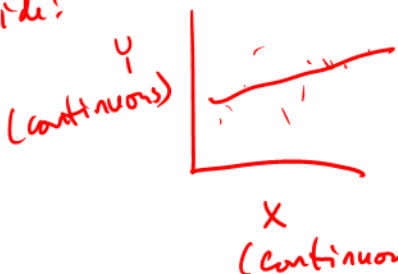
- If $R^2 \approx 1$, then even small differences among factor levels may be "statistically significant"

(even if inconsequential)

- How high must R^2 be for a "good" model?

- don't rely too much on this ;
it's usually best to just report R^2

Aside:



Linear regression also has R^2

where interpretation is same and R is $\text{corr}(X, Y)$