# STAT120C: Analysis of Variance (ANOVA)

So far we have considered only one or two samples. For one sample, we were concerned by the population mean. For two samples, we were concerned by the difference of two population means. What if there are more than two samples? In the two-sample t test, we considered only one factor. What if there is another factor that might affect the outcome variable? For example, we are interested to learn the effects of gender and smoking on body mass index. These questions can be answered using a more general framework: the analysis of variance. We will spend two weeks on this topic.

# 1   One-Way ANOVA

## 1.1   Introduction

**One-way layout/design**: a one-way layout / design is an experiment design in which independent measurements are made under each of several conditions. It generalizes the the design of two independent samples.

**Example Chlorpheniramine maleate in tablets (Kirchhoefer 1979)** : To study the level of chlorpheniramine maleate in tablets from seven labs, measurements of composites that had nominal dosages equal to 4mg from the seven labs. And 10 measurements were made by each lab.
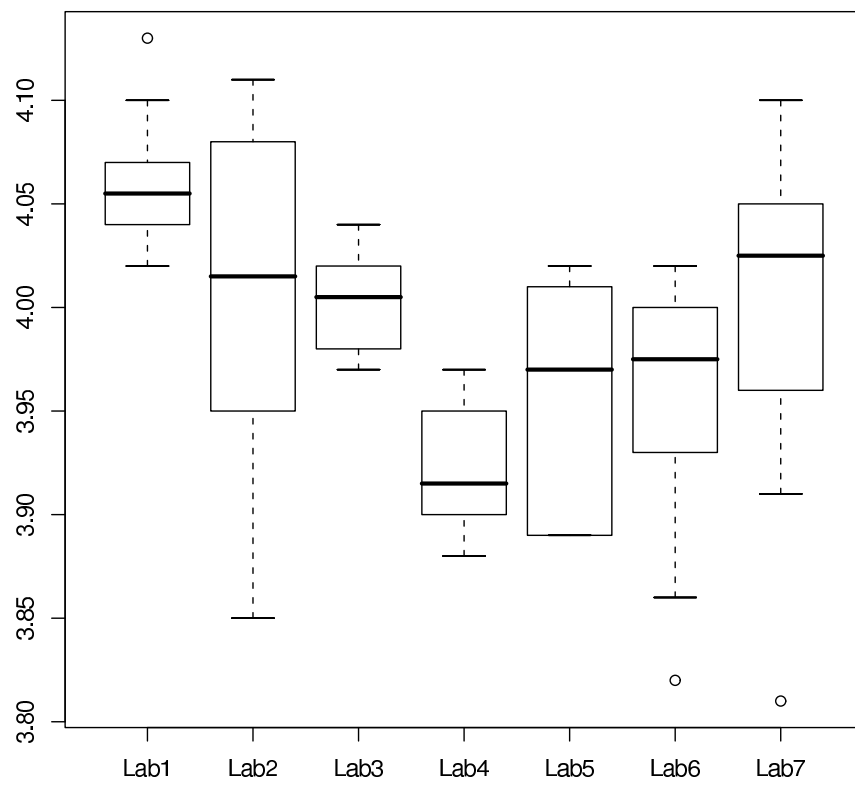
| | Observations(j) | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | $\bar{y}_i$ |
| Lab1 | 4.13 | 4.07 | 4.04 | 4.07 | 4.05 | 4.04 | 4.02 | 4.06 | 4.10 | 4.04 | 4.062 |
| Lab2 | 3.86 | 3.85 | 4.08 | 4.11 | 4.08 | 4.01 | 4.02 | 4.04 | 3.97 | 3.95 | 3.997 |
| Lab3 | 4.00 | 4.02 | 4.01 | 4.01 | 4.04 | 3.99 | 4.03 | 3.97 | 3.98 | 3.98 | 4.003 |
| Lab4 | 3.88 | 3.88 | 3.91 | 3.95 | 3.92 | 3.97 | 3.92 | 3.90 | 3.97 | 3.90 | 3.920 |
| Lab5 | 4.02 | 3.95 | 4.02 | 3.89 | 3.91 | 4.01 | 3.89 | 3.89 | 3.99 | 4.00 | 3.957 |
| Lab6 | 4.02 | 3.86 | 3.96 | 3.97 | 4.00 | 3.82 | 3.98 | 3.99 | 4.02 | 3.93 | 3.955 |
| Lab7 | 4.00 | 4.02 | 4.03 | 4.04 | 4.10 | 3.81 | 3.91 | 3.96 | 4.05 | 4.06 | 3.998 |

A boxplot of the data is shown in the figure below. There are two sources of variability in the data: variability within labs and variability between labs.

Let $\mu_i$ be the mean level of chlorpheniramine in tablets from Lab $i$. We might be interested in many questions. For example,

(a) Is the mean level of chlorpheniramine maleate in tablets from Lab 1 different from 4? (one-sample t test)

$$\mu_1 = 4$$

(b) Is the mean level of chlorpheniramine maleate in tablets from Lab 1 different from that from Lab 2? (two-sample t test)

$$\mu_1 = \mu_2$$

(c) Do mean levels differ across the seven labs? (which test?)

$$\mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5 = \mu_6 = \mu_7$$

Questions (a) and (b) can answered using one-sample and two-sample t-tests, respectively. To answer question (c), we need to use a more general framework.

## 1.2 The $F$ test under balanced designs

### 1.2.1 Notation

In two-sample t-test, we used notation $X_i$ and $Y_i$. When there are more than two conditions, we shall use $Y_{ij}$, which denotes the $j$th measurement under the $i$th condition.

| Treatment Group | Observations(j) | | | |
|---|---|---|---|---|
| 1 | $y_{11}$ | $y_{12}$ | $\cdots$ | $y_{1J}$ |
| 2 | $y_{21}$ | $y_{22}$ | $\cdots$ | $y_{2J}$ |
| . | . | . | $\cdots$ | . |
| . | . | . | $\cdots$ | . |
| . | . | . | $\cdots$ | . |
| I | $y_{I1}$ | $y_{I2}$ | $\cdots$ | $y_{IJ}$ |

For two-sample t test, we assumed that

$$X_i \overset{iid}{\sim} N(\mu_X, \sigma^2),$$

$$Y_i \overset{iid}{\sim} N(\mu_Y, \sigma^2)$$

and

$$X_1, \cdots, X_m \text{ and } Y_1, \cdots, Y_n \text{ are independent.}$$

Similar assumptions will be made for the one-way layout. Here $Y_{ij}$ is the $j$th observation under the $i$th treatment/condition. The **statistical model** we use for one-way ANOVA is

$$y_{ij} = \mu_i + \epsilon_{ij},$$

where $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$, for $i = 1, ..., I; j = 1, \cdots, J$.

3

In the model, $\mu_i$ is the mean of the $i$the treatment/condition. This is a balanced design. We say it is balanced because all groups have the same number of observations.

The model can also be written as

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

with $\sum \alpha_i = 0$ and $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$, for $i = 1, 2, ..., I; j = 1, \cdots, J$.

