

# Formulas

Stephanie Ranft S2459825

October 17, 2019

## Statistics 2 PSBE2-07

### Variance

#### Pooled variance

Consider a test between  $I$  independent samples. Whilst we cannot assume that they are all from the same population (and hence have the same variance), the Central Limit Theorem allows us to conclude that the pooled variance converges to the true variance. To see how this works, and in order to better understand when and where to use pooled variance:

$$s_p^2 = \frac{\sum_{i=1}^I (n_i - 1) s_i^2}{\sum_{i=1}^I (n_i - 1)}, \quad \text{where } I \text{ is the total number of groups.} \quad (1)$$

$$= \frac{(n_1 - 1) s_1^2 + (n_2 - 1) s_2^2 + \dots + (n_I - 1) s_I^2}{\underbrace{(n_1 - 1) + (n_2 - 1) + \dots + (n_I - 1)}_{n-I}} \quad (2)$$

If we were to assume that one of the samples was a lot bigger than the others, suppose sample  $k$  ( $n_k \gg n_i$  for every sample), then we can assume that the pooled variance will converge to  $s_k^2$ :

$$= \frac{\frac{n_1-1}{n_k-1} s_1^2 + \dots + s_k^2 + \dots + \frac{n_I-1}{n_k-1} s_I^2}{\frac{n_1-1}{n_k-1} + \dots + \frac{n_k-1}{n_k-1} + \dots + \frac{n_I-1}{n_k-1}} \quad (3)$$

$$\sim s_k^2. \quad (4)$$

The reason for this has something to do with the “power” of having a large  $n$ ; that the sample is reliably similar to the population (lower standard error). This is an important point in sample collection, because if all samples except one have a really small size ( $<30$ ) but one sample is substantially larger ( $>100$ ), then the statistician can more readily believe the results of the largest sample (due to CLT - recall formula for SE). If we have that all of the samples are nearly the same size (choose  $n_k \approx n_1 \approx \dots \approx n_I$ ), then there is no dominating variance and we can assume the following:

$$s_p^2 \approx \frac{(n_k - 1) s_1^2 + (n_k - 1) s_2^2 + \dots + (n_k - 1) s_I^2}{(n_k - 1) + (n_k - 1) + \dots + (n_k - 1)} \quad (5)$$

$$= \frac{\cancel{(n_k - 1)} \times (s_1^2 + s_2^2 + \dots + s_I^2)}{I \times \cancel{(n_k - 1)}} \quad (6)$$

$$= \frac{\sum_{i=1}^I s_i^2}{I} \quad (7)$$

$$= \bar{s}^2, \quad \text{which is the mean of the variances.} \quad (8)$$

This indicates that the pooled variances is the weighted average of variances, where the highest weight is given to the largest sample size. This is to ensure that our pooled variance is the best estimator of the population variance. Another way to look at this is to look at the fact that  $s_i^2 = \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 / (n_i - 1)$ , where  $j$  indexes the persons in a group.

$$\implies s_p^2 = \frac{\sum_{j=1}^{n_1} (y_{1j} - \bar{y}_1)^2 + \sum_{j=1}^{n_2} (y_{2j} - \bar{y}_2)^2 + \dots + \sum_{j=1}^{n_I} (y_{Ij} - \bar{y}_I)^2}{n - I} \quad (9)$$

$$= \frac{\sum_{i=1}^I \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2}{n - I} \quad (10)$$

$$= \frac{SSE}{df_E} = MSE. \quad (11)$$

It is important to know when exactly to use pooled v.s. unpooled variance, but if we can **assume that the variances of the populations are equal**, i.e.  $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_I^2$ , then we use pooled variance.

For  $I = 2$  (comparing two means), we can use the  $t$  statistic to determine the results of our test:

$$t = \frac{\bar{y}_1 - \bar{y}_2}{s_p \times \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t \left( \underbrace{n_1 + n_2 - 2}_n \right) \quad (12)$$

$$\implies t^2 = \frac{(\bar{y}_1 - \bar{y}_2)^2}{s_p^2 \times \underbrace{\frac{n_1 + n_2}{n_1 n_2}}_{\frac{1}{n_1} + \frac{1}{n_2}}} \quad (13)$$

$$= \frac{(\bar{y}_1 - \bar{y}_2)^2 / (n_1 + n_2)}{s_p^2 / n_1 n_2} \quad (14)$$

If  $n_1 \approx n_2$

$$\implies t^2 \approx \frac{\frac{n}{2} (\bar{y}_1 - \bar{y}_2)^2}{s_p^2} \sim F(1, n_1 + n_2 - 2). \quad (15)$$

For  $I > 2$  (comparing multiple means), we use the  $F$  statistic:

$$F = \frac{MSG}{MSE} \quad (16)$$

$$= \frac{SSG/df_G}{SSE/df_E} \quad (17)$$

$$= \frac{\sum_{i=1}^I \sum_{j=1}^{n_i} (\bar{y}_i - \bar{y})^2 / (I - 1)}{\sum_{i=1}^I \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 / \underbrace{(n - I)}_{=df_T - df_G}} \quad (18)$$

$$= \frac{\sum_{i=1}^I n_i \times (\bar{y}_i - \bar{y})^2 / (I - 1)}{\sum_{i=1}^I \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 / (n - I)} \sim F \left( \overbrace{I - 1}^{df_G}, \overbrace{n - I}^{df_T - df_G} \right) \quad (19)$$

More on this when we discuss ANOVA...

## Unpooled variance

If we cannot see from the data that the independent samples are drawn from the same population, or that the  $I$  variances are similar, then we use unpooled variance. In general, we do not consider using this for  $I > 2$  (for comparing more than two means). It is mathematically possible, but in practice it is seldom used and certainly not a part of the scope of this course. The main reason being that the calculations required to determine the degrees of freedom is quite complicated (see (21)), and most psychologists use a computer. So, the **only time you use unpooled variance is when you have observable differences between the two groups**. So, you could note that  $s_1^2 \gg s_2^2$  or perhaps the  $n_i$  of each group is low ( $< 30$ ) and unequal; it is something you need to determine for yourself. For example, if your test was about two machines from different manufacturers and whether they can complete the same task in the same amount of time, then you would use unpooled variance. If you were exploring behavioural data two culturally different countries, you would use unpooled variances.

**Rule of thumb:** if it's only two groups and you can assume/predict that the variances are unequal, use unpooled.

$$s_{up} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}. \quad (20)$$

$$\implies t = \frac{\bar{y}_1 - \bar{y}_2}{s_{up}} \sim t(k), \quad \text{where } k = \frac{(n_1 - 1) \times (n_2 - 1)}{(n_2 - 1) \times C^2 + (1 - C)^2 \times (n_1 - 1)} \quad (21)$$

$$\text{and } C = \frac{s_1^2/n_1}{s_{up}^2}. \quad (22)$$

## Paired data

I don't think I need to refer much to this, however if you are testing whether there is any effect before and then after, consider this to be paired data. In that case, you don't use either pooled or unpooled variance, but look instead at the variance of differences. That is, transform your data from  $x$  and  $y$  to  $d = x - y$ , for instance.

## ANOVA

### One-way

In Table 1 is an incomplete tabulated output for a test comparing the means between groups. Can you fill in the missing values?

	SS	df	MS	F	sig.
G	91.467		45.733		0.021
E	276.400	27			
T	367.867				

Table 1: ANOVA one-way table

Some background on the test: a teacher wants to know if the starting level of her pupils affects the mean length of time to complete the exam. Formulate the null and alternative hypotheses.

$$H_0 : \quad (23)$$

$$H_a : \quad (24)$$

**Summarise your findings of this test:** In Table 1, can have that the sum of squares between the groups (SSG) is 91.467 and that the mean square between groups (MSG) is 45.733. You want to find out how these two are related, and notice that  $45 \times 2 = 90$ , which is indeed the degrees of freedom for groups ( $df_G = 2$ ). Now, you can conclude that you have 3 groups in total ( $I = 3$ ).

Moving on the row marked 'E': there is an evident relationship between the mean squared error within each group (MSE) of 276.400 and the degrees of freedom for the error ( $df_E$ ) of 27. That is,  $df_E \times 10 \approx \text{MSE}$ ; using the calculator, you find that MSE is equal to 10.237. Given that  $df_E (= n - I; I = 3)$  is 27, we know that each group has 30 participants ( $n = 30$ ).

Now that we know  $n$ , and hence  $df_T$  is 29, we can calculate the variance of our data (MST) as 12.685.

In order to calculate our  $F$  statistic, we need to understand exactly what it is:  $F$  is the ratio of variation between and within groups. There are three distinct cases which we will look at now.

$$0 \leq F < 1$$

In this case (refer to Figure 1), we know that the variance in the data which can be explained by the differences between the groups, is less than the variance within the groups themselves. Either, there is not much difference between the groups, or there is a lot of variation within the groups. In both cases, you would need to perform post-hoc tests, such as contrasts, to confirm or deny your findings. Evidently if  $F = 0$ , then there is no variation between the groups, i.e.  $\bar{y}_1 = \bar{y}_2 = \dots = \bar{y}_I = \bar{y}$ .

$$F \approx 1$$

We can note the implications of an  $F$  near to one as 'good'; we conclude that the variance in the data is not solely due to variance between groups but equally within. We call this as 'good' because of the complications that arise when we find significant results (more on this in a moment). This  $F$  tells us that the variations between the groups is proportional to the variance within the groups themselves, so the effects of being in any particular group are not evident in the data. A very simple example of this is test scores between schools: you might want to test whether attending a more prestigious has any outcome on the results of the students themselves. So, you would have  $I$  equal to the number of schools (ranked in order of prestige) and  $J$  equal to the number of students at each school (in a particular graduating year). If you received an  $F$  value close to 1, then you would conclude that there is no significant advantage benefited to

$$F = \frac{\text{MSG}}{\text{MSE}} < 1 \quad (25)$$

$$\implies \frac{\text{SSG}}{\text{df}_G} < \frac{\text{SSE}}{\text{df}_E} \quad (26)$$

$$\implies \text{variance between} < \text{variance within} \quad (27)$$

Figure 1

students who attend a prestigious school, in terms of grades. Additionally, note the following relationship:

$$\left. \begin{aligned} F &= \frac{\text{Var } y - s_p^2}{s_p^2} = \frac{\text{Var } y}{s_p^2} - 1 \\ F &\approx 1 \end{aligned} \right\} \implies \frac{\text{Var } y}{s_p^2} \approx 2. \quad (28)$$

This says that, if the total variation in  $y$  is twice the size of the collective variation within each group, then any observed variation between the groups in our sample is acceptable (accept null hypothesis).

