

The One-Way Analysis of Variance (ANOVA) Table

Recall the data on the amount of fat absorbed by batches of doughnuts cooked in four different types of fat.

Fat Type	Fat Absorbed in grams						Average	Variance
1	65	73	69	78	57	96	73	178.0
2	78	91	97	82	85	77	85	60.4
3	75	93	78	71	63	76	76	97.6
4	55	66	49	64	70	68	62	67.6

Y_{ij} is often used to denote an observation (data point). The first subscript i denotes the group associated with the observation. The second subscript j denotes the observation number within its group. For example, $Y_{24} = 82$, the 4th observation from group 2. Find Y_{11} , Y_{34} , and Y_{43} .

Calculations related to the ANOVA F -test are often organized in an ANOVA table as follows:

Source	D.F.	Sum of Squares	Mean Squares	F
Between	$I - 1$	$SSB = \sum_{i=1}^I n_i (\bar{Y}_i - \bar{Y})^2$	$MSB = \frac{\sum_{i=1}^I n_i (\bar{Y}_i - \bar{Y})^2}{I - 1}$	$\frac{MSB}{MSW}$
Within	$n - I$	$SSW = \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$	$MSW = \frac{\sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2}{n - I}$	
Total	$n - 1$	$SSTO = \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y})^2$		

(Note that the formula for MSW in the table above is equal to $s_p^2 = \frac{\sum_{i=1}^I (n_i - 1) s_i^2}{n - I}$ as we learned earlier.)

We may partition the total sum of squares into two pieces.

$$\begin{aligned}
 \sum_{i=1}^I n_i (\bar{Y}_i - \bar{Y})^2 + \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2 &= \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y})^2 \\
 SSB + SSW &= SSTO \\
 SS_{BETWEEN} + SS_{WITHIN} &= SS_{TOTAL} \\
 SS_{MODEL} + SS_{ERROR} &= SS_{TOTAL}
 \end{aligned}$$

We call the table above an *analysis of variance (ANOVA)* table because the table partitions one measure of total variance in the data ($SSTO$) into two components. One component is a measure of variance in the data due to differences between group means (SSB). The other is a measure of variance in the data due to variance of the observations within groups (SSW). Compute the ANOVA table for the doughnut data.

The F-Test as a Comparison of Full and Reduced Models

We can view our F -test of $H_0 : \mu_1 = \mu_2 = \dots = \mu_I$ against H_A : *not all μ_i are equal* as a comparison of two models where one model is a special case of the other. The *full model* says that each group has its own mean μ_i . Each of these means can equal any value with no restrictions. When the null hypothesis H_0 is true, a simpler model holds where all the group means are equal to one common value, say μ . That common value μ can be anything, but the point is all the group means are equal to some unknown value μ . This simpler model is a special case of the full model. We call the simpler model the *reduced model*.

If the alternative hypothesis H_A is true, it makes sense to estimate the mean for group i by the sample mean for group i . If the null hypothesis H_0 is true, the reduced model holds. Thus it makes sense to estimate the mean for group i by the mean of all the observations \bar{Y} because all the observations come from a single distribution with some mean μ when the reduced model holds. In the doughnut example we can describe the full and reduced model parameters and estimates with the following tables.

Group Mean Parameters					Group Mean Estimates					Computed Estimates							
		Group						Group						Group			
Model		1	2	3	4	Model		1	2	3	4	Model		1	2	3	4
Full		μ_1	μ_2	μ_3	μ_4	Full		\bar{Y}_1	\bar{Y}_2	\bar{Y}_3	\bar{Y}_4	Full		73	85	76	62
Reduced		μ	μ	μ	μ	Reduced		\bar{Y}	\bar{Y}	\bar{Y}	\bar{Y}	Reduced		74	74	74	74

A *residual* is the value of an observation minus its estimated mean ($\hat{e}_{ij} = Y_{ij} - \hat{Y}_{ij}$). The residuals for the full and reduced models are given in the table below.

i	j	Y_{ij}	Full Model			Reduced Model		
			\hat{Y}_{ij}	\hat{e}_{ij}	\hat{e}_{ij}^2	\hat{Y}_{ij}	\hat{e}_{ij}	\hat{e}_{ij}^2
1	1	65	73	-8	64	74	-9	81
1	2	73	73	0	0	74	-1	1
1	3	69	73	-4	16	74	-5	25
1	4	78	73	5	25	74	4	16
1	5	57	73	-16	256	74	-17	289
1	6	96	73	23	529	74	22	484
2	1	78	85	-7	49	74	4	16
2	2	91	85	6	36	74	17	289
2	3	97	85	12	144	74	23	529
2	4	82	85	-3	9	74	8	64
2	5	85	85	0	0	74	11	121
2	6	77	85	-8	64	74	3	9
3	1	75	76	-1	1	74	1	1
3	2	93	76	17	289	74	19	361
3	3	78	76	2	4	74	4	16
3	4	71	76	-5	25	74	-3	9
3	5	63	76	-13	169	74	-11	121
3	6	76	76	0	0	74	2	4
4	1	55	62	-7	49	74	-19	361
4	2	66	62	4	16	74	-8	64
4	3	49	62	-13	169	74	-25	625
4	4	64	62	2	4	74	-10	100
4	5	70	62	8	64	74	-4	16
4	6	68	62	6	36	74	-6	36
			RSS(full)=2018			RSS(red.)=3638		

RSS stands for *residual sum of squares*. The RSS value for a particular model is simply the sum of the squared residuals for that model. The degrees of freedom associated with an RSS value is $df_{RSS} = n - p$, where p is the number of parameters estimated when computing the residuals. In general the statistic

$$F = \frac{[\text{RSS}(\text{red.}) - \text{RSS}(\text{full})] / [df_{\text{RSS}(\text{red.})} - df_{\text{RSS}(\text{full})}]}{\text{RSS}(\text{full}) / df_{\text{RSS}(\text{full})}}$$

can be used to determine if a full model fits significantly better than a reduced model. The null hypothesis of the test says that the reduced model is correct. The alternative hypothesis says that the reduced model is too simple and that the more complex full model is more appropriate. To determine a p -value, the F -statistic is compared to an F -distribution with numerator degrees of freedom equal $df_{\text{RSS}(\text{red.})} - df_{\text{RSS}(\text{full})}$ and denominator degrees of freedom equal $df_{\text{RSS}(\text{full})}$. Show that this F -statistic is the same as that computed previously for the doughnut data.