The design is balanced. In this situation, $\mu = \sum_{i=1}^{I} \mu_i / I$ and $\alpha_i = \mu_i - \mu$. Here $\mu$ is the overall mean level, $\alpha_i$ is the differential effect of the $i$th treatment, and $\epsilon_{ij}$ is the random error in the $j$th observation under the $i$th treatment. The errors are assumed to be iid $N(0, \sigma^2)$. Because $\alpha_i$'s are deviations from the overall mean:

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

The expected response to the $i$th treatment is $E(Y_{ij}) = \mu + \alpha_i$. Thus, if $\alpha_i = 0$, for $i = 1, 2, \cdots, I$, all treatment have the same expected response. In general, $\alpha_i - \alpha_j$ is the difference between the expected outcome values under treatments/conditions $i$ and $j$.

We introduce the following summary statistics:

$Y_{i\cdot} = \sum_{j=1}^{J} Y_{ij}$

$\bar{Y}_{i\cdot} = \sum_{j=1}^{J} Y_{ij} / J$

$Y_{\cdot\cdot} = \sum_{i=1}^{I} \sum_{j=1}^{J} Y_{ij}$

$\bar{Y}_{\cdot\cdot} = \sum_{i=1}^{I} \sum_{j=1}^{J} Y_{ij} / (IJ)$

### 1.2.2  Decomposition of the total sum of squares

The analysis of variance is based on the following decomposition of the Sum of Squares of TOtal (SSTO):

$$SSTO = \sum_{i=1}^{I} \sum_{j=1}^{J} (Y_{ij} - \bar{Y}_{\cdot\cdot})^2$$

4

In the following, we show that $SSTO$ can be decomposed into two sums of squares.

$$
\begin{aligned}
SSTO &= \sum_{i=1}^{I}\sum_{j=1}^{J}(Y_{ij}-\bar{Y}_{..})^2 \\
&= \sum_{i=1}^{I}\sum_{j=1}^{J}[(Y_{ij}-\bar{Y}_{i\cdot})+(\bar{Y}_{i\cdot}-\bar{Y}_{..})]^2 \\
&= \sum_{i=1}^{I}\sum_{j=1}^{J}(Y_{ij}-\bar{Y}_{i\cdot})^2 + \sum_{i=1}^{I}\sum_{j=1}^{J}(\bar{Y}_{i\cdot}-\bar{Y}_{..})^2 + 2\sum_{i=1}^{I}\sum_{j=1}^{J}(Y_{ij}-\bar{Y}_{i\cdot})(\bar{Y}_{i\cdot}-\bar{Y}_{..}) \\
&= \sum_{i=1}^{I}\sum_{j=1}^{J}(Y_{ij}-\bar{Y}_{i\cdot})^2 + J\sum_{i=1}^{I}(\bar{Y}_{i\cdot}-\bar{Y}_{..})^2 + 2\sum_{i=1}^{I}(\bar{Y}_{i\cdot}-\bar{Y}_{..})[\sum_{j=1}^{J}(Y_{ij}-\bar{Y}_{i\cdot})]
\end{aligned}
$$

The last term of the final expression vanished because the sum of deviations from a mean is zero.

This equation says that the sum of squares of total (SSTO) can be decomposed to the sum of squares within treatment groups (SSW) plus the sum of squares between treatment groups (SSB). Therefore, it is often written as

$$SSTO = SSB + SSW,$$

where $SSB = \sum_{i}\sum_{j}(\bar{Y}_{i\cdot}-\bar{Y}_{..})^2 = J\sum_{i}(\bar{Y}_{i\cdot}-\bar{Y}_{..})^2$ and $SSW = \sum_{i}\sum_{j}(Y_{ij}-\bar{Y}_{i\cdot})^2$.

**Note:**

- Sometimes $SSW$ is referred to as $SSE$, and $SSB$ is referred to as $SSTR$.

- Two useful terms: $MSW \equiv SSW/[I(J-1)]$, $MSB = SSB/(I-1)$

**Question: What do the values of $SSW$ and $SSB$ tell us?** We will see that the idea underlying ANOVA is the comparison of the size of various sums of squares. For example, what does the following tell us?

(1) Large within group variance and small between group variance

(2) Small within group variance and large between group variance.

draw two boxplots here

Intuitively, large within group variance and small between group variance implies small treatment difference; small within group and large between group variance indicates large treatment effects. To see how these sums of squares are related to our hypothesis testing, we need to derive some theoretical results.

**Theorem A** Under the assumptions of the one-way ANOVA model,

$$E(SSW) = I(J-1)\sigma^2$$

$$E(SSB) = J\sum_{i=1}^{I}\alpha_i^2 + (I-1)\sigma^2$$

Proof:

$$SSW = \sum_{i=1}^{I}\sum_{j=1}^{J}(Y_{ij} - \bar{Y}_{i\cdot})^2$$

$$= \sum_{i=1}^{I}(J-1)S_i^2$$

$$= (J-1)\sum_{i=1}^{I}S_i^2$$

Because the sample variance $S_i^2$ is an unbiased estimator for $\sigma^2$, i.e., $E(S_i^2) = \sigma^2$, $E(SSW) = (J-1)\sum_{i=1}^{I}E(S_i^2) = I(J-1)\sigma^2$.

An unbiased estimate of $\sigma^2$ is

$$MSW = \frac{SSW}{I(J-1)}$$

MSW is also called the mean squared error and denoted by $MSE$.

To prove the second part of the theorem, we introduce **Lemma A**

**Lemma A** Let $X_i$, where $i = 1, \cdots, n$, be independent random variables with $E(X_i) = \mu_i$ and $Var(X_i) = \sigma^2$. Then

$$E(X_i - \bar{X})^2 = (\mu_i - \bar{\mu})^2 + \frac{n-1}{n}\sigma^2$$

where

$$\bar{\mu} = \frac{1}{n}\sum_{i=1}^{n}\mu_i$$

Proof of Lemma A:

First, for any random variable $U$ with finite variance, by the definition of variance we have $E(U^2) = [E(U)]^2 + Var(U)$. Let $U = X_i - \bar{X}$, then

$$E(U^2) = E(X_i - \bar{X})^2$$
$$[E(U)]^2 = (\mu_i - \bar{\mu})^2$$

$$
\begin{aligned}
Var(U) = Var(X_i - \bar{X}) &= Var(X_i) + Var(\bar{X}) - 2Cov(X_i, \bar{X}) \\
&= \sigma^2 + \frac{1}{n}\sigma^2 - 2Cov(X_i, \frac{1}{n}\sum_{j=1}^{n} X_j) \\
&= \sigma^2 + \frac{1}{n}\sigma^2 - \frac{2}{n}\sigma^2 \\
&= \frac{n-1}{n}\sigma^2
\end{aligned}
$$

The above identities indicate the lemma.