### **F > 1**

In this case, we can conclude that we have a significant result: there is a great difference between the groups. With respect to the previous example, we would conclude that the prestige of a school has an effect on the grades of attendees; if  $F$  is **much** larger than 1, we could say that the effect is **profound** or immense. Within each school, there is not much variation in the data compared with the variation between the schools (e.g.  $I = 3$ ): Figure 2 is a pictorial example of  $F > 1$ , which would lead us to assume that further tests (e.g. contrasts) need to be performed in order to determine which school benefits the greatest advantages to it's attendees.

So, back to our table! We can discern that we will have an  $F > 1$ , as  $10 \times 4.5 = 45$ , and indeed we have  $F = 4.467$  (see the completed Table 2). Relating this to the significance of 0.021 found by the software, the probability of achieving an  $F$  more extreme than  $f_{(2,27|\alpha=0.05)}^* = 3.354$  in any repetition is 0.021, which is less than our  $\alpha = 0.05$ .

The usual method of formulating the null and alternate hypotheses (as you may have figured, by now), is

$$H_0 : \mu_1 = \mu_2 = \mu_3, \text{ i.e. there is no difference between the three groups;} \quad (29)$$

$$H_a : \text{there exists a difference somewhere between the groups.} \quad (30)$$

Using the  $F$  statistic and  $p$ -value, we reject the null hypothesis at a 5% significance level in favour of the alternate hypothesis, and conclude that at least one of the groups is different from the others. In order to determine which group differs the most, we might formulate some contrasts based on our plot of the data, similar to the one in Figure 2. For example if one of the groups, say group three, was distinctly further away from groups one and two, we might decide to find an  $a$ , such that  $0 \leq a \leq 1$  and test:

$$H_0 : \frac{a(\mu_1 + \mu_2)}{2} = (1 - a)\mu_3 \quad (31)$$

If we find that an  $a$  near to zero provides us with a insignificant results (i.e. not reject  $H_0$ ), then we know that group three dominates. However, if an  $a$  near to one has this provision, then we know that a groups one and two equally dominate. Here, dominate means that they greatly contribute to the difference between groups (i.e. large comparable mean/s).

The background on the data set tells us that this is indeed about schooling, but only one particular test (limited factors  $\implies$  ANOVA I, and not ANOVA II). We may conclude that the findings of the test imply that the starting level of the pupils does indeed affect the mean length of time to complete the exam. Without knowing any more information, such as individual group means  $\bar{y}_i$ , standard deviations  $sd_i$  and number per group  $n_i$ , we conclude our testing here. If this information was available to us, we could run some contrasts to find which particular starting level, beginner, intermediate or advanced, provided the biggest advantage on this particular exam.

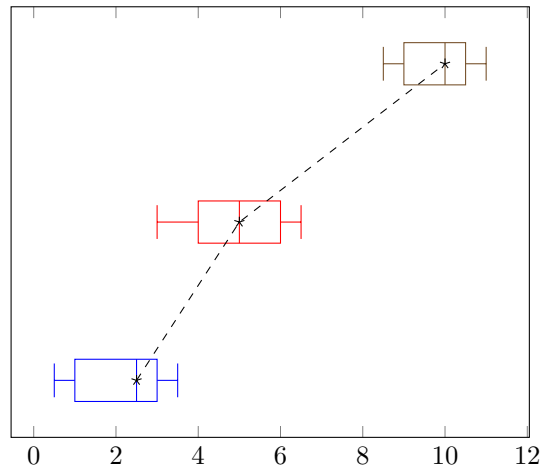


Figure 2: The box plots display the variance **within** groups, which is noticeably small, and the dashed line displays the variance **between** groups, which is noticeably large.

	SS	df	MS	F	sig.
G	91.467	2	45.733	4.467	0.021
E	276.400	27	10.237		
T	367.867	29	12.685		

Table 2: ANOVA one-way table (completed)

## Two-way

Fun time!!!

So, not only do you have two or more groups ( $I$ ), but we now consider that there might be multiple effects ( $A$ ,  $B$ ,  $C$ , ...) and their interactions ( $A \times B$ ,  $A \times B \times C$ , ...). We will start with a simple example, and progress from there:

Perhaps a farmer wants to investigate the effects of manure ( $A$ ) and nitrogen-based fertiliser ( $B$ ), and also the combination of both ( $A \times B$ ), on the yield of their corn. The output of an ANOVA II test is given below:

	SS	df	MS	F	sig.
A	15.842		15.842		0.029
B	17.298		17.298		0.035
$A \times B$	3.872		3.872		0.273
E	48.000		3.000		
T	85.012	19	4.474		

Table 3: My caption

Can you fill in the table and draw some conclusions?

**Summarise your findings:** There is a pretty obvious relationship between the columns “SS” and “MS”, so we can assume that the degrees of freedom for the factors and the interaction effect are 1. Meaning: that the factors  $A$  and  $B$  have two categories each, such as high and low levels (of manure and fertiliser). So, then we know that there are in total  $2 \times 2 = 4$  groups: high levels of both, high levels of manure v.s. low levels of fertiliser, low levels of manure v.s. high levels of fertiliser, and low levels of both. Further, we can examine the calculation of  $df_E$ :

$$df_E = \underbrace{(n-1) - (I-1) - (J-1) - \overbrace{(I-1)(J-1)}^{IJ-I-J+1}}_{df_T - (df_A + df_B + df_{A \times B})} = n - IJ. \quad (32)$$

We already have  $I = 2 = J$ , and  $n = 19 + 1$ , so  $df_E = 20 - 2 \times 2 = 16$ . We remember how to calculate the  $F$

statistic for each factor and interaction:

$$F = \frac{\text{variance between}}{\text{variance within}} = \begin{cases} \frac{MSA}{MSE}, & \text{when testing if } A \text{ has a significant effect;} \\ \frac{MSB}{MSE}, & \text{when testing if } B \text{ has a significant effect;} \\ \frac{MSAB}{MSE}, & \text{when testing if an interaction of } A \text{ and } B \text{ has a significant effect.} \end{cases} \quad (33)$$

When we are talking about ‘significant effect’, we want to know if the variation between the groups is relatively equal to the variation within the groups. So, in our table we can input  $F_A = 15.842/3 = 5.281$ ,  $F_B = 17.298/3 = 5.77$  and  $F_{A \times B} = 3.872/3 = 1.291$ . So, given that our  $F$  statistic for factors  $A$  and  $B$  are well above 1, we can conclude that the levels (each) of manure and fertiliser has a profound affect on corn yield. Further, this is evident by the  $p$ -values both being under our accepted 5% level. Now it becomes a little harder to discern the right answer concerning the existence on an interaction effect: we have that the  $F$  value is above 1.2 and perhaps would think of rejecting  $H_0$ , however we must consider the  $p$ -value! Repetitions of this study would yield more extreme  $F$  values for this interaction effect more than a quarter of time ( $\mathbb{P}(F > f_{(1,16|0.05)}) = 0.273 > 0.25$ ), so we may safely conclude that there is no interaction effect. You can see what (graphically) denotes an effect in the slides from lecture 5 (slide 26-31). If you were to advise the farmer, what advice would you give?

Previously we found that factors  $A$  and  $B$  both had a main effect, but there was no significant interaction effect. Now we are interested in conducting a basic two-way ANOVA without the interaction effect. In order to do this, we include the data from the interaction effect into the error terms:

$$SSE_{\text{new}} = SSE + SSAB \quad df_{E, \text{new}} = df_E + df_A \times B \quad (34)$$

$$\implies MSE_{\text{new}} = \frac{SSE_{\text{new}}}{df_{E, \text{new}}} \quad (35)$$

**Something important to remember for the exam:** the “spread of means” is not equal to the “mean spread of the data”. The former refers to the variance between means (MSG) and the latter, the mean variance (MSE). The professor will use language like this in order to confuse you!!

## Contrasts

Why do we perform contrasts? Simply put, it is to reduce the overall statistical error: if  $\alpha\%$  for each test and you have  $I$  groups then you will need to perform  $I \times (I - 1)/2$  tests, meaning that

$$\text{overall error rate} = \mathbb{P}(\text{at least one false rejection}) \quad (36)$$

$$= \mathbb{P}\left(\text{at least one Type I error} \left| \frac{I \times (I - 1)}{2} \text{ tests} \right.\right) \quad (37)$$

$$= 1 - \mathbb{P}\left(\text{no Type I errors} \left| \frac{I \times (I - 1)}{2} \text{ tests} \right.\right) \quad (38)$$

$$\approx 1 - (1 - \alpha)^{I \times (I - 1)/2} \quad (39)$$

So, for a simple ANOVA I with 5 groups at 1% significance level, this means we need to perform  $5 \times 4/2 = 10$  tests resulting in an error rate of  $1 - 0.99^{10} = 0.0956$ , i.e. 9.6% of a Type I error, which is higher than our accepted 5% level (chance capitalisation). It is possible to plan for this by introducing contrasts **prior to undertaking the test**.

