American University of Beirut STAT 210: Elementary Statistics for Sciences 2022–2023 Fall

Siamak Taati

Chapter 12 Analysis of variance

Analysis of variance (ANOVA for short) is an approach to testing the influence of factors (i.e., categorical variables) on a numerical variable.

For instance, we may want to test:

- ► Whether the <u>BMI</u> (body mass index) of adults is associated with their ethnicity.
- ▶ Whether there is a difference in <u>income</u> among three different professions.

Analysis of variance (ANOVA for short) is an approach to testing the influence of factors (i.e., categorical variables) on a numerical variable.

For instance, we may want to test:

- Whether the <u>BMI</u> (body mass index) of adults is associated with their <u>ethnicity</u>.
- Whether there is a difference in <u>income</u> among three different professions.

Here, we only consider the influence of a single factor. The method can be extended to multiple factors.

Analysis of variance (ANOVA for short) is an approach to testing the influence of factors (i.e., categorical variables) on a numerical variable.

We can view analysis of variance as

 an analogue of the "chi-squared test of independence" in which instead of testing the dependence of two categorical variables on each other, we test the <u>dependence of a</u> <u>numerical variable on one (or more) categorical variables.</u>

Analysis of variance (ANOVA for short) is an approach to testing the influence of factors (i.e., categorical variables) on a numerical variable.

We can view analysis of variance as

 an analogue of the "chi-squared test of independence" in which instead of testing the dependence of two categorical variables on each other, we test the <u>dependence of a</u> numerical variable on one (or more) categorical variables.

Alternatively, we can view analysis of variance as

 an extension of the "t-test for comparing the means of two populations based on independent samples", where we wish to compare the means of more than two populations.

The strategy will be as usual:

We use a statistic F and compare

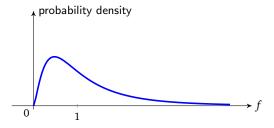
- 1. The observed value of F,
- 2. The distribution of F as suggested by (H_0) .

We ask ourselves:

 \bigcirc Is the observed value of F too extreme for \bigcirc to be plausible?

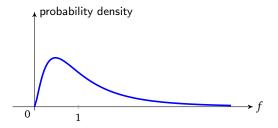
The statistics we will use turns out to have the so-called F-distribution under the null hypothesis.

The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



A continuous RV with an F-distribution is called an F-RV.

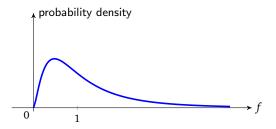
The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



A continuous RV with an F-distribution is called an F-RV.

The F-distribution is <u>unimodal</u> and <u>right-skewed</u>. An F-RV can only take non-negative values.

The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



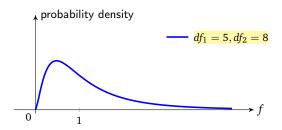
A continuous RV with an F-distribution is called an F-RV.

The F-distribution is <u>unimodal</u> and <u>right-skewed</u>. An F-RV can only take non-negative values.

The possible values of an F-RV are <u>all non-negative</u> numbers $0 \le f < +\infty$.



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



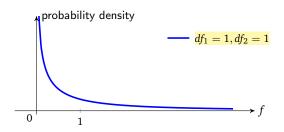
Remark 1

The mean of an F-RV is $\frac{df_2}{df_2+2}$. The mode of an F-RV is $\frac{df_1}{df_1} imes \frac{df_2}{df_2+2}$

[less than 1] [less than the mean]



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



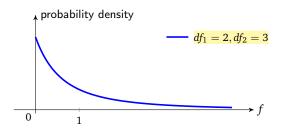
Remark 1

The mean of an F-RV is $\frac{df_2}{df_2+2}$. The mode of an F-RV is $\frac{df_1-2}{df_1} imes \frac{df_2}{df_2+2}$

[less than 1] [less than the mean]



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



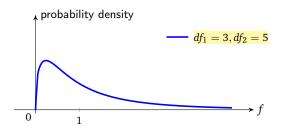
Remark 1

The mean of an F-RV is $\frac{df_2}{df_2+2}$. The mode of an F-RV is $\frac{df_1-2}{df_1} imes \frac{df_2}{df_2+2}$

[less than 1] [less than the mean]



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



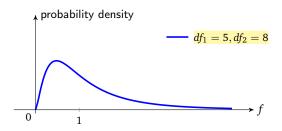
Remark 1

The mean of an F-RV is $\frac{df_2}{df_2+2}$. The mode of an F-RV is $\frac{df_1}{df_1} imes \frac{df_2}{df_2+2}$

[less than 1] [less than the mean]



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



Remark 1

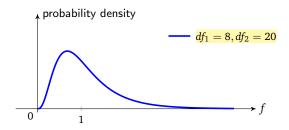
The mean of an F-RV is $\frac{df_2}{df_2+2}$. The mode of an F-RV is $\frac{df_1}{df_1} imes \frac{df_2}{df_2+2}$

[less than 1]

