

# Theoretical Proofs Related To ANOVA

BIO210 Biostatistics

Extra Reading Material for Lecture 31

Xi Chen

School of Life Sciences

Southern University of Science and Technology

Fall 2024

## Experimental Setup

Here we are going through three main proofs regarding ANOVA we introduced during the lecture. We provide some details that we did not cover. Most of the stuff here is basically some algebraic manipulations of expressions, which is not that important and too lengthy to show during the lecture.

First, let's clarify our data. Generally, we have drawn different samples from different populations. Let's say we have  $k$  different populations. If  $k = 2$ , we are essentially dealing with the case of  $t$ -tests. In ANOVA,  $k \geq 3$ .

Now we let  $\mu_i$  and  $\sigma_i^2$  be the mean and the variance, respectively, of the population  $i$ , where  $i = 1, 2, 3, \dots, k$ . We assume all those populations follow normal distributions:

$$\text{Population } i \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

Then we draw samples from each populations. We obtain sample  $i$  of size  $n_i$  from population  $i$ , and the sample mean and the sample variance are  $\bar{x}_i$  and  $s_i^2$ , respectively. In summary, the data is like this:

Sample	1	2	3	...	$k$
Data	$x_{11}$	$x_{21}$	$x_{31}$	...	$x_{k1}$
	$x_{12}$	$x_{22}$	$x_{32}$	...	$x_{k2}$
	$x_{13}$	$x_{23}$	$x_{33}$	...	$x_{k3}$
	$x_{14}$	$x_{24}$	$x_{34}$	...	$x_{k4}$
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
Sample mean	$\bar{x}_1$	$\bar{x}_2$	$\bar{x}_3$	...	$\bar{x}_k$
Sample variance	$s_1^2$	$s_2^2$	$s_3^2$	...	$s_k^2$
Sample size	$n_1$	$n_2$	$n_3$	...	$n_k$

We also denote the total number of data points as  $n$ , that is,  $n = \sum_{i=1}^k n_i$ .

# 1 SST = SSB + SSW

*Proof.* We start with the definition of SST:

$$\begin{aligned}
 \text{SST} &= \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} [(x_{ij} - \bar{x}_i) + (\bar{x}_i - \bar{x})]^2 \\
 &= \sum_{i=1}^k \sum_{j=1}^{n_i} [(x_{ij} - \bar{x}_i)^2 + (\bar{x}_i - \bar{x})^2 + 2(x_{ij} - \bar{x}_i)(\bar{x}_i - \bar{x})] \\
 &= \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2 + \sum_{i=1}^k \sum_{j=1}^{n_i} (\bar{x}_i - \bar{x})^2 + 2 \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)(\bar{x}_i - \bar{x})
 \end{aligned}$$

Note that the **red term** is just **SSW**. We also notice that the inner sum is with respect to  $j$ , so any terms regarding  $i$  can be treated as a constant term and taken to the front of the inner sum. Therefore, the **blue term** becomes:

$$\sum_{i=1}^k \sum_{j=1}^{n_i} (\bar{x}_i - \bar{x})^2 = \sum_{i=1}^k n_i (\bar{x}_i - \bar{x})^2$$

which is just **SSB**. The last term becomes:

$$\begin{aligned}
 2 \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)(\bar{x}_i - \bar{x}) &= 2 \sum_{i=1}^k (\bar{x}_i - \bar{x}) \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i) \\
 &= 2 \sum_{i=1}^k (\bar{x}_i - \bar{x}) \left[ \sum_{j=1}^{n_i} x_{ij} - \sum_{j=1}^{n_i} \bar{x}_i \right] \\
 &= 2 \sum_{i=1}^k (\bar{x}_i - \bar{x}) [n_i \cdot \bar{x}_i - n_i \cdot \bar{x}_i] \\
 &= 0
 \end{aligned}$$