Assume, again,  $I = 5$  and we want to know which group accounts for the largest deviation of the data. **Remember**,  $SST = SSG + SSE$  - this means that the total variance in the data set is either due to the variance between groups or within the groups themselves. If we have an  $F > 1$  (and  $p < \alpha$  given the sample size of each group is “large enough”), we know that it is attributable to variance between groups and now we want to know which particular group/s are different. We have that,

$$\bar{x}_1 > \bar{x}_2$$

Is group 1 different to group 2?

$$\begin{aligned} H_{01} : \mu_1 - \mu_2 &= 0 \\ H_{a1} : \mu_1 - \mu_2 &> 0 \end{aligned} \quad (40)$$

$$\text{coefficients: } (1, -1, 0, 0, 0) \quad (41)$$

$$\begin{array}{ll} \bar{x}_1 < \bar{x}_3 & \\ \bar{x}_2 < \bar{x}_3 & \end{array} \quad \text{Are groups 1 and 2 different to group 3?} \quad \begin{array}{l} H_{02} : \frac{\mu_1 + \mu_2}{2} - \mu_3 = 0 \\ H_{a2} : \frac{\mu_1 + \mu_2}{2} - \mu_3 < 0 \end{array} \quad (42)$$

$$\text{coefficients: } \left( \frac{1}{2}, \frac{1}{2}, -1, 0, 0 \right) \quad (43)$$

$$\bar{x}_5 > \max \{ \bar{x}_1, \bar{x}_2, \bar{x}_3, \bar{x}_4 \} \quad \text{Is group 5 the most different?} \quad \begin{array}{l} H_{03} : \frac{\mu_1 + \mu_2 + \mu_3 + \mu_4}{4} - \mu_5 = 0 \\ H_{a3} : \frac{\mu_1 + \mu_2 + \mu_3 + \mu_4}{4} - \mu_5 < 0 \end{array} \quad (44)$$

$$\text{coefficients: } \left( \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, -1 \right) \quad (45)$$

We don't know the true values of  $\mu_1, \dots, \mu_5$ , so we estimate with  $\bar{x}_1, \dots, \bar{x}_5$  to produce an estimated contrast value  $c$  and it's associated  $t$  statistic:  $c = \sum_{i=1}^I a_i \bar{y}_i$  and  $t = c/\text{SE}_c \sim t(n-I)$  (under  $H_0$ ). Let's look at an example:

Group	$\bar{y}_i$	$s_i$	$n_i$
1	33.31	3.63	9
2	28.72	2.91	15
3	32.46	3.98	9
Total	30.99	3.3	33

Table 4: My caption

Notice that group 2 has the lowest mean and sd, and the largest sample size? This means that we can infer already that group 2 is the most different. First, we should calculate the pooled variance:

$$s_p = \sqrt{\frac{\sum_{i=1}^3 (n_i - 1) s_i^2}{\sum_{i=1}^3 (n_i - 1)}} \quad (46)$$

$$= \sqrt{\frac{8 \times 3.63^2 + 14 \times 2.91^2 + 8 \times 3.98^2}{33 - 3}} \quad (47)$$

$$= \sqrt{\frac{350.69}{30}} \quad (48)$$

$$= \sqrt{11.69} \quad (49)$$

$$= 3.42. \quad (50)$$

So we test,

$$H_{01} : \mu_2 = \mu_1 \quad (51)$$

$$H_{a1} : \mu_2 < \mu_1 \quad (52)$$

$$\implies c = -1 \times \bar{y}_1 + 1 \times \bar{y}_2 + 0 \times \bar{y}_3 = -4.59. \quad (53)$$

$$\implies \text{SE}_c = s_p \sqrt{\sum_{i=1}^3 \frac{a_i^2}{n_i}} \quad (54)$$

$$= 3.42 \times \sqrt{\frac{1}{9} + \frac{1}{15}} \quad (55)$$

$$= 3.42 \times \sqrt{0.178} \quad (56)$$

$$= 3.42 \times 0.422 \quad (57)$$

$$= 1.44. \quad (58)$$

$$\implies t = \frac{c}{\text{SE}_c} \quad (59)$$

$$= \frac{-4.59}{1.44} \quad (60)$$

$$= -3.1875 \sim t(33-3)_{\alpha/2=0.025}. \quad (61)$$

Our critical  $t^*$  value for 32 df and (two-tailed) 5% significance is -2.042, which is larger than our  $t$ -statistic. So we must conclude to reject  $H_0$  in favour of  $H_a$ ; you can read the  $p$ -value from a table as around 0.003. If we were to construct a confidence interval about  $c$ :

$$\text{CI} = (c - t^* \times \text{SE}_c, c + t^{**} \times \text{SE}_c) \quad (62)$$

$$= (-3.1875 - 2.042 \times 1.44, -3.1875 + 2.042 \times 1.44) \quad (63)$$

$$= (-6.13, -0.25) \quad (64)$$

Note that the entire confidence interval is located to the left of zero? So, we certainly reject the null hypothesis, as even when it is assumed, we do not have zero contained in the interval!

## Confidence intervals

The general equation for confidence intervals involving multiple comparisons is given by

$$CI_{ij} = (\bar{y}_i - \bar{y}_j) \pm t^{**} \times s_p \sqrt{\frac{1}{n_i} + \frac{1}{n_j}} \quad (65)$$

For this point of the course, you have two options for your critical  $t^{**}$ :

1. **Bonferroni**: adjust the significance of each test to ensure that overall error rate is less than the specified  $\alpha$ :

$$k = \text{the number of tests} \quad (66)$$

$$= \frac{I \times (I - 1)}{2}, \quad \text{where } I \text{ is the number of groups.} \quad (67)$$

Let  $\alpha^*$  be the significance of **each** of the  $k$  tests, then

$$\alpha^* = \frac{\alpha}{k} \implies t^{**} = t_{1-\alpha/2k}^* \quad (68)$$

$\nu = \text{df}_E$

2. **Least significant differences (LSD)**: we use this for  $I = 3$  groups, otherwise use Bonferroni. This is because LSD method does not alter the significance, but only each individual test is improved, and  $I = k$  for 3 groups.

$$t^{**} = t_{1-\alpha/2}^* \quad (69)$$

$\nu = \text{df}_E$

**Important:**

$$p_{\text{Bonferroni}} = \mathbb{P}(|T| > t_{1-\alpha/2k}) = \frac{\mathbb{P}(|T| > t_{1-\alpha/2})}{k} = \frac{p_{\text{LSD}}}{k} \quad (70)$$

The  $p$ -values are scaled with respect to the number of tests performed.

## Kruskal-Wallis procedure

ANOVA assumptions are normality, homoskedasticity and independence, and if they are severely violated then we turn to the Kruskal-Wallis procedure (non-parametric ANOVA). The null hypothesis is similar to ANOVA, however the violation of assumptions leads us to a new direction: “the **distribution** is the same in all groups”. The alternative is that the scores in some groups are systematically larger.

Begin by ordering all  $n$  scores from lowest to highest, assigning rank 1 to the lowest. If some scores are equal, they receive the mean of the ranked score, e.g. there are two scores of 9 and it is the fifth lowest score (taking up spaces 5 and 6) then each of the ‘9’s receive a rank of  $(5 + 6)/2 = 5.5$ .

Now that each score has a rank, you need to separate the ranked scores back into their respective groups. Sum the rankings in each group to yield  $R_i$ . Furthermore,

1. if the sample sizes  $n_i$  are small ( $<5$ ), then we use the following test statistic in ANOVA I:

$$H = \frac{12}{n(n+1)} \sum_{i=1}^I \frac{R_i^2}{n_i} - 3(n+1) = \frac{12}{n(n+1)} \text{SSG}_{\text{ranks}} \quad (71)$$

2. if the sample sizes are not too small, say  $n_i \geq 5$ , then use the following approximation:  $H \sim \chi^2(I-1)$ .

For example, you are investigating the effects of exercise on depression. So you have three groups: no exercise, 20 minutes of jogging per day, and 60 minutes of jogging per day. In order to simplify, we assume that each participant is equivalently depressed, then at the end of the month you ask each participant how depressed they feel as a score out of one hundred, where 1 is totally miserable and 100 is ecstatically happy.

The method of testing is self-recorded and the scores (ordinal) given are non-parametric, so in order to draw conclusion we will need to use the Kruskal-Wallis procedure.



Table 5: This table contains the self-recorded ratings of 1 to 100 from 27 depressed people, where 1 is totally miserable and 100 is ecstatically happy. The ratings were ask for after performing daily jogging for a month, in groups of 20 and 60 minutes per day. The third group is the control group, who were asked to not exercise every day.

$i = 1, 2, 3;$ $n_i = 8$	no exercise	20 min/day	60 min/day
	23	22	59
	26	27	66
	51	39	38
	49	29	49
	58	46	56
	37	48	60
	29	49	56
	44	65	62
$\bar{x}_i$	<b>39.63</b>	<b>40.63</b>	<b>55.75</b>
$s_i$	<b>12.85</b>	<b>14.23</b>	<b>8.73</b>

Looking at Table 5, the minimum rating given was 22 and the maximum was 66, so 22 is given rank 1 and 66 is given the lowest rank. Tied scores get the average of the rankings they would've received. I've summarised the next few steps in the following ranked table:

Table 6: This table contains the rankings of the scores given, as well as their sums per group.