[less than the mean]



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



Remark 1

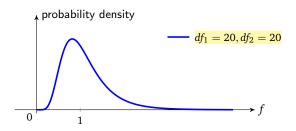
The mean of an F-RV is $\frac{df_2}{df_2+2}$. The mode of an F-RV is $\frac{df_2}{df_1} imes \frac{df_2}{df_2+2}$

[less than 1]

[less than the mean]



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



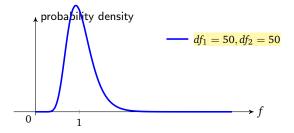
Remark 1

The mean of an F-RV is $\frac{df_2}{df_2+2}$. The mode of an F-RV is $\frac{df_1}{df_1} imes \frac{df_2}{df_2+2}$

[less than 1] [less than the mean]



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



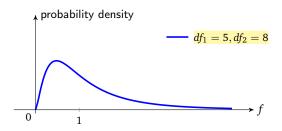
Remark 1

The mean of an F-RV is $\frac{df_2}{df_2+2}$. The mode of an F-RV is $\frac{df_2}{df_1} imes \frac{df_2}{df_2+2}$

[less than 1] [less than the mean]



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



Remark 1

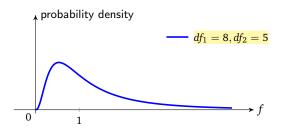
The mean of an F-RV is $\frac{df_2}{df_2+2}$. The mode of an F-RV is $\frac{df_1}{df_1} imes \frac{df_2}{df_2+2}$

[less than 1]

[less than the mean]



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



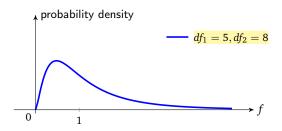
Remark 1

The mean of an F-RV is $\frac{df_2}{df_2+2}$. The mode of an F-RV is $\frac{df_1}{df_1} imes \frac{df_2}{df_2+2}$

[less than 1] [less than the mean]



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



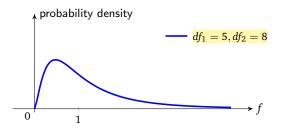
Remark 1

The mean of an F-RV is $\frac{df_2}{df_2+2}$. The mode of an F-RV is $\frac{df_1}{df_1} imes \frac{df_2}{df_2+2}$

[less than 1] [less than the mean]



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



Remark 2 (for those curious)

The F-distribution is inter-connected with the chi-squared distribution.

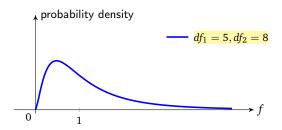
If X_1 and X_2 are independent chi-squared RVs with 5 and 8 degrees of freedom respectively, then

 $\frac{X_1/5}{X_2/8}$

has the F-distribution with parameters $df_1 = 5$ and $df_2 = 8$.



The F-distribution with parameters df_1 and df_2 (# of degrees of freedom for the numerator and denominator):



Remark 2 (for those curious)

The F-distribution is inter-connected with the chi-squared distribution.

If X_1 and X_2 are independent chi-squared RVs with df_1 and df_2 degrees of freedom respectively, then

 $\frac{X_1/df_1}{X_2/df_2}$

has the F-distribution with parameters df_1 and df_2 .



Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

Nutrition researchers have conducted an experiment on a sample of 29 rats. Each of the 29 subjects is randomly assigned one of four possible diets A, B, C, D. The following table contains the liver weight expressed as percentage of body weight.

diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84	
diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44
diet C	3.34	3.72	3.81	3.66	3.55	3.51		
diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91

The researchers would like to know whether the diet influences the liver weight.

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

Nutrition researchers have conducted an experiment on a sample of 29 rats. Each of the 29 subjects is randomly assigned one of four possible diets A, B, C, D. The following table contains the liver weight expressed as percentage of body weight.

diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84	
diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44
diet C	3.34	3.72	3.81	3.66	3.55	3.51		
diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91

The researchers would like to know whether the diet influences the liver weight.

Remark

Since this is a <u>randomized experiment</u>, we can potentially conclude causality from it. In this context, the categories are often called <u>treatments</u>.

Example (Effect of diet on liver weight)
[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]
The researchers would like to know whether the diet influences the liver weight.

Q What are the competing hypotheses?

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989] The researchers would like to know whether the diet influences the liver weight.

What are the competing hypotheses?

 $A \cap H_0$ Liver weight is not influenced by diet.

(H₁) Liver weight is influenced by diet. [i.e., is <u>not</u> independent of]

4 D > 4 P > 4 E > 4 E > 9 Q O

[i.e., is independent of]

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989] The researchers would like to know whether the diet influences the liver weight.

What are the competing hypotheses?

 $A \mid H_0$ Liver weight is not influenced by diet. [i.e., is independent of]

 $\overline{H_1}$ Liver weight is influenced by diet. [i.e., is <u>not</u> independent of]

We will assume that the population in each category is normally distributed, each with the same (unknown) standard deviation σ .