Now, we see that  $\text{SST} = \text{SSB} + \text{SSW}$ . □

## 2 The Distribution Related To SSW

Let's write SSW in an estimator format. Recall that another way of writing SSW is:

$$SSW = \sum_{i=1}^k df_i S_i^2 = \sum_{i=1}^k (n_i - 1) S_i^2$$

Remember from **Lecture 16**, we derived that:

$$\frac{(n_i - 1) S_i^2}{\sigma_i^2} \sim \chi^2(n_i - 1)$$

If we sum them up, we have:

$$\sum_{i=1}^k \frac{(n_i - 1) S_i^2}{\sigma_i^2} \sim \chi^2(n - k)$$

Again, let's consider the simpler case where all populations have an equal variance  $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2 = \sigma^2$ . The above formula becomes:

$$\sum_{i=1}^k \frac{(n_i - 1) S_i^2}{\sigma_i^2} = \sum_{i=1}^k \frac{(n_i - 1) S_i^2}{\sigma^2} = \frac{\sum_{i=1}^k (n_i - 1) S_i^2}{\sigma^2} \sim \chi^2(n - k)$$

Note the numerator is **SSW**, so:

$$\frac{SSW}{\sigma^2} \sim \chi^2(n - k) \quad (1)$$

One thing we want to emphasise is that the above formula (1) is always true regardless of whether the null hypothesis is true or not.

## 3 The Distribution Related To SSB

Similarly, let's first write SSB in an estimator format:

$$SSB = \sum_{i=1}^k n_i (\bar{X}_i - \bar{\bar{X}})^2$$

From the central limit theorem, we have:

$$\bar{X}_i \sim \mathcal{N}\left(\mu_i, \frac{\sigma_i^2}{n_i}\right)$$

Again, we do the same thing as SSW, that is, assuming all population variances are equal:  $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2 = \sigma^2$ . Therefore,

$$\bar{\mathbf{X}}_i \sim \mathcal{N}\left(\mu_i, \frac{\sigma^2}{n_i}\right) \quad (2)$$

Now let's have a look at  $\bar{\bar{\mathbf{X}}}$ , which is the grand mean calculated by all data points:

$$\bar{\bar{\mathbf{X}}} = \frac{\sum_{i=1}^k n_i \cdot \bar{\mathbf{X}}_i}{n} = \sum_{i=1}^k \frac{n_i}{n} \bar{\mathbf{X}}_i$$

Apparently,  $\bar{\bar{\mathbf{X}}}$  is a weighted average of each of the sample mean  $\bar{\mathbf{X}}_i$ , and it also follows a normal distribution. Now, let's figure out its mean and variance.

$$\begin{aligned} \mathbb{E}[\bar{\bar{\mathbf{X}}}] &= \mathbb{E}\left[\sum_{i=1}^k \frac{n_i}{n} \bar{\mathbf{X}}_i\right] = \sum_{i=1}^k \mathbb{E}\left[\frac{n_i}{n} \bar{\mathbf{X}}_i\right] = \sum_{i=1}^k \frac{n_i}{n} \mathbb{E}[\bar{\mathbf{X}}_i] \\ &= \sum_{i=1}^k \frac{n_i}{n} \mu_i = \bar{\mu} \end{aligned}$$

where  $\bar{\mu}$  is just a weighted average of each population mean  $\mu_i$ . In terms of its variance:

$$\begin{aligned} \text{Var}(\bar{\bar{\mathbf{X}}}) &= \text{Var}\left(\sum_{i=1}^k \frac{n_i}{n} \bar{\mathbf{X}}_i\right) = \sum_{i=1}^k \text{Var}\left(\frac{n_i}{n} \bar{\mathbf{X}}_i\right) \quad \text{due to i.i.d. of } \mathbf{X}_i \\ &= \sum_{i=1}^k \frac{n_i^2}{n^2} \text{Var}(\bar{\mathbf{X}}_i) = \sum_{i=1}^k \frac{n_i^2}{n^2} \cdot \frac{\sigma^2}{n_i} = \sum_{i=1}^k \frac{n_i \sigma^2}{n^2} \\ &= \frac{\sigma^2}{n^2} \sum_{i=1}^k n_i = \frac{\sigma^2}{n^2} \cdot n = \frac{\sigma^2}{n} \end{aligned}$$

