## STAT 705 Chapter 16: One-way ANOVA

#### **Timothy Hanson**

#### Department of Statistics, University of South Carolina

Stat 705: Data Analysis II

Analysis of variance (ANOVA) models are regression models with qualitative predictors, called <u>factors</u> or <u>treatments</u>.

Factors have different levels.

For example, the factor "education" may have the levels *high school, undergraduate, graduate.* The factor "gender" has two levels *female, male.* 

We may have several factors as predictors, e.g. race and gender may be used to predict annual salary in \$.

There are two types of factors:

- Classification (investigator cannot control).
- Experimental (investigator can control).

A <u>control treatment</u> (or control factor level) is sometimes used to measure effects of (new or experimental) treatments under investigation, relative to the "status quo."

e.g. ibuprofin, aspirin, and placebo. We have 3 factor levels. Without placebo, we do not know how iboprofin or aspirin does relative to no pain killer, only relative to each other.

Uses of ANOVA models: find best/worst treatment, measure effectiveness of new treatment, compare treatments.

Often interested in determining whether there is a *difference* in treatments.

Read Sections 16.1–16.8 in the text.

Have r different treatments or factor levels. At each level i, have  $n_i$  observations from group i.

Total number of observations is  $n_T = n_1 + n_2 + \cdots + n_r$ .

Example: Two factors: MS, PhD.  $Y_{ij}$  is age in years. Spring of 2014 we observe

$$Y_{11} = 28, Y_{12} = 24, Y_{13} = 24, Y_{14} = 22, Y_{15} = 26, Y_{16} = 23,$$

 $Y_{21} = 29, Y_{22} = 23, Y_{23} = 26, Y_{24} = 25, Y_{25} = 22, Y_{26} = 23, Y_{27} = 38, Y_{28} = 33, Y_{29} = 30, Y_{2,10} = 27.$ 

## One-way ANOVA model

$$Y_{ij} = \mu_i + \epsilon_{ij}, \ \ \epsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2).$$

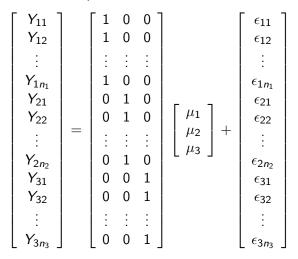
Can rewrite as

$$Y_{ij} \stackrel{ind.}{\sim} N(\mu_i, \sigma^2).$$

- Data are normal, data are independent, variance constant across groups.
- μ<sub>i</sub> is allowed to be different for each group. μ<sub>1</sub>,..., μ<sub>r</sub> are the r population means of the response. A picture helps.
- Questions: what is  $E\{Y_{ij}\}$ ? What is  $\sigma^2\{Y_{ij}\}$ ?

### Matrix formulation

(pp. 683–684, 710–712) For r = 3 we have



or

 $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.$ 

For 
$$r=3$$
, let  $Q(\mu_1,\mu_2,\mu_3)=\sum_{i=1}^3\sum_{j=1}^{n_i}(Y_{ij}-\mu_i)^2$  .

Need to minumize this over all possible  $(\mu_1, \mu_2, \mu_3)$  to find least-squares (LS) solution. Can easily show that  $Q(\mu_1, \mu_2, \mu_3)$  has minimum at

$$\hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\mu}_1 \\ \hat{\mu}_2 \\ \hat{\mu}_3 \end{bmatrix} = \begin{bmatrix} \bar{Y}_{1\bullet} \\ \bar{Y}_{2\bullet} \\ \bar{Y}_{3\bullet} \end{bmatrix}$$

where  $\bar{Y}_{i\bullet} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}$  is the sample mean from the *i*th group (pp. 687–688).

These  $\hat{oldsymbol{eta}}$  are also maximum likelihood estimates.