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989] The researchers would like to know whether the diet influences the liver weight.

Q What are the competing hypotheses?

 \overline{A} $\overline{H_0}$ Liver weight is not influenced by diet. [i.e., is independent of]

 $\overline{H_1}$ Liver weight is influenced by diet. [i.e., is <u>not</u> independent of]

We will assume that the population in each category is normally distributed, each with the same (unknown) standard deviation σ .

With these assumptions, the competing hypotheses can be rephrased as follows:

$$(H_0) \mu_A = \mu_B = \mu_C = \mu_D.$$

 (H_1) Not all population means are the same.

(μ_A is the population mean liver weight within category A, etc.)



Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

 (H_0) Liver weight is not influenced by diet. (H_1) Liver weight is influenced by diet.

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

 H_0 Liver weight is not influenced by diet. H_1 Liver weight is influenced by diet.

Q What is a suitable test statistic?

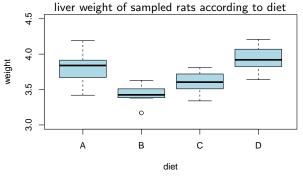
Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

- H_0 Liver weight is not influenced by diet. H_1 Liver weight is influenced by diet.
- Q What is a suitable test statistic?
- Where should we look to find evidence against the null hypothesis?

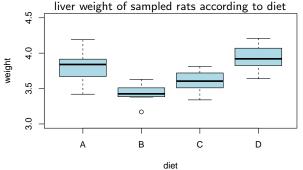
Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]



Example (Effect of diet on liver weight)

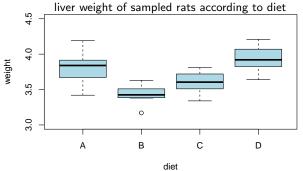
[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]



If the variations between categories are "large" compared to the variations within categories, then we have evidence in favor of the influence of diet on liver weight.

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

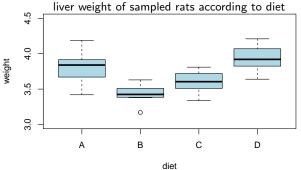


If the variations between categories are "large" compared to the variations within categories, then we have evidence in favor of the influence of diet on liver weight.

(Q) How can we quantify such a comparison?

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]



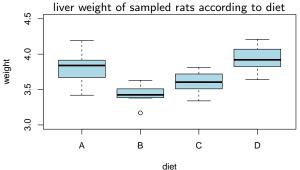
Two "variance-like" quantities to be introduced soon:

- MSB (mean square between samples)
 measures the variations between categories.
- MSW (mean square within samples)
 measures the variations within categories.



Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]



If MSB is "large" compared to MSW, then that counts as evidence against the null hypothesis.

[i.e., in favor of the influence of the diet on liver weight]

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

 H_0 Liver weight is not influenced by diet. H_1 Liver weight is influenced by diet.

Q What is a suitable test statistic?

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

- (H_0) Liver weight is not influenced by diet. (H_1) Liver weight is influenced by diet.
- What is a suitable test statistic?
- The statistic

$$F := \frac{\mathsf{MSB}}{\mathsf{MSW}}$$

measures the variations between categories relative to the variations within categories.

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

- (H_0) Liver weight is not influenced by diet. (H_1) Liver weight is influenced by diet.
- What is a suitable test statistic?
- The statistic

$$F := \frac{\mathsf{MSB}}{\mathsf{MSW}}$$

measures the variations between categories relative to the variations within categories.

What is the distribution of F according to (H_0)

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

- (H_0) Liver weight is not influenced by diet. (H_1) Liver weight is influenced by diet.
- What is a suitable test statistic?
- The statistic

$$F := \frac{\mathsf{MSB}}{\mathsf{MSW}}$$

measures the variations between categories relative to the variations within categories.

- What is the distribution of F according to (H_0)
- A It turns out that F has the F-distribution with $df_1 = 4 - 1$ (number of categories minus 1) and $df_2 = 29 - 4$ (total sample size minus number of categories).

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

 H_0 Liver weight is not influenced by diet. H_1 Liver weight is influenced by diet.

 \bigcirc What is the observed value of F?

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

 (H_0) Liver weight is not influenced by diet. (H_1) Liver weight is influenced by diet.

Q) What is the observed value of F?

We should first introduce MSB and MSW and learn how to calculate them.

Consider a sample of size n from a population with k distinct categories.

diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84	
diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44
diet C	3.34	3.72	3.81	3.66	3.55	3.51		
diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91

Consider a sample of size n from a population with k distinct categories.

▶ Let $X_{i,j}$ be the *j*th sampled entity within category *i*.

In our example, n = 29 and k = 4.

diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84	
diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44
diet C	3.34	3.72	3.81	3.66	3.55	3.51		
diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91

 $X_{2,3} = 3.38$ and $X_{4,1} = 3.64$.

Consider a sample of size n from a population with k distinct categories.

▶ Let $X_{i,j}$ be the *j*th sampled entity within category *i*.