$i = 1, 2, 3;$ $n_i = 8$	no exercise	20 min/day	60 min/day
	2	1	20
	3	4	24
	16	9	8
	14	5.5	14
	19	11	17.5
	7	12	21
	5.5	14	17.5
	10	23	22
$\bar{x}_i$	<b>39.63</b>	<b>40.63</b>	<b>55.75</b>
$s_i$	<b>12.85</b>	<b>14.23</b>	<b>8.73</b>
$R_i$	<b>76.5</b>	<b>79.5</b>	<b>144</b>

You can see in Table 5 that there is a rating of 49/100 in each column, which is the 14th, 15th and 16th lowest ranking. So they each get the ranking of  $(14 + 15 + 16)/3 = 14$ . As each  $n_i = 8 > 5$ , we assume that the  $H$  test statistic is approximated by the  $\chi^2$  distribution with 2 degrees of freedom. Calculating  $H$ :

$$H = \frac{12}{n(n+1)} \sum_{i=1}^I \frac{R_i^2}{n_i} - 3(n+1) \quad (72)$$

$$= \frac{12}{24(24+1)} \sum_{i=1}^3 \frac{R_i^2}{8} - 3(24+1) \quad (73)$$

$$= \frac{12}{24 \times 25} \left( \frac{76.5^2}{8} + \frac{79.5^2}{8} + \frac{144^2}{8} \right) - 3 \times 25 \quad (74)$$

$$\approx 7.271 \sim \chi^2(2) \quad (75)$$

Looking at the  $p$ -values for a  $\chi^2(2)$ , we can surmise that the probability of obtaining an  $H$  statistic more extreme than what we calculated is between 5% and 2.5%:

$$0.025 \leq \mathbb{P}(H > 7.271) \leq 0.05 \quad (76)$$

Given the small  $p$ -value for the  $H$  statistic, we can conclude that there is some difference between the three groups, i.e. the effects of exercise on depression is evident. Looking back at Table 6, you may notice that

the mean of group 3 is relatively larger than the others, and it's standard deviation is comparatively smaller. Meaning that we may infer that group three is the “different” group: 60 minutes of medium exercise (e.g. jogging) has the greatest impact on depression (in terms of improvement), compared to 20 minutes or no exercise (check this with the next section on effect size ☺).

## Post-hoc methods

Post-hoc multiple comparisons: tests, using confidence intervals, for differences between **all** pairs of means:

$$\text{CI}_{ij} \text{ for comparing the means of groups } i \text{ and } j: (\bar{y}_i - \bar{y}_j) \pm t_{df_E}^{**} s_p \sqrt{\frac{1}{n_i} + \frac{1}{n_j}} \quad (77)$$

Note the use of  $t^{**}$ ? This refers you to the section on LSD and Bonferroni!

All this is about, is two main things: do your CI's overlap; is your assumed value inside the CI.

1. So, if we are comparing groups (comparing means), then we would like our assumed value  $\mu_i - \mu_j = 0$  ( $H_0$ : all groups are the same, i.e. all means are the same) to be inside the confidence interval. If it does not, we might need to look at our test for any errors (e.g. low effect size, small sample size, chance capitulation). However, if there are no errors present (e.g. large effect size,  $n \geq 30$ , Bonferroni method used, testing procedure is widely regarded as effective, etc.) then you would **reject**  $H_0$ . Why? You assume something, then you construct your test in such a way that there is no possible room for error, however the result shows that  $(1 - \alpha)\%$  of all samples lead to an interval that does not cover the unknown parameter. So our assumed parameter must be incorrect ( $\mu_i - \mu_j \neq 0$ ). If our samples collected are large enough, the CLT allows us to conclude that true parameter lies somewhere closer to the estimated value  $a = \bar{y}_i - \bar{y}_j$ .
2. What does it mean if the CI's overlap? Referring back to Figure 2, you can see the boxplots represent variation within the groups, and that the red and blue boxplots overlap whilst the other box stands alone. Rather than adding another figure here, imagine that the boxplots are displaying 95% confidence intervals. So you have that the confidence intervals of groups one and two overlap, whilst group three's confidence interval is completely separated. This means that it is likely that the population means for groups one and two are equal, and the pop.mean of group three is likely greater. In relation to (77), if more than two CI's overlap sequentially (e.g. the upper bound of the CI for  $\mu_1 - \mu_2$  is contained in the CI for comparing  $\mu_2 - \mu_3$ , and the upper bound for the CI comparing  $\mu_2 - \mu_3$  is contained in the CI for comparing  $\mu_3 - \mu_4, \dots$ ), then you might assume that there it is likely that the groups are all (relatively) the same. How much they overlap is indicative of how likely it is that the groups are all the same. For example, if the upper bound of the CI for  $\mu_1 - \mu_2$  is more than the estimate for comparing  $\mu_2 - \mu_3$ , namely the midpoint  $\bar{y}_2 - \bar{y}_3$ , then you might assume that these groups all differ by the same amount (more than relatively). If the CI's (almost) coexist, then you can be sure that these groups certainly differ by the same amount - if zero is included in all the CI's, then the groups are the same (no difference).

## Effect size

What is “effect size” in the relation of statistics? It's an objective and standardised way to determine if there is indeed an observable effect in the data. The usual method is by way of ratio, so if the ratio is on the large end of the scale then you can say that the researcher will notice an effect by the factors/groups in the data. Another way to think about it, is to regard effect size as the statistical ‘yard stick’; “relative to the size of my hand, how big is this ant? So, will I notice it (see it) crawling onto my hand?”.

Whilst writing this, I came across a really great website that explains effect size in lay-mans terms. I think it is helpful to read, in order to give a ‘layman's response’ to a psychological research question. It may also help you to visualise effect size: <https://www.theanalysisfactor.com/effect-size/>.

## Cohen's $d$

Cohen's  $d$  is used to measure the standardised difference between means in

1. one random sample drawn from a normally distributed population ( $y \sim \mathcal{N}(\mu, \sigma^2)$ );

$$d = \frac{\bar{y} - \mu}{\sigma} \quad (78)$$

This is a  $z$ -score, and the old “68-95-99.7” rule gives an indication of how we might view this: if you have a  $z$ -score of less than 1, then you know that your estimate based on the data,  $\bar{y}$ , lies in the middle 68%.

Let's consider only  $|d|$  and assume  $|d| \leq 0.2$  (small effect size - see slide 15 from lecture 2) - checking the  $z$ -table we conclude that there is less than 16% of the population, which is less extreme than our sample estimate ( $\mathbb{P}(|D| < 0.2) = 0.15852$ ). The key point here is that this is a **standardised score**, so you can use it to measure the small thing across different groups of populations. The simplest example is comparing the group means for different classes within a school for the same test. If you know the results for the test is normally distributed with mean  $\mu$  and standard deviation  $\sigma$ , you can compare whether the means from different classes are noticeably different from what you expect ( $\mu$ ).

- two random samples drawn from normally distributed populations

$$d = \frac{\bar{y}_1 - \bar{y}_2}{s_p}. \quad (79)$$

**N.B.:** the use of pooled variance implies that we assume homoskedasticity; if this is true, then  $s_p$  is the best approximation for the true standard deviation  $\sigma$ .

Similar to the above, however we assume that  $\mu_1 = \mu_2$  and so  $\mu_1 - \mu_2 = 0$  is “invisible” in the equation. This is almost a  $t$ -statistic (larger standard error due to the omission of  $\sqrt{1/n_1 + 1/n_2}$  in the denominator), so it gives the standardised difference of sample means between two (assumed equal) groups. This will tell you if the standardised difference (effect size) is small ( $\leq 0.2$ ), medium ( $\approx 0.5$ ) or large ( $\geq 0.8$ ).

## Eta squared

$\eta^2$  is the proportion of the total sample variance explained by the effect ( $A, B, A \times B, \dots$ ).

$$\eta^2 = \frac{SS_{\text{effect}}}{SST} \quad (80)$$

Again, it's basically a ratio: if the variation in the data is due mostly to the effect, then this ratio is near (or more) 0.14. In ANOVA II, “effect” can be factor  $A$ , interaction effect  $A \times B$ , etc. In ANOVA I, the only “effect” we consider is group variation ( $SS_{\text{effect}} = SSG$ ) so  $\eta^2 = R^2$ , where  $R^2 = [\text{Cov}(\hat{y}, y)] / s_y^2$  is the squared ratio of covariance and variance of the data  $y$  and the predicted model  $\hat{y}$  - **percentage of variance explained by the model** (see the section on regression). So for ANOVA I (not considering regression models),  $\eta^2 \times 100$  is the percentage of variation in the data as explained by the variation between groups.

### Advantages:

- effects are additive for balanced groups, i.e.  $n_1 = n_2 = \dots = n_I$ :

$$\sum_{\text{all effects}} SS_{\text{effect}} = SSM \quad (81)$$

### Disadvantages:

- $\eta^2$  depends on the number and size of the remaining effects. So if you have a lot of factors/interactions, or if one factor/interaction contributes the greatest, then you might have an  $\eta^2$  which is small, despite the effect size being (actually) medium, for instance:

$$\eta^2 = \frac{SS_A}{SST} \quad (82)$$

$$= \frac{SS_A}{\sum_{\text{all effects}} SS_{\text{effect}} + SS_E} \quad (83)$$

$$= \frac{SS_A}{SS_A + SS_B + SS_C + \dots + SS_{A \times B} + \dots + SS_{A \times B \times C} + \dots + SS_E} \quad (84)$$

$$(85)$$

- $\eta^2$  does not estimate the proportion of variance accounted for in the **population** - it is a **biased estimator**. It always overestimates the explained variance in the population, despite being ‘good’ for the sample.

## Partial eta squared

This tries to solve the disadvantages of  $\eta^2$ , by restricting the ratio to only the effect you are interested in. The disproportional affect of adding effects to the denominator are cancelled out, however there is still the same problem of it being a biased estimator. Additionally, like  $\eta^2$ ,  $\eta_p^2 = R^2$  in ANOVA I.

$$\eta_p^2 = \frac{SS_{\text{effect}}}{SS_{\text{effect}} + SSE} \quad (86)$$

### Advantages:

- $\eta_p^2$  does not depend on the remaining effects, like  $\eta^2$  does. This is because the denominator omits the explained variation from other effects, and focuses solely on the ratio of explained variation (of the effect in question) vs. the combination of explained and unexplained, but restricted to a particular effect.