Therefore, we see that:

$$\bar{\bar{\mathbf{X}}} \sim \mathcal{N}\left(\bar{\mu}, \frac{\sigma^2}{n}\right) \quad (3)$$

If you think about it, it actually makes sense due to the central limit theorem.

Now let's get back to see SSB:

$$\begin{aligned}
 \text{SSB} &= \sum_{i=1}^k n_i (\bar{X}_i - \bar{\bar{X}})^2 = \sum_{i=1}^k n_i \left[ (\bar{X}_i - \mu_i) + (\mu_i - \bar{\bar{X}}) \right]^2 \\
 &= \sum_{i=1}^k \left[ n_i (\bar{X}_i - \mu_i)^2 + n_i (\bar{\bar{X}} - \mu_i)^2 + 2n_i (\bar{X}_i - \mu_i) (\mu_i - \bar{\bar{X}}) \right] \\
 &= \sum_{i=1}^k n_i (\bar{X}_i - \mu_i)^2 + \sum_{i=1}^k n_i (\bar{\bar{X}} - \mu_i)^2 + 2 \sum_{i=1}^k n_i (\bar{X}_i - \mu_i) (\mu_i - \bar{\bar{X}}) \quad (4)
 \end{aligned}$$

Now let's take a closer look at the blue term:

$$\begin{aligned}
 &\sum_{i=1}^k n_i (\bar{\bar{X}} - \mu_i)^2 + 2 \sum_{i=1}^k n_i (\bar{X}_i - \mu_i) (\mu_i - \bar{\bar{X}}) \\
 &= \sum_{i=1}^k \left( n_i \bar{\bar{X}}^2 - 2n_i \mu_i \bar{\bar{X}} + n_i \mu_i^2 \right) \\
 &\quad + 2 \sum_{i=1}^k \left( n_i \bar{X}_i \mu_i - n_i \bar{X}_i \bar{\bar{X}} - n_i \mu_i^2 + n_i \mu_i \bar{\bar{X}} \right) \\
 &= \bar{\bar{X}}^2 \sum_{i=1}^k n_i - 2\bar{\bar{X}} \sum_{i=1}^k n_i \mu_i + \sum_{i=1}^k n_i \mu_i^2 \\
 &\quad + 2 \sum_{i=1}^k n_i \bar{X}_i \mu_i - 2\bar{\bar{X}} \sum_{i=1}^k n_i \bar{X}_i - 2 \sum_{i=1}^k n_i \mu_i^2 + 2\bar{\bar{X}} \sum_{i=1}^k n_i \mu_i
 \end{aligned}$$

Note that

$$\sum_{i=1}^k n_i = n \quad \text{and} \quad \sum_{i=1}^k n_i \bar{X}_i = n \bar{\bar{X}}$$

so we can further simplify the expression as:

$$\begin{aligned}
 &n \bar{\bar{X}}^2 - 2\bar{\bar{X}} \sum_{i=1}^k n_i \mu_i + \sum_{i=1}^k n_i \mu_i^2 \\
 &\quad + 2 \sum_{i=1}^k n_i \bar{X}_i \mu_i - 2\bar{\bar{X}} \cdot n \bar{\bar{X}} - 2 \sum_{i=1}^k n_i \mu_i^2 + 2\bar{\bar{X}} \sum_{i=1}^k n_i \mu_i
 \end{aligned}$$

Merging the terms in the same colour, we can further simplify the expression as:

$$-n \bar{\bar{X}}^2 + 2 \sum_{i=1}^k n_i \bar{X}_i \mu_i - \sum_{i=1}^k n_i \mu_i^2$$