In our example, n = 29 and k = 4.

diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84	
diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44
diet C	3.34	3.72	3.81	3.66	3.55	3.51		
diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91

 $X_{2,3} = 3.38$ and $X_{4,1} = 3.64$.

Consider a sample of size n from a population with k distinct categories.

▶ Let $X_{i,j}$ be the *j*th sampled entity within category *i*.

In our example, n = 29 and k = 4.

diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84	
diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44
diet C	3.34	3.72	3.81	3.66	3.55	3.51		
diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91

 $X_{2,3} = 3.38$ and $X_{4,1} = 3.64$.

Consider a sample of size n from a population with k distinct categories.

- ▶ Let X_{i,j} be the jth sampled entity within category i.
- ▶ Let \overline{X}_i be the mean of the sampled entities within category i.

diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84	
diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44
diet C	3.34	3.72	3.81	3.66	3.55	3.51		
diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91

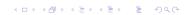
- $X_{2,3} = 3.38$ and $X_{4,1} = 3.64$.
- $\overline{X}_1 = \frac{3.42 + 3.96 + 3.87 + 4.19 + 3.58 + 3.76 + 3.84}{7} = 3.802857$

Consider a sample of size n from a population with k distinct categories.

- ▶ Let X_{i,j} be the jth sampled entity within category i.
- Let \overline{X}_i be the mean of the sampled entities within category i.
- Let n_i be the number of sampled entities within category i.

diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84	
diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44
diet C	3.34	3.72	3.81	3.66	3.55	3.51		
diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91

- $X_{2,3} = 3.38$ and $X_{4,1} = 3.64$.
- $\overline{X}_1 = \frac{3.42 + 3.96 + 3.87 + 4.19 + 3.58 + 3.76 + 3.84}{7} = 3.802857$
- $n_1 = 7$ and $n_2 = 8$.



Consider a sample of size n from a population with k distinct categories.

- ▶ Let X_{i,j} be the jth sampled entity within category i.
- Let \overline{X}_i be the mean of the sampled entities within category i.
- Let n_i be the number of sampled entities within category i.

diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84	
diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44
diet C	3.34	3.72	3.81	3.66	3.55	3.51		
diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91

- $X_{2,3} = 3.38$ and $X_{4,1} = 3.64$.
- $\overline{X}_1 = \frac{3.42 + 3.96 + 3.87 + 4.19 + 3.58 + 3.76 + 3.84}{7} = 3.802857$
- $n_1 = 7$ and $n_2 = 8$.



Consider a sample of size n from a population with k distinct categories.

Consider a sample of size n from a population with k distinct categories.

In this setting, the sample variance is

$$S^2 = \frac{1}{n-1} \sum_{i} \sum_{j} (X_{i,j} - \overline{X})^2$$
.

Consider a sample of size n from a population with k distinct categories.

In this setting, the sample variance is

$$S^2 = \frac{1}{n-1} \sum_{i} \sum_{j} (X_{i,j} - \overline{X})^2$$
.

A straightforward algebraic manipulation shows that

$$\underbrace{\sum_{i} \sum_{j} (X_{i,j} - \overline{X})^2}_{\mathsf{SST}} = \underbrace{\sum_{i} \sum_{j} (X_{i,j} - \overline{X}_i)^2}_{\mathsf{SSW}} + \underbrace{\sum_{i} n_i (\overline{X}_i - \overline{X})^2}_{\mathsf{SSB}} \ .$$

- SST: total sum of squares
- SSW:sum of squares within samples
- SSB: sum of squares between samples



Consider a sample of size n from a population with k distinct categories.

In this setting, the sample variance is

$$S^2 = \frac{1}{n-1} \sum_{i} \sum_{j} (X_{i,j} - \overline{X})^2$$
.

A straightforward algebraic manipulation shows that

$$\underbrace{\sum_{i} \sum_{j} (X_{i,j} - \overline{X})^2}_{\mathsf{SST}} = \underbrace{\sum_{i} \sum_{j} (X_{i,j} - \overline{X}_i)^2}_{\mathsf{SSW}} + \underbrace{\sum_{i} n_i (\overline{X}_i - \overline{X})^2}_{\mathsf{SSB}} \ .$$

This is a variant of the so-called law of total variance.