### Disadvantages:

- Effects are no longer additive for balanced designs (see the section on  $\eta^2$  for further explanation of this).
- $\eta_p^2$  only estimates the proportion of variance accounted for in the sample, and not in the population. Thus, it overestimates the population effects.

**N.B.:** SPSS only outputs  $\eta_p^2$ , so in the exam if you see an SPSS output with  $\eta^2$ , you can be sure which it is ☺.

## Omega squared

$\omega^2$  aims to **exactly** estimate population effects, rather than sample effects, so it is an unbiased estimator.

$$\omega^2 = \frac{SS_{\text{effect}} - df_{\text{effect}} \times MSE}{MSE + SST} \quad (87)$$

### Advantages:

- $\omega^2$  does not overestimate population effects, like  $\eta^2$  and  $\eta_p^2$  do.

### Disadvantages:

- Effects are no longer additive for balanced designs (see the section on  $\eta^2$  for further explanation of this). Furthermore,  $\omega^2$  can be negative! What does this mean?

$$\omega^2 < 0 \quad (88)$$

$$\iff SS_{\text{effect}} - df_{\text{effect}} \times MSE < 0 \quad (89)$$

I should note here that the symbol  $\iff$  means ‘if and only if’ and denotes equivalency. So in this case, ‘ $\omega^2$  is negative’ is equivalent to:

$$SS_{\text{effect}} < df_{\text{effect}} \times MSE \quad (90)$$

$$\iff \frac{SS_{\text{effect}}}{df_{\text{effect}}} = MS_{\text{effect}} < MSE \quad (91)$$

$$\iff F = \frac{MS_{\text{effect}}}{MSE} < 1. \quad (92)$$

I found a great website which is written in a fun and bitchy tone about the benefits of which effect size to use: <http://daniellakens.blogspot.com/2015/06/why-you-should-use-omega-squared.html>. I hope you have a laugh, too!

## Power

We will run quickly over statistical power: **power is how well your test works**. You, as a psychologist/statistician, want to ensure that you only make necessary changes to society and do so at the right time. What is the right time? When the data says you have to!

Consider that you are running a test on whether a new psychological disorder should be included in the new DSM. You structure your test as a hypothesis test at a significance level of  $\alpha = 0.015$  (one-tailed, two-sample  $t$ -test), and find a  $p$ -value which is smaller than your acceptance level.

$$H_0 : \mu = \mu_1 - \mu_2 = 0 \quad (93)$$

$$H_a : \mu > 0 \quad (94)$$

If you have already rejected your null hypothesis, that means your sample statistic  $t$  is more extreme than the accepted value for a particular significance level **and** degrees of freedom:

$$t = \frac{\bar{x}_1 - \bar{x}_2 - \mu_0}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} > t_{\alpha=0.015, \nu=n_1+n_2-2}^* \quad (95)$$

$$\implies \bar{x}_1 - \bar{x}_2 > \mu_0 + t_{\alpha=0.015, \nu=n_1+n_2-2}^* \times s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \quad (96)$$

Suppose you later find out a closer estimate that  $\mu_0 = 0$  for the difference in means:  $\mu_a > 0$ . What is probability that you correctly reject  $H_0$ , given that you know  $H_a$ ?

$$\implies \mathbb{P} \left( \bar{x}_1 - \bar{x}_2 > \mu_0 + t_{\alpha=0.015, \nu=n_1+n_2-2}^* \times s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \mid \mu_a \right) = \mathbb{P} \left( \bar{x}_1 - \bar{x}_2 - \mu_a > \mu_0 - \mu_a + t_{\alpha=0.015, \nu=n_1+n_2-2}^* \times s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right) \quad (97)$$

$$= \mathbb{P} \left( \frac{\bar{x}_1 - \bar{x}_2 - \mu_a}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} > \frac{\mu_0 - \mu_a + t_{\alpha=0.015, \nu=n_1+n_2-2}^* \times s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \right) \quad (98)$$

$$= \mathbb{P} \left( T > \frac{\mu_0 - \mu_a}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} + t_{\alpha=0.015, \nu=n_1+n_2-2}^* \right) \quad (99)$$

$$= 1 - \beta = \text{power}. \quad (100)$$

You can find the power of the test using any population, if you know  $\mu_a$  (or a value which is extremely close to it - can never be 100% certain in science!).

## Correlation and regression

Remember! correlation and regression are two different things, despite being heavily interrelated. Correlation is the measurement of the association between two (or more) variables, whereas regression uses the information provided by the association to predict or extrapolate data.

### Pearson's $r$

Pearson's  $\rho$  is the measure of a linear association **in a population** between two random variables  $X$  and  $Y$ , given by

$$\rho_{X,Y} = \frac{\text{Cov}(X,Y)}{\sigma_X \times \sigma_Y}, \quad (101)$$

where  $\sigma_X$  and  $\sigma_Y$  are the standard deviations of  $X$  and  $Y$ , respectively. Researchers do not always have the luxury of knowing the magnitude or direction of a linear relationship in a population, so they must estimate it using sample data. Pearson's  $r$  approximates  $\rho$ , so that we may draw inferences about the population.

$$r_{x,y} = \frac{\text{Cov}(x,y)}{s_x \times s_y} \quad (102)$$

We use subscript  $x,y$  to denote the variables which  $r$  is measuring the relationship of. We substitute in the equation for sample covariance:

$$= \frac{\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}}{s_x \times s_y} \quad (103)$$

We can rearrange the denominators so that the standard deviations for  $x$  and  $y$  are grouped with their respective expressions:

$$= \frac{\sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \times \left( \frac{y_i - \bar{y}}{s_y} \right)}{n-1} \quad (104)$$

You may notice the familiar expressions for calculating the  $z$  scores for  $x$  and  $y$  for a particular  $i$ :

$$= \frac{\sum_{i=1}^n z_{x_i} \times z_{y_i}}{n-1} \quad (105)$$

So, we have that Pearson's  $r$  is the expected value of the product of  $z$ -scores, where the degrees of freedom is  $n - 1$ . If we want to see this equation in another way, we can return to the first expression, and elaborate:

$$r_{x,y} = \frac{\text{Cov}(x, y)}{s_x \times s_y} \quad (106)$$

$$= \frac{\frac{\sum_{i=1}^n (x_i - \bar{x}) \times (y_i - \bar{y})}{n-1}}{s_x \times s_y} \quad (107)$$

$$= \frac{\frac{\sum_{i=1}^n (x_i - \bar{x}) \times (y_i - \bar{y})}{n-1}}{\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}} \times \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}} \quad (108)$$

The  $n - 1$  in all of the denominators equate out to 1 ('cancel out'):

$$= \frac{\sum_{i=1}^n (x_i - \bar{x}) \times (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \times \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (109)$$

Next, we expand the brackets in the denominator:

$$= \frac{\sum_{i=1}^n (x_i y_i - \bar{x} y_i - x_i \bar{y} + \bar{x} \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \times \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (110)$$

Using that  $n\bar{x} = \sum x_i$ , and similarly for  $y$ , results in the following expression for  $r$ :

$$= \frac{\sum_{i=1}^n x_i y_i - n \times \bar{x} \bar{y}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \times \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (111)$$

$r$  is bounded between  $-1$  and  $1$ , where an  $r = 1$  represents a positive correlation between  $x$  and  $y$ :  $x \propto y$  they are proportional, e.g.  $y = a + bx$  ( $a, b$  are constants with  $b > 0$ ); alternatively,  $r = -1$  represents a negative correlation between  $x$  and  $y$ :  $\frac{1}{x} \propto y$  they are inversely proportional, e.g.  $y = a - bx$ .

If  $r \approx 0$ , then there is no linear relationship between your variables  $x_i$  and  $y_i$ : **they are linearly independent**.

## Simple linear regression (SLR)

Now that you know what Pearson's  $r$  is, you can begin learning about simple linear regression and how they interrelate. You collect a data sample from a population, and after calculating Pearson's  $r$  for your sample you conclude that  $x$  and  $y$  are correlated (later you will want to know if they are correlated in the population - Fisher  $Z$  transformation). Now, you wish to predict future outcomes associated with your data sample, and so you construct a model

$$\underbrace{y_i}_{\text{data}} = \underbrace{\beta_0 + \beta_1 x_i}_{\mu_{y_i} \text{ model}} + \underbrace{\epsilon_i}_{\text{error}} \quad (112)$$

for your **population** which predicts the value  $y_i$  given your input  $x_i$ , based on the constants  $\beta_0$  and  $\beta_1$ . Given your input  $x$ , the  $y$  values are normally distributed with population mean  $\mu_y$  and variance  $\sigma^2$ , and the error term  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$  is independent of  $x_i$ . This is an incredibly important assumption, as you do not reasonably expect that your model can capture/test all possible contributions to the data. What you cannot test/capture is  $\epsilon_i$  - typical examples may be infant life experiences or innate ability. Also, you expect that your model is 'good', so you assume that the mean of  $\epsilon_i$  is zero.