*Back to the theorem ...*

Now, note that $\bar{Y}_{i.} \sim N(\mu + \alpha_i, \sigma^2/J)$, with the role of $X_i$ being played by $\bar{Y}_{i.}$ and that of $\bar{X}$ being played by $\bar{Y}_{..}$, we have

$$
\begin{aligned}
SSB &= J\sum_{i=1}^{J}(\bar{Y}_{i.} - \bar{Y}_{..})^2 \\
E(SSB) &= J\sum_{i=1}^{I}(\bar{Y}_{i.} - \bar{Y}_{..})^2 = J\sum_{i=1}^{I}[\frac{I-1}{I}\frac{\sigma^2}{J} + (\mu + \alpha_i - \mu)^2] \\
&= J\sum_{i=1}^{I}[\alpha_i^2 + \frac{I-1}{I}\sigma^2/J] \\
&= J\sum_{i=1}^{I}\alpha_i^2 + (I-1)\sigma^2
\end{aligned}
$$

Similar to $MSW$, we define $MSB = SSB/(I-1)$. Under the null, i.e., all $\alpha_i$'s are zero, $E(MSW) = E(MSB) = \sigma^2$. If some $\alpha_i$ are nonzero, then $MSB$ will be inflated. Using $MSW$ and $MSB$, theorem A can be restated as

$$
\begin{aligned}
E(MSW) &= \sigma^2 \\
E(MSB) &= \frac{J}{I-1}\sum_{i=1}^{I}\alpha_i^2 + \sigma^2
\end{aligned}
$$

If the null hypothesis is true, i.e., $\alpha_1 = \alpha_2 = \cdots = \alpha_I = 0$, the two mean squares have the same expectation, i.e., $E(MSB) = E(MSW)$. Under the alternative, we have $E(MSB) > E(MSW)$. This motivates us to construct a test statistic based on $MSB$ and $MSE$. In fact, the ratio of them is called the $F$ statistic.

**Note:** We can use **Lemma A** to calculate $E(SSTO)$

$$
\begin{aligned}
E(SSTO) &= \sum_{i=1}^{I}\sum_{j=1}^{J} E[(Y_{ij} - \bar{Y}_{..})^2] \\
&= \sum_{i=1}^{I}\sum_{j=1}^{J}[(\mu + \alpha_i - \mu)^2 + \frac{IJ-1}{IJ}\sigma^2] \\
&= J\sum_{i=1}^{I}\alpha_i^2 + (IJ-1)\sigma^2
\end{aligned}
$$

This result agrees with $E(SSW)$, $E(SSB)$, and the decomposition of $SSTO$.

### 1.2.3 The F-test

Distributions of sum of squares under the one-way ANOVA model: (**Theorem B**)

- $SSW/\sigma^2 \sim \chi_{I(J-1)}$ (part I of Theorem B)

- If the null is true, $SSB/\sigma^2 \sim \chi_{I-1}$ (part II of Theorem B)

- $SSW$ and $SSB$ are independent (part III of Theorem B)

**Proof of B.1**

$$
\begin{aligned}
SSW/\sigma^2 &= \frac{1}{\sigma^2}\sum_{i=1}^{I}\sum_{j=1}^{J}(Y_{ij} - \bar{Y}_{i.})^2 \\
&= \sum_{i=1}^{I}\frac{1}{\sigma^2}\sum_{j=1}^{J}(Y_{ij} - \bar{Y}_{i.})^2 \\
&= \sum_{i=1}^{I}\frac{(J-1)S_i^2}{\sigma^2}
\end{aligned}
$$

Here $S_i^2$ is the sample variance from the $i$th sample. From 120B, we know that

$$
\frac{(J-1)S_i^2}{\sigma^2} \sim \chi_{J-1}^2
$$

. Because the samples are independent, the sample variances $S_i^2$'s are independent; therefore, $SSW/\sigma^2 \sim \chi^2_{I(J-1)}$.

Note

(1) Part I implies that $E(SSW) = I(J-1)\sigma^2$.

(2) $S_p^2 = MSW = \frac{SSW}{I(J-1)}$ is called the pooled sample variance.

**Proof of B.2**

Consider the sample with the sample means:

$$\{\bar{Y}_{1\cdot}, \bar{Y}_{2\cdot}, \cdots, \bar{Y}_{I\cdot}\}$$

Since (1) each of them is a linear combination of independent normal random variables and (2) they are calculated from independent samples, they are independent normal random variables. In fact, $\bar{Y}_{i\cdot} \overset{independent}{\sim} N(\mu + \alpha_i, \sigma^2/J)$. Under the null hypothesis, they have the same population mean and variance. Thus, we can treat

$$\{\bar{Y}_{1\cdot}, \bar{Y}_{2\cdot}, \cdots, \bar{Y}_{I\cdot}\}$$

as a random sample from $N(\mu, \sigma^2/I)$ when the null hypothesis is true.

The corresponding sample mean and sample variances are $\bar{Y}_{\cdot\cdot} = \frac{\bar{Y}_{1\cdot}+\cdots+\bar{Y}_{I\cdot}}{I}$ and $S_{TR}^2 = \frac{1}{I-1}\sum_{i=1}^{I}(\bar{Y}_{i\cdot}-\bar{Y}_{\cdot\cdot})^2$, respectively. Based on the properties of sample variance (see stat120B), we have

$$\frac{(I-1)S_{TR}^2}{\sigma^2/J} \overset{H_0}{\sim} \chi^2_{I-1}$$

Now return to $SSB$.

$$
\begin{aligned}
SSB/\sigma^2 &= \frac{J}{\sigma^2}\sum_{i=1}^{I}(\bar{Y}_{i\cdot} - \bar{Y}_{\cdot\cdot})^2 \\
&= \frac{(I-1)S_{TR}^2}{\sigma^2/J} \overset{H_0}{\sim} \chi^2_{I-1}
\end{aligned}
$$

**Proof of B.3**

$SSW$ is a function of $S_i^2, i = 1, \cdots, I$, where $S_i^2 = \frac{1}{J-1}\sum_{j=1}^{J}(Y_{ij} - \bar{Y}_{i\cdot})^2$.

$SSB$ is a function of $\bar{Y}_{i\cdot}$'s ($\bar{Y}_{\cdot\cdot}$ is also a function of $\bar{Y}_{i\cdot}$'s).

We claim that $\bar{Y}_{1\cdot}, \cdots, \bar{Y}_{I\cdot}$ and $S_1^2, \cdot, S_I^2$ are independent with each other.

When $i \neq i'$, $S_i^2$ and $\bar{Y}_{i'\cdot}$ are independent because they are functions of different observations.

9

When $i = i'$, by 120statB, $S_i^2$ and $\bar{Y}_{i'.}$ are independent. (for a normal random sample, the sample mean and sample variance are independent)

Since $SSB$ is a function of $\bar{Y}_{1.}, \cdots, \bar{Y}_{I.}$, and $SSW$ is a function of $S_1^2, \cdots, S_I^2$, $SSB$ and $SSW$ are independent from each other.