# Matrix formula of least-squares estimators (r = 3)

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} 1 & \cdots & 1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 1 & \cdots & 1 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 & 1 & \cdots & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} n_1 & 0 & 0 \\ 0 & n_2 & 0 \\ 0 & 0 & n_3 \end{bmatrix},$$
$$(\mathbf{X}'\mathbf{X})^{-1} = \begin{bmatrix} n_1^{-1} & 0 & 0 \\ 0 & n_2^{-1} & 0 \\ 0 & n_2^{-1} & 0 \\ 0 & 0 & n_3^{-1} \end{bmatrix}, \quad \mathbf{X}'\mathbf{Y} = \begin{bmatrix} \mathbf{Y}_{1\bullet} \\ \mathbf{Y}_{2\bullet} \\ \mathbf{Y}_{3\bullet} \end{bmatrix},$$
$$\Rightarrow \hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = \begin{bmatrix} \tilde{\mathbf{Y}}_{1\bullet} \\ \tilde{\mathbf{Y}}_{2\bullet} \\ \tilde{\mathbf{Y}}_{3\bullet} \end{bmatrix}.$$

As in regression (STAT 704),

$$e_{ij}=Y_{ij}-\hat{Y}_{ij}=Y_{ij}-\hat{\mu}_i=Y_{ij}-ar{Y}_{iullet}.$$

As usual,  $\hat{Y}_{ij}$  is the estimated mean response under the model. Note that  $\sum_{j=1}^{n_i} e_{ij} = 0$ . [check this!]

In matrix terms

$$\mathbf{e} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{Y} - \hat{\mathbf{Y}}.$$

- r = 4 box designs for a new breakfast cereal.
- 20 stores w/ roughly equal sales volumes picked to participate; n<sub>i</sub> = 5 is planned for each.
- A fire occurred at one store that had design 3, so ended up with  $n_T = 19$  instead of 20, and  $n_1 = n_2 = n_4 = 5$  and  $n_3 = 4$ .

```
data kenton;
input sales design @0;
datalines;
11 1 17 1 16 1 14 1 15 1 12 2 10 2 15 2 19 2 11 2
23 3 20 3 18 3 17 3 27 4 33 4 22 4 26 4 28 4
;
proc sgscatter;
plot sales*design;
run;
proc glm plots=all; * zero/one dummy variables, but recover cell means via lsmeans;
class design;
model sales=design;
lsmeans design;
run;
```

Define the following

$$Y_{i\bullet} = \sum_{j=1}^{n_i} Y_{ij} = i \text{ group sum},$$
$$\bar{Y}_{i\bullet} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} = i \text{th group mean}$$
$$Y_{\bullet\bullet} = \sum_{i=1}^r \sum_{j=1}^{n_i} Y_{ij} = \sum_{i=1}^r Y_{i\bullet} = \text{sum all obs.}$$
$$1 \sum_{i=1}^r \sum_{j=1}^{n_i} Y_{ij} = \frac{1}{r} \sum_{i=1}^r Y_{i\bullet} = \text{sum all obs.}$$

$$\bar{Y}_{\bullet\bullet} = \frac{1}{n_T} \sum_{i=1}^r \sum_{j=1}^r Y_{ij} = \frac{1}{n_T} \sum_{i=1}^r Y_{i\bullet} = \text{mean all obs.}$$

### Sums of squares for treatments, error, and total

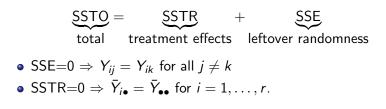
SSTO = 
$$\sum_{i=1}^{r} \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{\bullet \bullet})^2 = \text{variability in } Y_{ij}\text{'s}$$
SSTR = 
$$\sum_{i=1}^{r} \sum_{j=1}^{n_i} (\hat{Y}_{ij} - \bar{Y}_{\bullet \bullet})^2 = \sum_{i=1}^{r} \sum_{j=1}^{n_i} (\hat{\mu}_{ij} - \bar{Y}_{\bullet \bullet})^2$$

$$= \sum_{i=1}^{r} \sum_{j=1}^{n_i} (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet \bullet})^2 = \sum_{i=1}^{r} n_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet \bullet})^2$$

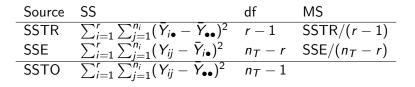
$$= \text{variability explained by ANOVA model}$$
SSE = 
$$\sum_{i=1}^{r} \sum_{j=1}^{n_i} (Y_{ij} - \hat{Y}_{ij})^2 = \sum_{i=1}^{r} \sum_{j=1}^{n_i} e_i^2$$