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad (113)$$

$$= \mu_{y_i} + \epsilon_i \sim \mathcal{N}(\mu_{y_i}, \sigma^2) \quad (114)$$

Notice that the mean of  $y_i$  also has subscript  $i$ ? This is because the mean for a particular value  $y_i$  is dependent on the input  $x_i$ , so it is really  $\mu_{y \text{ given } x_i}$ . If we rearrange this expression in favour of  $\epsilon_i$ , the population error term, we find that it is normally distributed with mean zero and the same variance as  $y$ :

$$\implies y_i - \mu_{y_i} = \epsilon_i \sim \mathcal{N}(0, \sigma^2). \quad (115)$$

You estimate  $\beta_0$  and  $\beta_1$  (population coefficients) using your sample, where  $b_0$  and  $b_1$  are the **ordinary least squares estimates** (OLS estimates) of  $\beta_0$  and  $\beta_1$ . They are called OLS estimators because they minimise the sum of squared errors between the actual data  $y_i$  and the predicted data  $\hat{y}_i$ .

$$\min \left( \sum_{i=1}^n e_i^2 \right) = \min \left( \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right) \quad (116)$$

$$= \min \left[ \sum_{i=1}^n (y_i - (b_0 + b_1 x_i))^2 \right] \quad (117)$$

The  $b_0$  and  $b_1$  which satisfy this are:

$$b_0 = \bar{y} - b_1 \bar{x}; \quad b_1 = r_{x,y} \times \frac{s_y}{s_x} = \frac{\text{Cov}(x, y)}{s_x^2} \quad (118)$$

where  $r_{x,y}$  is the Pearson's  $r$  for your data sample,  $s_y$  is the standard deviation of the sample for output  $y$ ,  $s_x$  is the standard deviation of the sample for the input  $x$ , and  $\text{Cov}(x, y)$  is the covariance of  $x$  and  $y$ .  $b_0$  is your intercept, i.e. the value of  $y_i$  when  $x_i = 0$ , so this just tells you the minimum of the range for your predicted output. Under  $H_0$ , you assume  $\beta_0$  to be zero, which is what you will test: is my sample coefficient  $b_0$  significantly different from zero? Note that if your  $b_0$  is forced to be equal to zero, then the regression line does not run through the centre mass point of  $(x_i, y_i) = (\bar{x}, \bar{y})$ :

$$\hat{y}_i = b_0 + b_1 x_i \quad (119)$$

$$= (\bar{y} - b_1 \bar{x}) + b_1 x_i \quad (120)$$

$$= \bar{y} + b_1 (x_i - \bar{x}) \quad (121)$$

$$= \bar{y}. \quad (122)$$

If we have a predicted output  $\hat{y}_i$  equal to the sample mean, there are two possibilities; either,

$$\begin{cases} b_1 = 0: & \text{This implies that } r_{x,y} = 0, \text{ and so } x \text{ and } y \text{ are} \\ & \text{linearly independent.} \\ x_i = \bar{x}: & \text{This implies that the regression line passes} \\ & \text{through the point } (\bar{x}, \bar{y}) \text{ if and only if } b_0 \text{ is} \\ & \text{not equal to zero.} \end{cases} \quad (123)$$

It is also cool to note:

$$\hat{y}_i = (\bar{y} - b_1 \bar{x}) + b_1 x_i \quad (124)$$

$$= \bar{y} + b_1 (x_i - \bar{x}) \quad (125)$$

$$= \bar{y} + \left( r_{x,y} \times \frac{s_y}{s_x} \right) \times (x_i - \bar{x}) \quad (126)$$

$$\implies \underbrace{\frac{\hat{y}_i - \bar{y}}{s_y}}_{\text{z-score for } \hat{y}_i} = r_{x,y} \times \underbrace{\frac{x_i - \bar{x}}{s_x}}_{\text{z-score for } x_i} \quad (127)$$

$b_1$  is your slope coefficient: when my  $x_i$  increases, does my  $y_i$  increase or decrease and at what rate? You can see from (127) that  $r_{x,y}$  is the slope of the regression line of the standardized data points. If your variables are strongly linearly correlated, such as  $r = 1$ , then  $b_1 \approx s_y/s_x$  which is the ratio of standard deviations between your independent and dependent variables; if  $r = -1$ , then  $b_1 \approx -s_y/s_x$ , meaning that  $\hat{y}_i$  will decrease as  $x_i$  increases, respective to the rate  $s_y/s_x$ . It is important to note that the sign of  $b_1$  (i.e. negative or positive) is strictly governed by the sign of the covariance as  $s_y$  and  $s_x$  are always strictly greater than zero (equal to zero if and only you have a single data point for your sample). So  $r$  not only tells you about the observed relationship in the sample, but what to expect of your predictive regression line.

So, now you have your prediction model  $\hat{y}_i = b_0 + b_1 x_i + e_i$ , you can calculate the (squared) **standard error of the estimate**, which estimates  $\sigma^2$ , the variance of the population  $y$ :

$$s^2 = \frac{\sum_{i=1}^n e_i^2}{\underbrace{n-2}_{\text{df}}} \quad (128)$$

$$= \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2} \quad (129)$$

If you are asked, “what is the percentage of explained variance ...” then you know that they are asking you for  $R^2$  (or adjusted  $R^2$ ). If the regression model is linear, then  $R^2 = r^2$  and the percentage of explained variance is  $R^2 \times 100\%$ .

$$R^2 = 1 - \frac{\text{SS}_{\text{residuals}}}{\text{SS}_{\text{total}}} \quad (130)$$

Recall that,

$$SS_{\text{total}} = \sum_{i=1}^n (y_i - \bar{y})^2 \quad (131)$$

$$= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (132)$$

$$= SS_{\text{regression}} + SS_{\text{residuals}}. \quad (133)$$

$$\implies R^2 = \frac{SS_{\text{total}} - SS_{\text{residuals}}}{SS_{\text{total}}} \quad (134)$$

$$= \frac{SS_{\text{regression}}}{SS_{\text{total}}} \quad (135)$$

$$= \frac{\text{explained variance}}{\text{total variance}} \quad (136)$$

If you have more than one independent variable, e.g.  $y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \epsilon_i$ , then you consider adjusting for the increase in parameters:

$$\bar{R}^2 = 1 - (1 - R^2) \cdot \frac{n-1}{n-p-1} \quad (137)$$

$$= 1 - \frac{SS_{\text{residuals}}}{SS_{\text{total}}} \cdot \frac{n-1}{n-p-1} \quad (138)$$

$$= 1 - \frac{SS_{\text{residuals}}/(n-p-1)}{SS_{\text{total}}/(n-1)} \quad (139)$$

Comparing the equation for  $R^2$  and  $\bar{R}^2$  (adjusted  $R^2$ ), the ‘adjustment’ is the degrees of freedom for the residuals.  $R^2$  assumes that  $\text{df}_{\text{residuals}} = n - 1$ , which leads to a biased estimate for the population variance of the residuals. By adjusting the degrees of freedom to  $n - p - 1$ , where  $p$  is the number of independent variables (not including constant  $\beta_0$ ), we gain an unbiased estimate. This will be further explored during multivariate regression.

If you are asked, “what is the variability about the line of regression”, then you know that they are asking you for the standard error of the estimate  $s = \sqrt{\sum e_i^2 / (n - 2)}$ .

If you are asked to construct a  $(1 - \alpha)\%$  CI for the population regression parameter  $\beta_0$  (or  $\beta_1$ ), then you know that they are asking you for a  $t$ -statistic based CI (see Table 7):

Table 7

	$\beta_0$	$\beta_1$
<b>Estimate</b>	$b_0 = \bar{y} - b_1 \bar{x}$	$b_1 = r_{x,y} \times \frac{s_y}{s_x}$
<b>Mean of estimate</b>	$\mu_{b_0} = \beta_0$	$\mu_{b_1} = \beta_1$
<b>CI</b>	$b_0 \pm t_{n-2}^* \times SE_{b_0},$	$b_1 \pm t_{n-2}^* \times SE_{b_1},$
	where $SE_{b_0}$ approximates $\sigma_{b_0}$ and is computed by SPSS	where $SE_{b_1}$ approximates $\sigma_{b_1}$ and is computed by SPSS
<b>Test</b>	$H_0: \beta_0 = 0$ vs $H_a: \beta_0 \neq 0$	$H_0: \beta_1 = 0$ vs $H_a: \beta_1 \neq 0$
<b>Statistic</b>	$t = \frac{b_0 - \beta_0^{=0}, \text{ under } H_0}{SE_{b_0}} \sim t(n-2)$	$t = \frac{b_1 - \beta_1^{=0}, \text{ under } H_0}{SE_{b_1}} \sim t(n-2)$
<b>Assumptions</b>	$b_0 \sim \mathcal{N}(\beta_0, \sigma_{b_0}^2)$ $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$	$b_1 \sim \mathcal{N}(\beta_1, \sigma_{b_1}^2)$

In general, we are most interested in the results of testing done on the true parameter  $\beta_1$ : if we do not reject  $H_0$ , then we conclude that  $x$  and  $y$  do **not** have a linearly correlated relationship in the population (could be independent variables or another relationship is more appropriate, e.g. semi-logarithmic).



## SLR pop-quiz

1. Why do we perform (linear) regression?
2. Based on an SPSS output, what inferences can we make?
3. When might we examine the residual plot, and what inferences might we draw from it?
4. What is the “point of centre mass” on a scatter plot?
5. What tests might we perform on the output of a (linear) regression?
6. What is  $R^2$ , and why/when do we adjust it?
7. If we compute the Pearson’s correlation coefficient, what can we infer from this statistic? Think about:
  - scatterplot  $(x_i, y_i)$ ;
  - residual plot  $(y_i, \hat{y}_i)$ ;
  - roles in regression equation  $\hat{y}_i = b_0 + b_1 x_i$ .

## Fisher Z-transformation

In the previous sections, we looked at the population model and tested goodness of fit for a linear model. In order to further investigate the existence of a linear relationship between  $x$  and  $y$  in the population, we look distinctly at the population correlation coefficient, Pearson’s  $\rho$ , and its sample counterpart, Pearson’s  $r$ . The following  $t$  statistic is used to test  $H_0: \rho = 0$  (linear independence) against  $H_a: \rho \neq 0$ .