**Summary of the proof:**

to prove B.1, we consider samples $\{Y_{i1}, Y_{i2}, \cdots, Y_{iJ}\}$ for each $i$;

to prove B.2, we consider the sample of means $\{\bar{Y}_1, \bar{Y}_2, \cdots, \bar{Y}_I\}$.

Finally, we have all the elements for **Theorem C**

**Theorem C: The F-test for the One-Way ANOVA** Assume that the assumptions of one-way ANOVA hold (normality, independence, and equal variance), when the null hypothesis $H_0 : \alpha_2 = \alpha_2 = \cdots = \alpha_I = 0$ (or $H_0 : \mu_1 = \mu_2 = \cdots = \mu_I$) is true,

$$F \equiv \frac{MSB}{MSW} = \frac{SSB/(I-1)}{SSW/[I(J-1)]} \sim F_{I-1, I(J-1)}$$

Proof:

$$
\begin{aligned}
F &= \frac{MSB}{MSW} = \frac{SSB/(I-1)}{SSW/[I(J-1)]} \\
&= \frac{\frac{SSB}{\sigma^2}/[I-1]}{\frac{SSW}{\sigma^2}/[I(J-1)]}
\end{aligned}
$$

Since $\frac{SSW}{\sigma^2} \sim \chi_{I(J-1)}$ and $\frac{SSB}{\sigma^2} \overset{H_0}{\sim} \chi_{I-1}$ and $SSW$ and $SSB$ are independent, $F \overset{H_0}{\sim} F_{I-1, I(J-1)}$.

We reject $H_0$ at a significance level $\alpha$ if the test statistic $F$ is greater than $F_{I-1, I(J-1), 1-\alpha}$. Here $F_{I-1, I(J-1), 1-\alpha}$ is the upper $\alpha$ point of $F_{I-1, I(J-1)}$.

The ANOVA table:

| Source | SS | df | MS | F |
|---|---|---|---|---|
| Treatment | $SSB = \sum_i^I J(\bar{Y}_{i.} - \bar{Y}_{..})^2$ | $I-1$ | $MSB = \frac{SSB}{I-1}$ | $MSB/MSW$ |
| Error | $SSW = \sum_{i=1}^I \sum_{j=1}^J (Y_{ij} - \bar{Y}_{i.})^2$ | $I(J-1)$ | $MSW = SSW/[I(J-1)]$ | |
| Total | $SSTO = \sum_{i=1}^I \sum_{j=1}^J (Y_{ij} - \bar{Y}_{..})^2$ | $IJ-1$ | | |

The ANOVA table for the example

| Source | $SS$ | $df$ | $MS$ | $F$ |
|--------|------|------|------|-----|
| Labs | .125 | 6 | .021 | 5.66 |
| Error | .231 | 63 | .0037 | |
| Total | .356 | 69 | | |

$F = 5.66 > F_{6,63,0.95} = 2.246408$. And the p-value is much smaller than 0.01. So we conclude that the levels of the chemical compound are different across different labs.

Another useful result based on **Theorem B**

$$\frac{(\bar{Y}_{i_1\cdot} - \bar{Y}_{i_2\cdot}) - (\alpha_{i_1} - \alpha_{i_2})}{\sqrt{S_p^2(\frac{1}{J} + \frac{1}{J})}} \sim t_{I(J-1)}$$

for $i_1 \neq i_2$. Here $S_p^2 = MSW$.

Proof: $\bar{Y}_{i_1\cdot} \sim N(\mu + \alpha_{i_1}, \sigma^2/J)$, $\bar{Y}_{i_2\cdot} \sim N(\mu + \alpha_{i_2}, \sigma^2/J)$.

For different $i_1, i_2$, $\bar{Y}_{i_1\cdot}$ and $\bar{Y}_{i_2\cdot}$ are also independent. Therefore

$$\frac{(\bar{Y}_{i_1\cdot} - \bar{Y}_{i_2\cdot}) - (\alpha_{i_1} - \alpha_{i_2})}{\sqrt{(\frac{1}{J} + \frac{1}{J})\sigma^2}} \sim N(0, 1)$$

Based on B.2, we have $SSW/\sigma^2 = I(J-1)MSW/\sigma^2 = I(J-1)S_p^2/\sigma^2 \sim \chi^2_{I(J-1)}$.

Also, in the proof of B.3, we have shown that the two vectors are independent. So $\bar{Y}_{i_1\cdot} - \bar{Y}_{i_2\cdot}$ are independent of $S_p^2$.

Based on the above facts, we have

$$\frac{[(\bar{Y}_{i_1\cdot} - \bar{Y}_{i_2\cdot}) - (\alpha_{i_1} - \alpha_{i_2})]/[\sqrt{(\frac{1}{J} + \frac{1}{J})\sigma^2}]}{\sqrt{I(J-1)S_p^2/[\sigma^2 I(J-1)]}} \sim t_{I(J-1)}$$

Simplify the left hand side,

$$\frac{(\bar{Y}_{i_1\cdot} - \bar{Y}_{i_2\cdot}) - (\alpha_{i_1} - \alpha_{i_2})}{\sqrt{S_p^2(\frac{1}{J} + \frac{1}{J})}} \sim t_{I(J-1)}$$

**A special case:**

When $I = 2$, the two-sample t-test statistic $t \sim t_{J-1}$. In the homework you have shown that $F = t^2 \sim F_{1,J-1}$. This agrees with the fact that $Z \sim t_n \Rightarrow Z^2 \sim F_{1,n}$.

11

## 1.3   The $F$ test under unbalanced designs

The test for unbalanced designs is very similar - just replacing $J$ with $J_i$. In this case,

$$SSW = \sum_{i=1}^{I} \sum_{j=1}^{J_i} (Y_{ij} - \bar{Y}_{i\cdot})^2$$

$$SSB = \sum_{i=1}^{I} J_i (\bar{Y}_{i\cdot} - \bar{Y}_{\cdot\cdot})^2$$

$$F = \frac{MSB}{MSW} = \frac{SSB/(I-1)}{SSW/[\sum_i (J_i - 1)]} \sim F_{I-1, \sum_i (J_i - 1)}$$

| Source | SS | df | MS | F |
|---|---|---|---|---|
| Treatment | $SSB = \sum_i^I J_i (\bar{Y}_{i\cdot} - \bar{Y}_{\cdot\cdot})^2$ | $I-1$ | $MSB = \frac{SSB}{I-1}$ | $MSB/MS$ |
| Error | $SSW = \sum_{i=1}^{I} \sum_{j=1}^{J_i} (Y_{ij} - \bar{Y}_{i\cdot})^2$ | $\sum_{i=1}^{I} (J_i - 1)$ | $MSW = SSW/\sum_{i=1}^{I} (J_i - 1)$ | |
| Total | $SSTO = \sum_{i=1}^{I} \sum_{j=1}^{J_i} (Y_{ij} - \bar{Y}_{\cdot\cdot})^2$ | $\sum_{i=1}^{I} J_i - 1$ | | |

Review about the assumptions we used to construct the $F$ test:

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

(2) equal variance

(3) independent $\epsilon_{ij}$.