Mean squares between and within samples

Consider a sample of size n from a population with k distinct categories.

$$\underbrace{\sum_{i} \sum_{j} (X_{i,j} - \overline{X})^2}_{\mathsf{SST}} = \underbrace{\sum_{i} \sum_{j} (X_{i,j} - \overline{X}_i)^2}_{\mathsf{SSW}} + \underbrace{\sum_{i} n_i (\overline{X}_i - \overline{X})^2}_{\mathsf{SSB}} \ .$$

Mean squares between and within samples

Consider a sample of size n from a population with k distinct categories.

$$\underbrace{\sum_{i} \sum_{j} (X_{i,j} - \overline{X})^2}_{\mathsf{SST}} = \underbrace{\sum_{i} \sum_{j} (X_{i,j} - \overline{X}_i)^2}_{\mathsf{SSW}} + \underbrace{\sum_{i} n_i (\overline{X}_i - \overline{X})^2}_{\mathsf{SSB}} \ .$$

The mean squares between and within samples are

$$\mathsf{MSB} \coloneqq \frac{\mathsf{SSB}}{k-1}$$

$$\boxed{\mathsf{MSW} \coloneqq \frac{\mathsf{SSW}}{n-k}}$$

The sums of squares between and within sample can alternatively be calculated as follows:

$$SSB = \sum_{i} \frac{1}{n_i} \left(\sum_{j} X_{i,j} \right)^2 - \frac{1}{n} \left(\sum_{i} \sum_{j} X_{i,j} \right)^2$$
$$SSW = \sum_{i} \sum_{j} X_{i,j}^2 - \sum_{i} \frac{1}{n_i} \left(\sum_{j} X_{i,j} \right)^2$$

The above formulas are somewhat more efficient than the original definitions. [i.e., they take less computation]

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$SSB = \sum_{i} \frac{1}{n_i} \left(\sum_{j} X_{i,j} \right)^2 - \frac{1}{n} \left(\sum_{i} \sum_{j} X_{i,j} \right)^2$$

$$\mathsf{SSW} = \sum_{i} \sum_{j} X_{i,j}^2 - \sum_{i} \frac{1}{n_i} \Big(\sum_{j} X_{i,j} \Big)^2$$

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$\begin{aligned} \mathsf{SSB} &= \sum_{i} \frac{1}{n_i} \Big(\sum_{j} X_{i,j} \Big)^2 - \frac{1}{n} \Big(\sum_{i} \sum_{j} X_{i,j} \Big)^2 \\ &= \left[\frac{1}{7} \times 26.62^2 + \frac{1}{8} \times 27.44^2 + \frac{1}{6} \times 21.59^2 + \frac{1}{8} \times 31.48^2 \right] - \left[\frac{1}{29} \times 107.13^2 \right] \end{aligned}$$

$$SSW = \sum_{i} \sum_{j} X_{i,j}^2 - \sum_{i} \frac{1}{n_i} \left(\sum_{j} X_{i,j}\right)^2$$

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$SSB = \sum_{i} \frac{1}{n_{i}} \left(\sum_{j} X_{i,j} \right)^{2} - \frac{1}{n} \left(\sum_{i} \sum_{j} X_{i,j} \right)^{2}$$
$$= \left[\frac{1}{7} \times \frac{26.62^{2}}{1} + \frac{1}{8} \times 27.44^{2} + \frac{1}{6} \times 21.59^{2} + \frac{1}{8} \times 31.48^{2} \right] - \left[\frac{1}{29} \times 107.13^{2} \right]$$

$$SSW = \sum_{i} \sum_{j} X_{i,j}^2 - \sum_{i} \frac{1}{n_i} \left(\sum_{j} X_{i,j} \right)^2$$

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$\begin{aligned} \mathsf{SSB} &= \sum_{i} \frac{1}{n_i} \Big(\sum_{j} X_{i,j} \Big)^2 - \frac{1}{n} \Big(\sum_{i} \sum_{j} X_{i,j} \Big)^2 \\ &= \left[\frac{1}{7} \times 26.62^2 + \frac{1}{8} \times \frac{27.44^2}{6} + \frac{1}{6} \times 21.59^2 + \frac{1}{8} \times 31.48^2 \right] - \left[\frac{1}{29} \times 107.13^2 \right] \end{aligned}$$

$$SSW = \sum_{i} \sum_{j} X_{i,j}^2 - \sum_{i} \frac{1}{n_i} \left(\sum_{j} X_{i,j} \right)^2$$

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$SSB = \sum_{i} \frac{1}{n_{i}} \left(\sum_{j} X_{i,j} \right)^{2} - \frac{1}{n} \left(\sum_{i} \sum_{j} X_{i,j} \right)^{2}$$
$$= \left[\frac{1}{7} \times 26.62^{2} + \frac{1}{8} \times 27.44^{2} + \frac{1}{6} \times \frac{21.59^{2}}{6} + \frac{1}{8} \times 31.48^{2} \right] - \left[\frac{1}{29} \times 107.13^{2} \right]$$

$$SSW = \sum_{i} \sum_{j} X_{i,j}^2 - \sum_{i} \frac{1}{n_i} \left(\sum_{j} X_{i,j}\right)^2$$

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$SSB = \sum_{i} \frac{1}{n_{i}} \left(\sum_{j} X_{i,j} \right)^{2} - \frac{1}{n} \left(\sum_{i} \sum_{j} X_{i,j} \right)^{2}$$
$$= \left[\frac{1}{7} \times 26.62^{2} + \frac{1}{8} \times 27.44^{2} + \frac{1}{6} \times 21.59^{2} + \frac{1}{8} \times \frac{31.48^{2}}{31.48^{2}} \right] - \left[\frac{1}{29} \times 107.13^{2} \right]$$

$$SSW = \sum_{i} \sum_{j} X_{i,j}^2 - \sum_{i} \frac{1}{n_i} \left(\sum_{j} X_{i,j} \right)^2$$