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \sim t(n-2) \quad (140)$$

The problem we face is when we reject the null hypothesis, and consider population correlation coefficients which are not equal to zero -  $H_0: \rho = a$  where  $a$  is a non-zero constant against  $H_a: \rho \neq a$ . Under this assumption,  $r$  is not normally distributed, so how do we construct a confidence interval around  $r$  for  $\rho$ ? The answer is the Fisher  $Z$  transformation to  $r_z$  (approximately normal):

$$r_z = \frac{1}{2} \log \left( \frac{1+r}{1-r} \right) \sim \mathcal{N} \left( \rho_z = \frac{1}{2} \log \left( \frac{1+\rho}{1-\rho} \right), \sigma_{r_z}^2 = \frac{1}{n-3} \right) \quad (141)$$

This says that our transformed sample correlation coefficient,  $r_z$ , is normally distributed with mean  $\rho_z$ , which is the transformed population correlation coefficient (just insert  $\rho$  in place of  $r$  in the equation for  $r_z$ ), and standard deviation  $1/\sqrt{n-3}$ .

$$Z = \frac{r_z - \rho_z}{1/\sqrt{n-3}} \sim \mathcal{N}(0, 1) \quad (142)$$

If we want to find out if  $r_z$  is significant, i.e. the probability of achieving a more extreme value, we can do a  $z$ -test:

$$\mathbb{P} \left( |Z| > \frac{r_z - \rho_z}{1/\sqrt{n-3}} \right) = p \leftarrow \text{look this up from the } z\text{-table.} \quad (143)$$

The above  $z$ -test is two tailed, because our alternative hypothesis is that  $\rho \neq a$ , so  $\rho < a$  or  $\rho > a$ . Now we can construct a  $(1 - \alpha)\%$  confidence interval about  $r_z$  for  $\rho_z$ , using the critical  $z$ -value.

$$\implies \text{CI for } \rho_z: \quad r_z \pm z^* \frac{1}{\sqrt{n-3}} \quad (144)$$

So we have an upper and lower bound now for our confidence interval for  $\rho_z$ :

$$\text{LB}_{\rho_z} = r_z - z^* \frac{1}{\sqrt{n-3}} \quad \text{UB}_{\rho_z} = r_z + z^* \frac{1}{\sqrt{n-3}} \quad (145)$$

This does not tell us anything about  $\rho$ , only  $\rho_z$ , so in order to have a confidence interval for  $\rho$  we can transform  $r_z$  back to  $r$  using:

$$r_z = \frac{1}{2} \log \left( \frac{1+r}{1-r} \right) \quad (146)$$

$$\implies 2r_z = \log\left(\frac{1+r}{1-r}\right) \quad (147)$$

$$\implies e^{2r_z} = \frac{1+r}{1-r} \quad (148)$$

$$\implies r = \frac{e^{2r_z} - 1}{e^{2r_z} + 1} \quad (149)$$

The above is the Inverse Fisher  $Z$  transformation. We can now construct a confidence interval for our population correlation parameter.

$$\implies \text{CI for } \rho: \left( \frac{e^{2\text{LB}_{\rho_z}} - 1}{e^{2\text{LB}_{\rho_z}} + 1}, \frac{e^{2\text{UB}_{\rho_z}} - 1}{e^{2\text{UB}_{\rho_z}} + 1} \right) \quad (150)$$

What do we know about confidence intervals? Well if we have constructed them well (effect size, sample size  $n$ , etc.) and the hypothesised population parameter is not contained in the interval **then we reject the null hypothesis**. So if we get a CI for  $\rho$ , and  $\rho_0$  is not included in the interval, then we must reject  $\rho = \rho_0$ , and look for an alternative value for  $\rho$ .

$$\begin{array}{ccc} r & \Rightarrow & b_0 \text{ and } b_1 \\ \Downarrow & & \\ r_z & \Rightarrow & \text{CI for } \rho_z \\ \Downarrow & & \Downarrow \\ p\text{-value for } r_z & & \text{CI for } \rho \end{array} \quad (151)$$

## What, when and why (assumptions)

The usual question that statisticians ask is, “what formula do I need to use for this?” and hopefully Tables 8 to 10 will help to clear that up.

Usually, we start by talking about some random variable  $y$  which has an approximately normal population distribution, i.e.  $y \sim \mathcal{N}(\mu, \sigma^2)$ . If the population of  $y$  is **not** distributed normally, the Central Limit Theorem allows us to conclude that the sampling distribution of the mean of  $y$  (many many samples) is! So even if the distribution of the population of  $y$  is skewed, or has no observable pattern, the CLT states that the mean of all possible samples (i.e. large  $n$ ) has an observable pattern, namely normal. So if  $y$  has mean  $\mu$  and standard deviation  $\sigma$ , then  $\bar{y} \sim \mathcal{N}(\mu, \sigma^2/n)$  is the sampling distribution of the mean of  $y$ .

Table 8

Number of groups (I)	Assumptions	Name	Statistic	Confidence interval
1	Normality CLT $\mu$ known $\sigma$ known independence	$z$ -score standardised score	$z = \frac{y - \mu}{\sigma} \sim \mathcal{N}(0, 1)$	
1	CLT $\mu$ known $\sigma$ <b>unknown</b> independence	One-sample independent $t$ -test for estimating population mean	$t = \frac{\bar{y} - \mu}{s/\sqrt{n}} \sim t(n-1)$	CI for $\mu$ : $\bar{y} \pm t^* \frac{s}{\sqrt{n}}$
2	Normality CLT independence $\mu_1, \mu_2$ unknown $\sigma_1, \sigma_2$ unknown <b>Homoskedasticity</b> : Homogeneity of variances $\sigma_1 \approx \sigma_2$ $H_0 : \mu = \mu_1 - \mu_2 = \mu_0 = 0$	Two-sample independent $t$ -test for comparing means (assumed equal variance)	$t = \frac{(\bar{y}_1 - \bar{y}_2) - \mu_0}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t(n-2)$	CI for $\mu = \mu_1 - \mu_2$ : $(\bar{y}_1 - \bar{y}_2) \pm t^* s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$
2	Normality CLT <b>dependence</b> $\mu_1, \mu_2$ unknown $\sigma_1, \sigma_2$ unknown <b>Homoskedasticity</b> : Homogeneity of variances $\sigma_1 \approx \sigma_2$ Construct a new variable $d_i = y_i - x_i$ (difference) from the dependent sample $(x_i, y_i)$ Equal sample sizes $n_x = n = n_y$ $H_0 : \mu_d = \mu_0 = 0$	Two-sample <b>dependent</b> $t$ -test for comparing means ( <b>paired data</b> , e.g. before and after a treatment)	$\bar{d} = \frac{\sum_{i=1}^n (y_i - x_i)}{n} = \bar{y} - \bar{x}$ $s_d^2 = \frac{\sum_{i=1}^n d_i^2 - \bar{d}^2/n}{n-1} \implies t = \frac{\bar{d} - \mu_d}{s_d/\sqrt{n}} \sim t(n-1)$	CI for $\mu_d = \mu_{0y} - \mu_{0x}$ : $\bar{d} \pm t^* \frac{s_d}{\sqrt{n}}$
2	Normality CLT independence $\mu_1, \mu_2$ unknown $\sigma_1, \sigma_2$ unknown <b>Homogeneity of variances violated</b> $\sigma_1 \neq \sigma_2$ $H_0 : \mu = \mu_1 - \mu_2 = \mu_0 = 0$	Two-sample independent $t$ -test for comparing means (assumed unequal variance)	$t = \frac{(\bar{y}_1 - \bar{y}_2) - \mu_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \sim t(k)$ , where $k$ is approximated by a computer	CI for $\mu = \mu_1 - \mu_2$ : $(\bar{y}_1 - \bar{y}_2) \pm t^* \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$

Table 9

ANOVA I	
Number of groups	3+
Assumptions	<p>Normality CLT independence <math>\mu_i, i = 1, 2, \dots, I</math>, unknown <math>\sigma_i, i = 1, 2, \dots, I</math>, unknown <b>Homoskedasticity:</b> Homogeneity of variances <math>\sigma_1 = \sigma_2 = \dots = \sigma_I</math> <math>H_0: \mu_1 = \mu_2 = \dots = \mu_I</math> all of the groups are the same <math>H_a</math>: one of the groups is different from the rest</p>
Name	<p>ANOVA I “is the observed variance in the data attributable to the variation between the groups, or within the groups”</p> <p>Compare the means of 3 or more groups; if 2 groups, then <math>t</math>-test suffices</p>
Statistic	$  \left. \begin{aligned}  &i = 1, 2, \dots, I \text{ indexes the } \mathbf{group\ number} \\  &j = 1, 2, \dots, n_i \text{ indexes the } \mathbf{score} \text{ within a group} \\  &\text{SST} = \text{SSG} + \text{SSE} \\  &\text{df}_T = \text{df}_G + \text{df}_E \\  &\text{MST} = \frac{\text{SST}}{\text{df}_T} = \text{Var } y \\  &\text{SST} = \sum_{i=1}^I \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2 \\  &\text{df} = n - 1 \\  &\text{MSG} = \frac{\text{SSG}}{\text{df}_G} \\  &\text{SSG} = \sum_{i=1}^I \sum_{j=1}^{n_i} (\bar{y}_i - \bar{y})^2 = \sum_{i=1}^I n_i (\bar{y}_i - \bar{y})^2 \\  &\text{df}_G = I - 1 \\  &\text{MSE} = \frac{\text{SSE}}{\text{df}_E} = s_p^2 \\  &\text{SSE} = \sum_{i=1}^I \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 = \sum_{i=1}^I (n_i - 1) s_i^2 \\  &\text{df}_E = n - I  \end{aligned} \right\} \implies F = \frac{\text{MSG}}{\text{MSE}} \sim F(\text{df}_G, \text{df}_E)  $
Confidence interval	<p>CI for mean of group <math>i</math>:</p> <p><b>(homoskedasticity met)</b>: <math>\bar{y}_i \pm t_{n-1}^* \frac{s_p}{\sqrt{n_i}}</math></p> <p><b>(homoskedasticity violated)</b>: <math>\bar{y}_i \pm t_{n_i-1}^* \frac{s_i}{\sqrt{n_i}}</math></p>