For large samples, assumption (1) is not very important. If the all groups have the same number of observations, violation of (2) does not have a strong impact on results. Assumption (3) is very challenging.

## 1.4   multiple comparisons

In one-way ANOVA we were concerned about whether all the means are the same. If the null hypothesis of equal mean is rejected, we still have no idea about how the means differ. Usually it is interesting to ask whether pairs of a subgroup of treatments show any difference. A naive approach would be to compare all pairs of treatment means using t-tests.

Here is the problem:

When considering a pair, with a significance level $\alpha$, we have a type I error rate $\alpha$. But this is for a pair.

If there there are multiple comparisons, the overall type I error across all comparisons will be inflated.

For instance, suppose that we perform $K$ independent level $\alpha$ tests and that in each case the null hypothesis is true. The overall type I error rate is defined as

$$
\begin{aligned}
&= \quad Pr(\text{reject at least one } H_0 | H_0 \text{ true on all tests}) \\
&= \quad 1 - Pr(\text{no } H_0 \text{ was rejected } | H_0 \text{ true on all tests}) \\
&= \quad 1 - (1 - \alpha)^K
\end{aligned}
$$

For $\alpha = 0.05$ and $K = 10$, the overall type I error rate is $1 - 0.95^{10} = 0.40$.

It is often desirable to control the overall type I error rate at a given level, say 0.05. How to do that? There are many options. Here we focus on a popular method: the Bonferroni correction.

**The Bonferroni Correction**

The Bonferroni Correction comes from **the Bonferroni inequality** (also known as Bool's inequality) , which states that:

$$
P(\cup_{i=1}^{K} A_i) \leq \sum_{i=1}^{K} Pr(A_i)
$$

for event $A_1, A_2, \cdots, A_K$.

When $K = 2$, it is easy to verify that: From STAT120A, we have $Pr[A_1 \cup A_2] = Pr[A_1] + Pr[A_2] - Pr[A_1 \cap A_2]$. But $P[A_1 \cap A_2] \geq 0$. Therefore,

$$
Pr[A_1 \cup A_2] \leq Pr[A_1] + Pr[A_2] - Pr[A_1 \cap A_2].
$$

For a general K, we can use mathematical induction to prove the inequality. (homework assignment)

For test $i$, we define

$A_i = \{\text{reject } H_i | \text{all null are true}\}$

Suppose that $\alpha^* = Pr(A_i) = Pr(\text{reject the ith null} | \text{the ith null is true})$.

According to the Bonferrroni inequality,

$$
\begin{aligned}
\text{overall type I error rate} \quad &= \quad Pr(\text{there is at least one positive --- all null hypotheses are true}) \\
&= \quad Pr(\cup_{i=1}^{K} A_i) \\
&\leq \quad \sum_{i=1}^{K} Pr(A_i) \\
&= \quad K\alpha^*
\end{aligned}
$$

Therefore, the overall type I error rate is $\le K\alpha*$. This implies that if we choose $\alpha* = \alpha/K$, the overall type I error rate is no greater than $\alpha$. For example, if we have 10 tests and we want to control the overall type I error rate $\le 0.05$, we reject each null hypothesis at a significance level of 0.005.
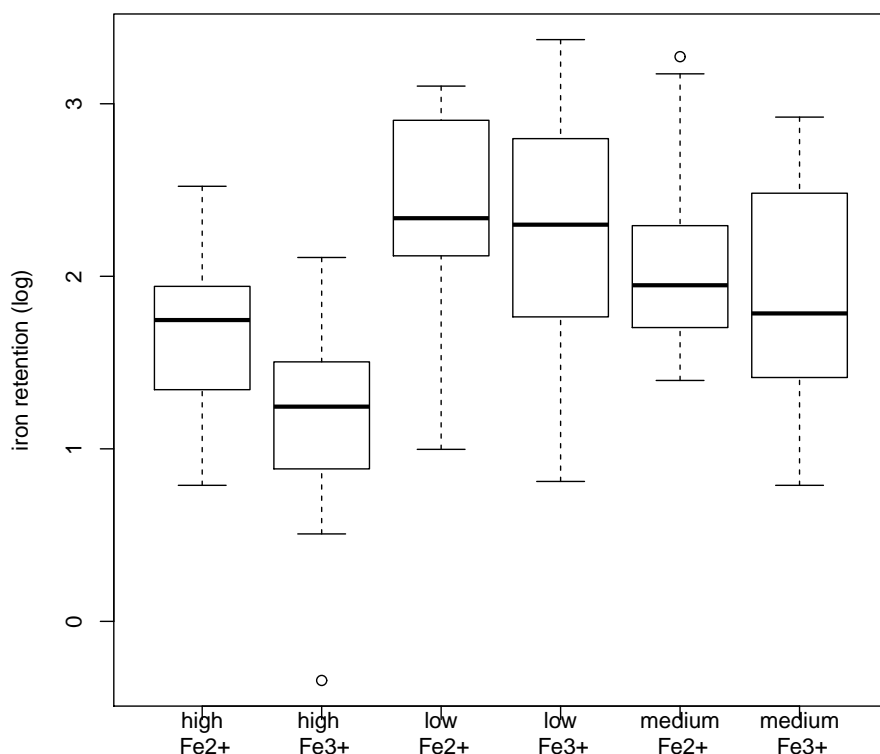
Recall that based on **Theorem B** we have

$$\frac{(\bar{Y}_{i_1 \cdot} - \bar{Y}_{i_2 \cdot}) - (\alpha_{i_1} - \alpha_{i_2})}{\sqrt{S_p^2(\frac{1}{J} + \frac{1}{J})}} \sim t_{I(J-1)}$$

Definition. A set of <u>simultaneous 95% confidence intervals</u> for the pairwise comparison is

$$(\bar{Y}_{i_1 \cdot} - \bar{Y}_{i_2 \cdot}) \pm \sqrt{S_p^2 \frac{2}{J}} t_{I(J-1), 1-0.05/(2K)}$$

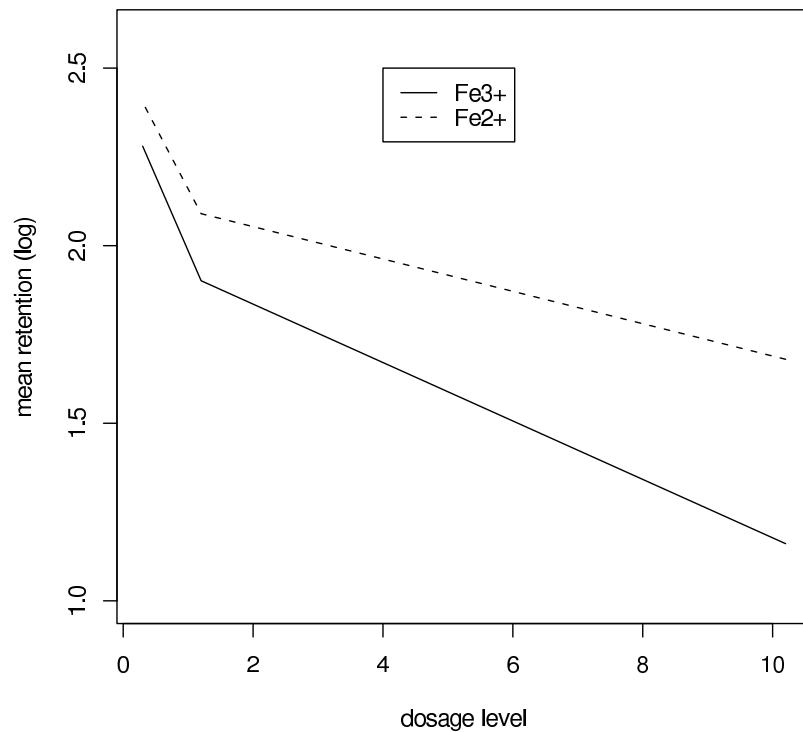Example. If I=7, we have $K = \binom{7}{2} = 21$ pairwise comparisons.