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$SSB = \sum_{i} \frac{1}{n_{i}} \left(\sum_{j} X_{i,j} \right)^{2} - \frac{1}{n} \left(\sum_{i} \sum_{j} X_{i,j} \right)^{2}$$
$$= \left[\frac{1}{7} \times 26.62^{2} + \frac{1}{8} \times 27.44^{2} + \frac{1}{6} \times 21.59^{2} + \frac{1}{8} \times 31.48^{2} \right] - \left[\frac{1}{29} \times \frac{107.13^{2}}{2} \right]$$

$$SSW = \sum_{i} \sum_{j} X_{i,j}^2 - \sum_{i} \frac{1}{n_i} \left(\sum_{j} X_{i,j}\right)^2$$

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$\begin{split} \mathsf{SSB} &= \sum_i \frac{1}{n_i} \Big(\sum_j X_{i,j} \Big)^2 - \frac{1}{n} \Big(\sum_i \sum_j X_{i,j} \Big)^2 \\ &= \left[\frac{1}{7} \times 26.62^2 + \frac{1}{8} \times 27.44^2 + \frac{1}{6} \times 21.59^2 + \frac{1}{8} \times 31.48^2 \right] - \left[\frac{1}{29} \times 107.13^2 \right] \\ &= \boxed{1.160077} \; . \\ \mathsf{SSW} &= \sum_i \sum_j X_{i,j}^2 - \sum_i \frac{1}{n_i} \Big(\sum_j X_{i,j} \Big)^2 \end{split}$$

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$\begin{split} \mathsf{SSB} &= \sum_i \frac{1}{n_i} \Big(\sum_j X_{i,j} \Big)^2 - \frac{1}{n} \Big(\sum_i \sum_j X_{i,j} \Big)^2 \\ &= \left[\frac{1}{7} \times 26.62^2 + \frac{1}{8} \times 27.44^2 + \frac{1}{6} \times 21.59^2 + \frac{1}{8} \times 31.48^2 \right] - \left[\frac{1}{29} \times 107.13^2 \right] \\ &= \boxed{1.160077} \; . \\ \mathsf{SSW} &= \sum_i \sum_j X_{i,j}^2 - \sum_i \frac{1}{n_i} \Big(\sum_j X_{i,j} \Big)^2 \\ &= 397.8143 - \left[\frac{1}{7} \times 26.62^2 + \frac{1}{8} \times 27.44^2 + \frac{1}{6} \times 21.59^2 + \frac{1}{8} \times 31.48^2 \right] \end{split}$$

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$\begin{split} \text{SSB} &= \sum_i \frac{1}{n_i} \Big(\sum_j X_{i,j} \Big)^2 - \frac{1}{n} \Big(\sum_i \sum_j X_{i,j} \Big)^2 \\ &= \left[\frac{1}{7} \times 26.62^2 + \frac{1}{8} \times 27.44^2 + \frac{1}{6} \times 21.59^2 + \frac{1}{8} \times 31.48^2 \right] - \left[\frac{1}{29} \times 107.13^2 \right] \\ &= \boxed{1.160077} \; . \\ \text{SSW} &= \sum_i \sum_j X_{i,j}^2 - \sum_i \frac{1}{n_i} \Big(\sum_j X_{i,j} \Big)^2 \\ &= \boxed{397.8143} - \left[\frac{1}{7} \times 26.62^2 + \frac{1}{8} \times 27.44^2 + \frac{1}{6} \times 21.59^2 + \frac{1}{8} \times 31.48^2 \right] \end{split}$$

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$\begin{split} \text{SSB} &= \sum_i \frac{1}{n_i} \Big(\sum_j X_{i,j} \Big)^2 - \frac{1}{n} \Big(\sum_i \sum_j X_{i,j} \Big)^2 \\ &= \left[\frac{1}{7} \times 26.62^2 + \frac{1}{8} \times 27.44^2 + \frac{1}{6} \times 21.59^2 + \frac{1}{8} \times 31.48^2 \right] - \left[\frac{1}{29} \times 107.13^2 \right] \\ &= \boxed{1.160077} \; . \\ \text{SSW} &= \sum_i \sum_j X_{i,j}^2 - \sum_i \frac{1}{n_i} \Big(\sum_j X_{i,j} \Big)^2 \\ &= 397.8143 - \left[\frac{1}{7} \times 26.62^2 + \frac{1}{8} \times 27.44^2 + \frac{1}{6} \times 21.59^2 + \frac{1}{8} \times 31.48^2 \right] \end{split}$$