If the **null hypothesis**  $H_0$  is true, which means  $\mu_1 = \mu_2 = \dots = \mu_k$ , then we

have  $\mu_i = \bar{\mu}$ , where  $i = 1, 2, 3, \dots, k$ . Therefore, the above expression can be simplified again as:

$$\begin{aligned} -n\bar{\bar{X}}^2 + 2\sum_{i=1}^k n_i \bar{X}_i \bar{\mu} - \sum_{i=1}^k n_i \bar{\mu}^2 &= -n\bar{\bar{X}}^2 + 2\bar{\mu} \sum_{i=1}^k n_i \bar{X}_i - \bar{\mu}^2 \sum_{i=1}^k n_i \\ &= -n\bar{\bar{X}}^2 + 2\bar{\mu} \cdot n\bar{\bar{X}} - n\bar{\mu}^2 \\ &= -n \left( \bar{\bar{X}} - \bar{\mu} \right)^2 \end{aligned} \quad (5)$$

which is the final form of the **blue term** when  $H_0$  is true. Putting the formula (5) back to formula (4), we have:

$$\text{SSB} = \sum_{i=1}^k n_i (\bar{X}_i - \mu_i)^2 - n \left( \bar{\bar{X}} - \bar{\mu} \right)^2 \quad (6)$$

which is only true **if the null hypothesis  $H_0$**  is true.

Now things become clearer. From formula (2), we can get:

$$\begin{aligned} \frac{\sqrt{n_i} (\bar{X}_i - \mu_i)}{\sigma} &\sim \mathcal{N}(0, 1) \Rightarrow \frac{n_i (\bar{X}_i - \mu_i)^2}{\sigma^2} \sim \chi^2(1) \\ \Rightarrow \sum_{i=1}^k \frac{n_i (\bar{X}_i - \mu_i)^2}{\sigma^2} &\sim \chi^2(k) \end{aligned}$$

Similarly, from formula (3), we can get:

$$\frac{\sqrt{n} (\bar{\bar{X}} - \bar{\mu})}{\sigma} \sim \mathcal{N}(0, 1) \Rightarrow \frac{n (\bar{\bar{X}} - \bar{\mu})^2}{\sigma^2} \sim \chi^2(1)$$

Now we divide by  $\sigma^2$  at both sides of equation (6), we get:

$$\begin{aligned} \frac{\text{SSB}}{\sigma^2} &= \underbrace{\sum_{i=1}^k \frac{n_i (\bar{X}_i - \mu_i)^2}{\sigma^2}}_{\sim \chi^2(k)} - \underbrace{\frac{n (\bar{\bar{X}} - \bar{\mu})^2}{\sigma^2}}_{\sim \chi^2(1)} \\ \Rightarrow \frac{\text{SSB}}{\sigma^2} &\sim \chi^2(k-1) \end{aligned} \quad (7)$$

Here we want to emphasise that equation (7) is valid **only when** the null hypothesis  $H_0$  is true.

## 4 The $F$ -test

By definition, we have:

$$\text{MSB} = \frac{\text{SSB}}{k-1} \text{ and } \text{MSW} = \frac{\text{SSW}}{n-k}$$

Therefore, under the null hypothesis  $H_0$ , we can see that:

$$\frac{\text{MSB}}{\text{MSW}} = \frac{\frac{\text{SSB}}{k-1}}{\frac{\text{SSW}}{n-k}} = \frac{\frac{\text{SSB}}{\sigma^2} \cdot \frac{1}{k-1}}{\frac{\text{SSW}}{\sigma^2} \cdot \frac{1}{n-k}} = \frac{\frac{\chi^2(k-1)}{k-1}}{\frac{\chi^2(n-k)}{n-k}} \sim \mathcal{F}(k-1, n-k)$$

That's why we use the  $F$ -tests for ANOVA.