## 2 Two-way ANOVA

In the one-way design there is only one factor. What if there are several factors? Often, we are interested to know the simultaneous effects of multiple factors, e.g, gender and smoking on hypertension. The statistical approach to analyze data from a two-way design is two-way ANOVA.

Iron retention. An experiment was performed to determine the factors that affect iron retention. The experiment was done on 108 mice, which were randomly divided into 6 groups of 18 each. These groups consist of the six combinations of two forms of iron and three concentrations.

The table below shows the mean logged retention percentage.

| Form | low | medium | high |
|------|-----|--------|------|
| Fe3+ | 2.28 ($\bar{Y}_{11.}$) | 1.90 ($\bar{Y}_{12.}$) | 1.16 ($\bar{Y}_{13.}$) |
| Fe2+ | 2.40 ($\bar{Y}_{21.}$) | 2.09 ($\bar{Y}_{22.}$) | 1.68 ($\bar{Y}_{23.}$) |

The figure below shows the mean retention levels (log).

Several scientific questions:

(1) Are the means different by iron form?

(2) Are the means different by dosage?

(3) Does difference in the means between the two iron forms change as a function of dosage?

The first two questions can be answered by one-way ANOVA. Here we use the two-way anova to answer all the questions simultaneously.

## 2.1 The two-way layout

A two-way layout is an experiment design involving two factors, each at two or more levels. If there are I levels of factor A and J of factor B, there are $I - by - J$ combinations. We assume that $K$ independent observations are taken for each of the $IJ$ combinations.

Let $Y_{ijk}$ denote the $kth$ observation in cell $ij$. The statistical model is a natural extension of the one-way anova:

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \delta_{ij} + \epsilon_{ijk}, \text{ with } \epsilon_{ijk} \sim N(0, \sigma^2)$$

with the constraints

$$
\begin{aligned}
\sum_{i=1}^{I} \alpha_i &= 0 \\
\sum_{j=1}^{J} \beta_i &= 0 \\
\sum_{i=1}^{I} \delta_{ij} &= \sum_{j=1}^{J} \delta_{ij} = 0
\end{aligned}
$$

Interpretations of the parameters:

$\mu$ is the grand mean over all cells

$\alpha_i$ is the main effect of the $ith$ level of factor A

$\beta_j$ is the main effect of the $jth$ level of factor B

$\delta_{ij}$ is the interaction effect between the $ith$ level of factor A and the $jth$ level of factor B, i.e., it denotes the effect of combining the $ith$ level of factor A and the $jth$ level of factor B.

**interpretation**

Under the two-way layout,

the mean response in the $ith$ level of factor A is $\mu + \alpha_i$

the mean response in the $jth$ level of factor B is $\mu + \beta_j$

the mean response in the cell $ij$ is $\mu + \alpha_i + \beta_j + \delta_{ij}$

We say the two factors have an interaction effect on response if any of $\delta_{ij}$ is not zero.

When there is no interaction, the means are

17

| Form | low | medium | high |
|---|---|---|---|
| Fe2+ | $\mu + \alpha_1 + \beta_1$ | $\mu + \alpha_1 + \beta_2$ | $\mu + \alpha_1 + \beta_3$ |
| Fe3+ | $\mu + \alpha_2 + \beta_1$ | $\mu + \alpha_2 + \beta_2$ | $\mu + \alpha_2 + \beta_3$ |
| Difference | $\alpha_1 - \alpha_2$ | $\alpha_1 - \alpha_2$ | $\alpha_1 - \alpha_2$ |

The graph of the means would look like

## 2.2 Two-way ANOVA

When there is an interaction between iron form and dosage, the means are

| Form | low | medium | high |
|---|---|---|---|
| Fe3+ | $\mu + \alpha_1 + \beta_1 + \delta_{11}$ | $\mu + \alpha_1 + \beta_2 + \delta_{12}$ | $\mu + \alpha_1 + \beta_3 + \delta_{13}$ |
| Fe2+ | $\mu + \alpha_2 + \beta_1 + \delta_{12}$ | $\mu + \alpha_2 + \beta_2 + \delta_{22}$ | $\mu + \alpha_2 + \beta_3 + \delta_{23}$ |
| Difference | $\alpha_1 - \alpha_2 + (\delta_{11} - \delta_{21})$ | $\alpha_1 - \alpha_2 + (\delta_{12} - \delta_{22}))$ | $\alpha_1 - \alpha_2 + (\delta_{13} - \delta_{23})$ |

The graph of the means would look like

### 2.2.1 Notations

$$\bar{Y}_{...} = \sum_{i=1}^{I}\sum_{j=1}^{J}\sum_{k=1}^{K} Y_{ijk}/(IJK)$$

$$\bar{Y}_{i..} = \sum_{j=1}^{J}\sum_{k=1}^{K} Y_{ijk}/(JK)$$

$$\bar{Y}_{\cdot j \cdot} = \sum_{i=1}^{I}\sum_{k=1}^{K} Y_{ijk}/(IK)$$

$$\bar{Y}_{ij\cdot} = \sum_{k=1}^{K} Y_{ijk}/K$$

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ij} \ , \ \epsilon_{ijk} \overset{iid}{\sim} N(0, \sigma^2) \ , \ \sum_i \alpha_i = \sum_j \beta_j = \sum_i \gamma_{ij} = \sum_j \gamma_{ij} = 0$$

Assume that the observations are independent and normally distributed with equal variance, the log likelihood is

$$
\begin{aligned}
l &= -\frac{IJK}{2} log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^{I} \sum_{j=1}^{J} \sum_{k=1}^{K} (Y_{ijk} - \mu - \alpha_i - \beta_j - \delta_{ij})^2 \\
&= -\frac{IJK}{2} log(2\pi) - \frac{IJK}{2} log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^{I} \sum_{j=1}^{J} \sum_{k=1}^{K} (Y_{ijk} - \mu - \alpha_i - \beta_j - \delta_{ij})^2
\end{aligned}
$$