i										total	$\sum x^2$
1	diet A	3.42	3.96	3.87	4.19	3.58	3.76	3.84		26.62	101.6106
2	diet B	3.17	3.63	3.38	3.47	3.39	3.41	3.55	3.44	27.44	94.2474
3	diet C	3.34	3.72	3.81	3.66	3.55	3.51			21.59	77.8283
4	diet D	3.64	3.93	3.77	4.18	4.21	3.88	3.96	3.91	31.48	124.1280
total										107.13	397.8143

$$\begin{aligned} \text{SSB} &= \sum_{i} \frac{1}{n_{i}} \Big(\sum_{j} X_{i,j} \Big)^{2} - \frac{1}{n} \Big(\sum_{i} \sum_{j} X_{i,j} \Big)^{2} \\ &= \Big[\frac{1}{7} \times 26.62^{2} + \frac{1}{8} \times 27.44^{2} + \frac{1}{6} \times 21.59^{2} + \frac{1}{8} \times 31.48^{2} \Big] - \Big[\frac{1}{29} \times 107.13^{2} \Big] \\ &= \boxed{1.160077} \ . \\ \text{SSW} &= \sum_{i} \sum_{j} X_{i,j}^{2} - \sum_{i} \frac{1}{n_{i}} \Big(\sum_{j} X_{i,j} \Big)^{2} \\ &= 397.8143 - \Big[\frac{1}{7} \times 26.62^{2} + \frac{1}{8} \times 27.44^{2} + \frac{1}{6} \times 21.59^{2} + \frac{1}{8} \times 31.48^{2} \Big] \\ &= \boxed{0.9012262} \ . \end{aligned}$$

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

 (H_0) Liver weight is not influenced by diet. (H_1)

 (H_1) Liver weight is influenced by diet.

 \bigcirc What is the observed value of F?

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

- (H_0) Liver weight is not influenced by diet. (H_1) Liver weight is influenced by diet.
- What is the observed value of F?
- The mean squares between and within samples are

$$\mathsf{MSB} = \frac{\mathsf{SSB}}{k-1} = \frac{1.160077}{4-1} \approx 0.3866924$$

$$\mathsf{MSW} = \frac{\mathsf{SSW}}{n-k} = \frac{0.9012262}{29-4} \approx 0.03604905$$

Therefore, the observed value of *F* is

$$f = \frac{\text{MSB}}{\text{MSW}} \approx \frac{0.3866924}{0.03604905} \approx \boxed{10.72684}$$
.

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

 H_0 Liver weight is not influenced by diet. H_1 Liver weight is influenced by diet.

 \bigcirc At significance level 1%, how should the researchers conclude?

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

- H_0 Liver weight is not influenced by diet. H_1 Liver weight is influenced by diet.
- ${f Q}$ At significance level 1%, how should the researchers conclude?
- A We have

$$p
-value = \mathbb{P}(F \ge 10.72684 \mid H_0) \approx$$

using the app for the F-distribution with $df_1=4-1$ and $df_2=29-4$.

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

- H_0 Liver weight is not influenced by diet. H_1 Liver weight is influenced by diet.
- ${f Q}$ At significance level 1%, how should the researchers conclude?
- A We have

$$\mathsf{p\text{-}value} = \mathbb{P}(F \geq 10.72684 \,|\: \textcolor{red}{(\mathsf{H}_0)}) \approx \boxed{0.0001}$$

using the app for the F-distribution with $df_1=4-1$ and $df_2=29-4$.

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

- (H_0) Liver weight is not influenced by diet. (H_1) Liver weight is influenced by diet.
- At significance level 1%, how should the researchers conclude?
- A We have

$$\mathsf{p\text{-}value} = \mathbb{P}(F \geq 10.72684 \,|\: \textcolor{red}{(\mathsf{H}_0)}) \approx \boxed{0.0001}$$

using the app for the F-distribution with $df_1 = 4 - 1$ and $df_2 = 29 - 4$.

Since p-value $< \alpha$, the researchers must reject the null hypothesis.

Example (Effect of diet on liver weight)

[This example is taken from S.C. Campbell, Statistics for Biologists, Cambridge, 1989]

- (H_0) Liver weight is not influenced by diet. (H_1) Liver weight is influenced by diet.
- At significance level 1%, how should the researchers conclude?
- A | We have

$$\mathsf{p\text{-}value} = \mathbb{P}(F \geq 10.72684 \,|\: \textcolor{red}{(\mathsf{H}_0)}) \approx \boxed{0.0001}$$

using the app for the F-distribution with $df_1 = 4 - 1$ and $df_2 = 29 - 4$.

Since p-value $< \alpha$, the researchers must reject the null hypothesis.

Conclusion: The evidence is statistically significant (p-value $\approx 0.01\%$) in favor of the claim that the liver weight of the rats is influenced by their diets.

Suppose C is a categorical variable with k possible values, and X is a numerical variable.

Suppose C is a categorical variable with k possible values, and X is a numerical variable.

Competing hypotheses (v1)

 (H_0) X is independent on C.

 H_1 X is dependent on C.