= variability NOT explained by ANOVA model

• As before in regression,



## ANOVA table (p. 694)



- SSTO has  $n_T 1$  df because there are  $n_T Y_{ij} \bar{Y}_{\bullet\bullet}$  terms in the sum, but they add up to zero (1 constraint).
- SSE has n<sub>T</sub> − r df because there are n<sub>T</sub> Y<sub>ij</sub> − Y
  <sub>i•</sub> terms in the sum, but there are r constraints of the form
   ∑<sub>j=1</sub><sup>n<sub>i</sub></sup> (Y<sub>ij</sub> − Y
  <sub>i•</sub>) = 0.
- SSTR has r 1 df because there are r terms  $n_i(\bar{Y}_{i\bullet} \bar{Y}_{\bullet\bullet})$  in the sum, but they sum to zero (1 constraint).

$$E\{\mathsf{MSE}\} = \sigma^2, \quad \mathsf{MSE} \text{ is unbiased estimate of } \sigma^2$$
$$E\{\mathsf{MSTR}\} = \sigma^2 + \frac{\sum_{i=1}^r n_i(\mu_i - \mu_{\bullet})^2}{r - 1},$$

where  $\mu_{\bullet} = \sum_{i=1}^{r} \frac{n_{i}\mu_{i}}{n_{T}}$  is weighted average of  $\mu_{1}, \ldots, \mu_{r}$  (pp. 696–698).

If  $\mu_i = \mu_j$  for all  $i, j \in \{1, ..., r\}$  then  $E\{MSTR\} = \sigma^2$ , otherwise  $E\{MSTR\} > \sigma^2$ .

Hence, if any group means are different then  $\frac{E\{MSTR\}}{E\{MSE\}} > 1$ .

## 16.6 F test of $H_0: \mu_1 = \cdots = \mu_r$

<u>Fact</u>: If  $\mu_1 = \cdots = \mu_r$  then

$$F^* = rac{\mathsf{MSTR}}{\mathsf{MSE}} \sim F(r-1, n_T - r).$$

To perform  $\alpha$ -level test of  $H_0: \mu_1 = \cdots = \mu_r$  vs.  $H_a$ : some  $\mu_i \neq \mu_j$  for  $i \neq j$ ,

- Accept if  $F^* \leq F(1 \alpha, r 1, n_T r)$  or p-value  $\geq \alpha$ .
- Reject if  $F^* > F(1 \alpha, r 1, n_T r)$  or p-value  $< \alpha$ .

p-value =  $P\{F(r-1, n_T - 1) \ge F^*\}$ .

Example: Kenton Foods

- If r = 2 then F\* = (t\*)<sup>2</sup> where t\* is t-statistic from 2-sample pooled-variance t-test.
- The F-test may be obtained from the general nested linear hypotheses approach (big model / little model). Here the full model is  $Y_{ij} = \mu_i + \epsilon_{ij}$  and the reduced is  $Y_{ij} = \mu + \epsilon_{ij}$ .

$$F^* = \frac{\left[\frac{SSE(R) - SSE(F)}{dfE_R - dfE_F}\right]}{\frac{SSE(F)}{dfE_F}} = \frac{MSTR}{MSE}$$

## 16.7 Alternative formulations

SAS will fit the cell means model (discussed so far) with a noint option in model statement; however, the F-test will not be correct. Your textbook discusses an alternative parameterization that is not easy to get out of the SAS procedures we will use.

By default, SAS fits the model

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

where  $\alpha_r = 0$ .

- $E\{Y_{rj}\} = \mu$ ;  $\mu$  is the cell-mean for the *r*th level.
- For i < r,  $E\{Y_{ij}\} = \mu + \alpha_i$ ;  $\alpha_i$  is *i*'s offset to group *r*'s mean  $\mu$ .

Note that SAS's default corresponds to a regression model where categorical predictors are modeled using the usual zero-one dummy variables. In class, let's find the design **X** for SAS's model for r = 3 and  $n_1 = n_2 = n_3 = 2$ .

Even though SAS parameterizes the model differently, with the rth level as baseline, the ANOVA table and F-test is the same as the cell means model.

Also  $\hat{\mu} = \bar{Y}_{r\bullet}$  and  $\hat{\alpha}_i = \bar{Y}_{i\bullet} - \bar{Y}_{r\bullet}$  are the OLS and MLE estimators. These are reported in SAS. Use, e.g. model sales=design / solution;

The cell means  $\hat{\mu}_i$  are obtained in SAS by adding lsmeans to glm or glimmix.