The mles are

$$
\begin{aligned}
\hat{\mu} &= \bar{Y}_{...} \\
\hat{\alpha}_i &= \bar{Y}_{i..} - \bar{Y}_{...}, \ i = 1, \cdots, I \\
\hat{\beta}_j &= \bar{Y}_{.j.} - \bar{Y}_{...}, \ j = 1, \cdots, J \\
\hat{\delta}_{ij} &= \bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{.j.} + \bar{Y}_{...}
\end{aligned}
$$

Similar to the definition used in the one-way ANOVA, the sum of squares of total is defined as

$$
\begin{aligned}
SSTO &= \sum_{i=1}^{I} \sum_{j=1}^{J} \sum_{k=1}^{K} (Y_{ijk} - \bar{Y}_{...})^2 \\
&= \sum_{i=1}^{I} \sum_{j=1}^{J} \sum_{k=1}^{K} (Y_{ijk} - \hat{\mu})^2
\end{aligned}
$$

$SSTO$ can be decomposed to

$$SSTO = SSA + SSB + SSAB + SSE$$

where

19

$$
\begin{aligned}
SSA &= JK\sum_{i=1}^{I}(\bar{Y}_{i\cdot\cdot} - \bar{Y}_{\cdots})^2 \\
&= JK\sum_{i=1}^{I}\hat{\alpha}_i^2 \\
SSB &= IK\sum_{j=1}^{J}(\bar{Y}_{\cdot j\cdot} - \bar{Y}_{\cdots})^2 \\
&= IK\sum_{j=1}^{J}\hat{\beta}_j^2 \\
SSAB &= K\sum_{i=1}^{I}\sum_{j=1}^{J}(\bar{Y}_{ij\cdot} - \bar{Y}_{i\cdot\cdot} - \bar{Y}_{\cdot j\cdot} + \bar{Y}_{\cdots})^2 \\
&= K\sum_{i=1}^{I}\sum_{j=1}^{J}\hat{\delta}_{ij}^2 \\
SSE &= \sum_{i=1}^{I}\sum_{j=1}^{J}\sum_{k=1}^{K}(Y_{ijk} - \bar{Y}_{ij\cdot})^2 = (K-1)\sum_{i}\sum_{j}S_{ij}^2 \\
&= \sum_{i=1}^{I}\sum_{j=1}^{J}\sum_{k=1}^{K}(Y_{ijk} - (\hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j + \hat{\delta}_{ij}))^2
\end{aligned}
$$

The proof is also similar to that of one-way ANOVA. Essentially, we need to use the identity

$$
Y_{ijk} - \bar{Y}_{\cdots} = (Y_{ijk} - \bar{Y}_{ij\cdot}) + (\bar{Y}_{i\cdot\cdot} - \bar{Y}_{\cdots}) + (\bar{Y}_{\cdot j\cdot} - \bar{Y}_{\cdots}) + (\bar{Y}_{ij\cdot} - \bar{Y}_{i\cdot\cdot} - \bar{Y}_{\cdot j\cdot} + \bar{Y}_{\cdots})
$$

**Theorem A: Expectations of Sums of Squares** Under the two-way ANOVA model,

$$
\begin{aligned}
(1)\ E(MSE) &= E(SSE/[IJ(K-1)]) = \sigma^2 \\[2mm]
(2)\ E(MSA) &= E(SSA/(I-1)) = \sigma^2 + \frac{JK}{I-1}\sum_{i=1}^{I}\alpha_i^2 \\[2mm]
(3)\ E(MSB) &= E(SSB/(J-1)) = \sigma^2 + \frac{IK}{J-1}\sum_{j=1}^{J}\beta_i^2 \\[2mm]
(4)\ E(MSAB) &= E(SSAB/[(I-1)(J-1)]) = \sigma^2 + \frac{K}{(I-1)(J-1)}\sum_{i=1}^{I}\sum_{j=1}^{J}\delta_{ij}^2
\end{aligned}
$$

Proof: They can be proved by using Lemma A:

Let $X_i$, where $i = 1, \cdots, n$ be independent random variables with $E(X_i) = \mu_i$ and $Var(X_i) = \sigma^2$. Then

$$
E(X_i - \bar{X})^2 = (\mu_i - \bar{\mu})^2 + \frac{n-1}{n}\sigma^2
$$

Apply it to (1): $E(Y_{ijk} - \bar{Y}_{ij.})^2 = (0-0)^2 + (K-1)/K\sigma^2$, thus

$$
E(SSE) = E[\sum_{ij}(K-1)[S_{ij}^2] = IJKE[S_{ij}^2] = IJ(K-1)
$$

Apply it (2):

Apply it to (3):

Notice that $Y_{ijk} \sim N(\mu + \alpha_i + \beta_j + \delta_{ij}, \sigma^2)$ and apply the lemma to $SSTO$,

$$
\begin{aligned}
E(SSTO) &= \sum_{i=1}^{I}\sum_{j=1}^{J}\sum_{k=1}^{K}E(Y_{ijk} - \bar{Y}_{...})^2 \\[2mm]
&= \sum_{i=1}^{I}\sum_{j=1}^{J}\sum_{k=1}^{K}[(\mu + \alpha_i + \beta_j + \delta_{ij} - \mu)^2 + \frac{IJK-1}{IJK}\sigma^2] \\[2mm]
&= (IJK-1)\sigma^2 + \sum_{i=1}^{I}\sum_{j=1}^{J}\sum_{k=1}^{K}(\alpha_i + \beta_j + \delta_{ij})^2 \\[2mm]
&= (IJK-1)\sigma^2 + JK\sum_{i=1}^{I}\alpha_i^2 + IK\sum_{j=1}^{J}\beta_j^2 + K\sum_{i=1}^{I}\sum_{j=1}^{J}\delta_{ij}^2
\end{aligned}
$$

21

The last step is true because of the constraints on the parameters. For example,

$$\sum_{i=1}^{I}\sum_{j=1}^{J}\sum_{k=1}^{K}\alpha_i\beta_j = K(\sum_{i=1}^{I}\alpha_i)(\sum_{j=1}^{J}\beta_j) = 0$$

Based on (1)-(3) and $E(SSTO)$, we can prove (4).

**Theorem B: Distributions of Sums of Squares** Assume that the errors are independent and normally distributed with means zero and variances $\sigma^2$. Then

(1) $SSE/\sigma^2$ follows a chi-squared distribution with $IJ(K-1)$ degrees of freedom.

(2) Under the null

$$H_0 : \alpha_i = 0, i = 1, \cdots, I$$

$SSA/\sigma^2$ follows a chi-squared distribution with $I-1$ degrees of freedom.

(3) Under the null

$$H_0 : \beta_j = 0, j = 1, \cdots, J$$

$SSB/\sigma^2$ follows a chi-squared distribution with $J-1$ degrees of freedom.

(4) Under the null

$$H_0 : \delta_{ij} = 0, i = 1, \cdots, I, j = 1, \cdots, J$$

$SSAB/\sigma^2$ follows a chi-squared distribution with $(I-1)(J-1)$ degrees of freedom.