[X is not influenced by C]

[X is influenced by C]

Suppose C is a categorical variable with k possible values, and X is a numerical variable.

Competing hypotheses (v1)

 (H_0) X is independent on C.

[X is not influenced by C]

 (H_1) X is dependent on C.

[X is influenced by C]

Evidence

The values of C and X for a random sample of size n.

Suppose C is a categorical variable with k possible values, and X is a numerical variable.

Competing hypotheses (v1)

 (H_0) X is independent on C.

[X is not influenced by C]

 H_1 X is dependent on C.

[X is influenced by C]

Evidence

The values of C and X for a random sample of size n.

Test statistic

$$F := \frac{\mathsf{MSB}}{\mathsf{MSW}}$$

where

- MSB := SSB/(k-1) is the mean square between samples,
- MSW := SSW/(n k) is the mean square within samples.

Suppose C is a categorical variable with k possible values, and C is a numerical variable.

Competing hypotheses (v1)

- (H_0) X is independent on C.
- (H_1) X is dependent on C.

[X is not influenced by C]

[X is influenced by C]

Covered scenario

The population in each category is normal with the same (unknown) variance σ^2 .

Suppose C is a categorical variable with k possible values, and C is a numerical variable.

Competing hypotheses (v1)

- (H_0) X is independent on C.
- (H_1) X is dependent on C.

- [X is not influenced by C]
 - [X is influenced by C]

Covered scenario

The population in each category is normal with the same (unknown) variance σ^2 .

Strategy

We compare the observed value of F against the F-distribution with $df_1 = k - 1$ and $df_2 = n - k$. We can follow either the rejection region approach or the p-value approach.

Suppose C is a categorical variable with k possible values, and C is a numerical variable.

Competing hypotheses (v1)

 (H_0) X is independent on C.

[X is not influenced by C]

 (H_1) X is dependent on C.

[X is influenced by C]

Covered scenario

The population in each category is normal with the same (unknown) variance σ^2 .

Strategy

We compare the observed value of F against the F-distribution with $df_1 = k - 1$ and $df_2 = n - k$. We can follow either the rejection region approach or the p-value approach.

Remark

We always use a right-tailed test.



Suppose X is a numerical variable, and the means of X in k distinct populations are $\mu_1, \mu_2, \dots, \mu_k$.

Suppose X is a numerical variable, and the means of X in k distinct populations are $\mu_1, \mu_2, \dots, \mu_k$.

Competing hypotheses (v2)

 (H_1) Not all k means are equal.

Suppose X is a numerical variable, and the means of X in k distinct populations are $\mu_1, \mu_2, \dots, \mu_k$.

Competing hypotheses (v2)

- $(H_0) \mu_1 = \mu_2 = \cdots = \mu_k.$
- (H_1) Not all k means are equal.

Evidence

k random samples, one from each population, with total size n.

Suppose X is a numerical variable, and the means of X in k distinct populations are $\mu_1, \mu_2, \dots, \mu_k$.

Competing hypotheses (v2)

- $H_0 \mu_1 = \mu_2 = \cdots = \mu_k.$
- (H_1) Not all k means are equal.

Evidence

k random samples, one from each population, with total size n.

Test statistic

$$F := \frac{\mathsf{MSB}}{\mathsf{MSW}}$$

where

- MSB := SSB/(k-1) is the mean square between samples,
- MSW := SSW/(n k) is the mean square within samples.

Suppose X is a numerical variable, and the means of X in k distinct populations are $\mu_1, \mu_2, \dots, \mu_k$.

Competing hypotheses (v2)

- $(H_0) \ \mu_1 = \mu_2 = \dots = \mu_k.$
- (H_1) Not all k means are equal.

Covered scenario

The samples are independent of one another. Each population is normal. All populations have the same (unknown) variance σ^2 .

Suppose X is a numerical variable, and the means of X in k distinct populations are $\mu_1, \mu_2, \dots, \mu_k$.

Competing hypotheses (v2)

- $(H_0) \mu_1 = \mu_2 = \cdots = \mu_k.$
- (H_1) Not all k means are equal.

Covered scenario

The samples are independent of one another. Each population is normal. All populations have the same (unknown) variance σ^2 .

Strategy

We compare the observed value of F against the F-distribution with $df_1 = k - 1$ and $df_2 = n - k$. We can follow either the rejection region approach or the p-value approach.

Suppose X is a numerical variable, and the means of X in k distinct populations are $\mu_1, \mu_2, \dots, \mu_k$.

Competing hypotheses (v2)

- $(H_0) \mu_1 = \mu_2 = \cdots = \mu_k.$
- (H_1) Not all k means are equal.

Covered scenario

The samples are independent of one another. Each population is normal. All populations have the same (unknown) variance σ^2 .

Strategy

We compare the observed value of F against the F-distribution with $df_1 = k - 1$ and $df_2 = n - k$. We can follow either the rejection region approach or the p-value approach.

Remark

We always use a right-tailed test.