(5) The sums of squares are independently distributed.

Proof.

We only provide the proof for (1).

$$SSE/\sigma^2 = \frac{1}{\sigma^2}\sum_{i=1}^{I}\sum_{j=1}^{J}\sum_{k=1}^{K}(Y_{ijk} - \bar{Y}_{ij\cdot})^2$$

Consider the sample in group $(i, j)$:

$$\{Y_{ij1}, Y_{ij2}, \cdots, Y_{ijK}\}.$$

Here we fixed the level of factor A at $i$ and the level of factor $B$ at $j$. All the observations have the same distribution: $N(\mu + \alpha_i + \beta_j + \delta_{ij}, \sigma^2$ and they are independent. This is becuase we assumed that $Y_{ijk} = \mu + \alpha_i + \beta_j + \delta_{ij} + \epsilon_{ijk}$ where $\epsilon_{ijk} \overset{iid}{\sim} N(0, \sigma^2)$. So

$$\{Y_{ij1}, Y_{ij2}, \cdots, Y_{ijK}\}$$

is a random sample from $N(\mu + \alpha_i + \beta_j + \delta_{ij}, \sigma^2)$. In 120B we learned that the distribution of the sample variance of a random sample from a Normal distribution. Let $S_{ij}^2$ denote the sample variance of the random sample, we have

$$\frac{1}{\sigma^2} \sum_{k=1}^{K} (Y_{ijk} - \bar{Y}_{ij.})^2 = \frac{K-1}{\sigma^2} S_{ij}^2 \sim \chi_{K-1}^2.$$

Since $S_{ij}^2$'s are from independent samples, they are independent. In addition, it is not difficult to verfiy that

$$\frac{SSE}{\sigma^2} = \sum_i \sum_j \frac{(K-1)S_{ij}^2}{\sigma^2}$$

Therefore,

$$SSE/\sigma^2 \sim \chi_{IJ(K-1)}^2$$

Proofs for the (2,3) are similar. Proofs for (4,5) are not required in this course.

For (2) Think about the sample variance of this random sample under the null

$$\{\bar{Y}_{1..}, \bar{Y}_{2..}, \cdots, \bar{Y}_{I..}\}$$

For (3) Think about the sample variance of this random sample under the null

$$\{\bar{Y}_{.1.}, \bar{Y}_{.2.}, \cdots, \bar{Y}_{.I.}\}$$

For (4) The rationale for $df = (I-1)(J-1)$: there are $IJ$ $\delta_{ij}$'s but they are redundant: given the constraints $\sum_i = \sum_j = 0$ the last row and the last columns are not needed, which results in $(I-1)(J-1)$ unique parameters.

Implication of very large SSA. what does a large $SSA$ tell us? Under the null

$$H_A : \alpha_i = 0, i = 1, \cdot, I$$

$$E(MSA) = E(MSE)$$

If the null hypothesis is not true, $E[MSA] > E[MSE]$. Therefore,we can use $MSA/MSE$ as a test statistic. In fact, under the null, it follows $F_{I-1,IJ(K-1)}$. Summary,

(1) Under the null

$$H_A : \alpha_i = 0, i = 1, \cdot, I$$

$$F = \frac{MSA}{MSE} = \frac{SSA/(I-1)}{SSE/[IJ(K-1)]}$$

$$= \frac{SSA/[\sigma^2(I-1)]}{SSE/[\sigma^2 IJ(K-1)]}$$

$$=_d \frac{\chi^2_{I-1}/(I-1)}{\chi^2_{IJ(K-1)}/[IJ(K-1)]}$$

$$\sim F_{I-1,IJ(K-1)}$$

(2) Under the null

$$H_A : \beta_j = 0, j = 1, \cdot, J$$

$$F = \frac{MSB}{MSE} \sim F_{J-1,IJ(K-1)}$$

(3) Under the null

$$H_A : \delta_{ij} = 0, i = 1, \cdots, I, j = 1, \cdots, J$$

$$F = \frac{MSAB}{MSE} \sim F_{(I-1)(J-1),IJ(K-1)}$$

All the distributions have the same denominator degrees of freedom. This is because all the sums of squares are compared to the mean sum of squares of error.

ANOVA table for two-way models

| Source | $SS$ | $df$ | $MS$ | $F$ |
|--------|------|------|------|-----|
| A | $SSA$ | $I-1$ | $MSB = \frac{SSA}{I-1}$ | $\frac{MSA}{MSE}$ |
| B | $SSB$ | $J-1$ | $MSB = \frac{SSB}{J-1}$ | $\frac{MSB}{MSE}$ |
| AB | $SSAB$ | $(I-1)(J-1)$ | $MSAB = \frac{SSAB}{(I-1)(J-1)}$ | $\frac{MSAB}{MSE}$ |
| Error | $SSE$ | $IJ(K-1)$ | $MSE = \frac{SSE}{IJ(K-1)}$ | |
| Total | $SSTO$ | $IJK-1$ | | |

The ANOVA table for the iron retention example

| Source | SS | df | MS | F |
|--------|------|-----|-------|-------|
| Iron form | 2.074 | 1 | 2.074 | 5.99 |
| Dosage | 15.588 | 2 | 7.794 | 22.53 |
| Interaction | .810 | 2 | .405 | 1.17 |
| Error | 35.396 | 102 | .346 | |
| Total 53.768 | | 107 | | |

To test the effect of iron form, we test

$$H_0 : \alpha_1 = \alpha_2 = 0 \text{ v.s. } H_1 : \text{ not all } \alpha_i \text{ equal } 0$$

using the statistic

$$F = \frac{SS_{Iron}/1}{SS_E/102} = \frac{2.074}{0.346} = 5.99$$

Because the test statistic is greater than the upper 0.05 point of $F_{1,102}$ (=3.93. If use R, do "qf(0.95, 1, 102)"), we conclude that there is an effect due to iron form. Equivalently, you can calculate the corresponding P-value using R "1-pf(5.99, 1, 102)";

Similarly, dosage also has an effect on the iron retention.

Last, we want to test whether there is an interaction effect by considering the following test statistic

$$F = \frac{SS_{AB}/(I-1)(J-1)}{SS_E/IJ(K-1)} = \frac{2.074}{0.346} = 1.17.$$

It is less than the upper 0.05 point of $F_{2,102}$. So there is not enough evidence for an interaction effect between iron form and dosage on iron retention.

qf(0.95,2,102)=3.09 (critical value)

1-pf(1.17,2,102)=0.31 (p-value)


A brief review of all the tests we have discussed:


- One group: $X_1, \cdots, X_n$, $H_0 : \mu = \mu_0$. We use the one-sample t-test.

- Two groups but paired: the t-test for paired samples, which is the one-sample t-test.

- Two groups: $X_1, \cdots, X_m, Y_1, \cdots, Y_n$. We use the two-sample t-test.

- I groups: $Y_{ij}$. We use one-way ANOVA.

- Two factors. We use two-way ANOVA.