Table 10

ANOVA II	
Number of groups	3+ or 2+ factors
Assumptions	<p>Normality CLT independence <math>\mu_i, i = 1, 2, \dots, I</math>, unknown <math>\sigma_i, i = 1, 2, \dots, I</math>, unknown <b>Homoskedasticity:</b> Homogeneity of variances <math>\sigma_1 = \sigma_2 = \dots = \sigma_I</math> <math>H_0: \mu_1 = \mu_2 = \dots = \mu_I</math> all of the groups are the same <math>H_a</math>: one of the groups is different from the rest</p>
Name	<p>ANOVA II</p> <p>“is the observed variance in the data attributable to the variation between the groups (across factors), or within the groups”</p> <p>Compare the means of 3 or more groups, where group membership is defined by 2 factors, <math>A</math> and <math>B</math>; if 2 groups, then <math>t</math>-test suffices</p>
Statistic	$  \left. \begin{aligned}  &i = 1, 2, \dots, I \text{ indexes } \mathbf{factor } A \\  &j = 1, 2, \dots, J \text{ indexes } \mathbf{factor } B \\  &k = 1, 2, \dots, n \text{ indexes the individual score} \\  &SST = SSA + SSB + SSAB + SSE \\  &df_T = df_A + df_B + df_{A \times B} + df_E \\  &MST = \frac{SST}{df_T} = \text{Var } y \\  &SST = \sum_{k=1}^n \sum_{i=1}^I \sum_{j=1}^J (y_{ijk} - \bar{y})^2 \\  &df = n - 1 \\  &MSA = \frac{SSA}{df_A} \\  &SSA = \sum_{k=1}^n \sum_{i=1}^I \sum_{j=1}^J (\bar{y}_i - \bar{y})^2 = \sum_{i=1}^I n \times J \times (\bar{y}_i - \bar{y})^2 \\  &df_A = I - 1 \\  &MSB = \frac{SSB}{df_B} \\  &SSB = \sum_{k=1}^n \sum_{i=1}^I \sum_{j=1}^J (\bar{y}_j - \bar{y})^2 = \sum_{j=1}^J n \times I \times (\bar{y}_j - \bar{y})^2 \\  &df_B = J - 1 \\  &MSAB = \frac{SSAB}{df_{A \times B}} \\  &SSAB = \sum_{k=1}^n \sum_{i=1}^I \sum_{j=1}^J (\bar{y}_{ij} - \bar{y}_i - \bar{y}_j + \bar{y})^2 \\  &df_{A \times B} = (I - 1) \times (J - 1) \\  &MSE = \frac{SSE}{df_E} = s_p^2 \\  &SSE = \sum_{k=1}^n \sum_{i=1}^I \sum_{j=1}^J (y_{ijk} - \bar{y})^2 \\  &df_E = n - I \times J  \end{aligned} \right\} \implies \begin{cases} F_A = \frac{MSA}{MSE} \sim F(df_A, df_E) \\ F_B = \frac{MSB}{MSE} \sim F(df_B, df_E) \\ F_{A \times B} = \frac{MSAB}{MSE} \sim F(df_{A \times B}, df_E) \end{cases}  $
Confidence interval	<p><math>CI_{ij}</math> for comparing the means of groups <math>i</math> and <math>j</math>: <math>(\bar{y}_i - \bar{y}_j) \pm t_{df_E}^{**} s_p \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}</math> (See the section on post-hoc methods)</p>

## Exercises

### Regression - ANOVA analysis

1. The “Healthy Breakfast” dataset contains, among other variables, the Consumer Reports ratings of 77 cereals, the number of grams of sugar contained in each serving, and the number of grams of fat contained in each serving.

Considering “Sugars” as the explanatory variable and “Rating” as the response variable generated the following regression line:

$$\text{Rating} = 59.3 - 2.40 \text{ Sugars}$$

Source	DF	SS	MS	$F$	$p$
Regression	1	8654.7	8654.7	102.35	0.000
Error	75	6342.1	84.6		
Total	76	14996.8	194.76		

Table 11: Analysis of Variance - rating ~ sugar

As a simple linear regression model, we previously considered “Sugars” as the explanatory variable and “Rating” as the response variable.

The regression line generated by the inclusion of “Sugars” and “Fat” is the following:

$$\text{Rating} = 61.1 - 2.21 \text{ Sugars} - 3.07 \text{ Fat}$$

Source	DF	SS	MS	$F$	$p$
Regression	2	9325.3	4662.6	60.84	0.000
Error	74	5671.5	76.6		
Total	76	14996.8	194.76		

Source	DF	Seq SS
Sugars	1	8654.7
Fat	1	670.5

Table 12: Analysis of Variance - rating ~ sugar + fat

- (a) Define the population regression model using table 12. If two cereals have the same fat content but different sugar content, what can you say about the rating?
- (b) What does VIF stand for? Compute the VIF using table 12.
- (c) What does VAF stand for? Compute the VAF using tables 11 and 12.
- (d) How do the ANOVA results change when “FAT” is added as a second explanatory variable?
- (e) Formulate appropriate hypotheses, make a decision and explain your reasoning.

*Solution.*

- (a) The population regression model used in table 12 is  $\text{rating} = \beta_0 + \beta_1 \text{ sugars} + \beta_2 \text{ fat} + \varepsilon$ , where  $\beta_j$ 's are approximated by  $b_j$ 's such that  $\vec{b} = (61.1, -2.21, -3.07)^T$ . If variable fat is kept constant, then the marginal difference in rating is -2.21 per gram of sugar. This says that the rating of the breakfast cereal will decrease by 2.21 points per additional gram of sugar, under the condition that fat content is kept constant.
- (b) VIF stands for variance inflation factor, and is given by  $\text{VIF}_j = 1/(1 - R_j^2)$ , where  $R_j^2$  is the coefficient of determination of the regression equation  $X_j = \alpha_0 + \alpha_1 X_{-j} + \delta$  (regress the explanatory variables on the others). The square root of the VIF indicates how much larger the standard error increases compared to if that variable had 0 correlation to other predictor variables in the model. For example, if the variance inflation factor of a predictor variable were 5.27 ( $\sqrt{5.27} = 2.3$ ), this means that the standard error for the coefficient of that predictor variable is 2.3 times larger than if that predictor variable had 0 correlation with the other predictor variables. Rule of thumb:  $\text{VIF}_j > 10$  indicates multicollinearity in the model, i.e. explanatory variables are dependent on each other.

It is not possible to compute the VIF using table 12, as we need the partial SS information.

- (c) VAF stands from variance accounted for and is given by the  $R^2$  coefficient for linear regression, where  $R^2 = \text{SSR}/\text{SST}$ . From table 11,  $R^2 = 8654.7/14996.8 = 0.577$ , and from table 12,  $R^2 = 9325.3/14996.8 = 0.622$ .

- (d) Comparing the VAF's tells us that the model is improved with the addition of fat to the model. This is further shown by the column Seq SS, which shows that fat reduces the SSE by 670.5, which in turn reduces the MSE, indicating less deviation between the observed and fitted values.
- (e)  $H_0 : \beta_2 = 0$  and  $H_A : \beta_2 \neq 0$ .  $F$  is significant with  $p < 0.05$ , i.e. reject  $H_0$  in favour of  $H_A$  and conclude that fat is in the population model.

2. Answer the following questions using the tables and graphs below.

Table 13: Descriptive Statistics

	sales	adverts	airplay	attract
Valid	200	200	200	200
Missing	0	0	0	0
Mean	193.200	614.412	27.500	6.770
Std. Error of Mean	5.706	34.341	0.868	0.099
Std. Deviation	80.699	485.655	12.270	1.395
Minimum	10.000	9.104	0.000	1.000
Maximum	360.000	2271.860	63.000	10.000

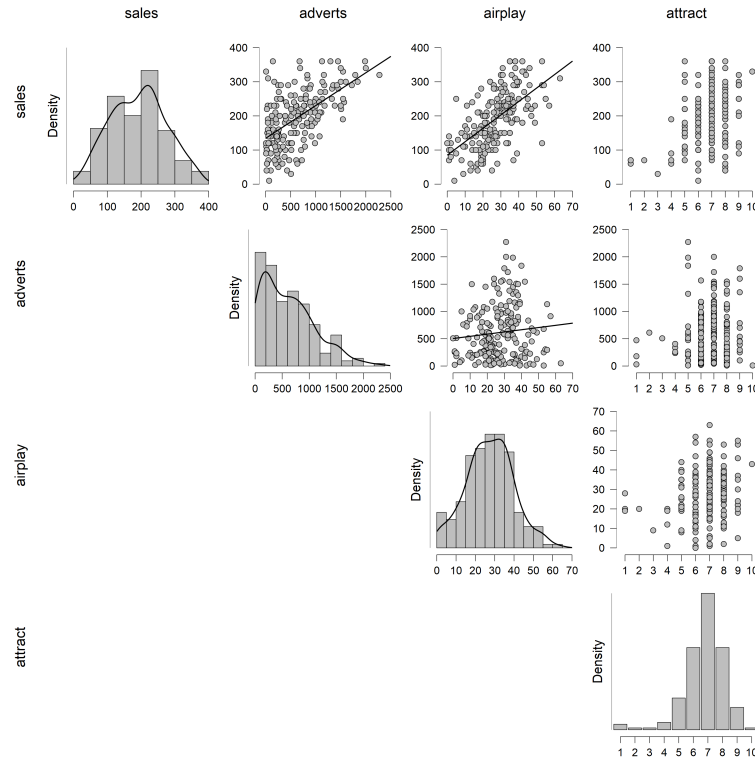


Figure 3

Table 14: Model Summary

Model	R	R <sup>2</sup>	Adjusted R <sup>2</sup>	RMSE	R <sup>2</sup> Change	F Change	df1	df2	p
0	0.578	0.335	0.331	65.991	0.335	99.587	1	198	< .001
1	0.815	0.665	0.660	47.087	0.330	96.447	2	196	< .001

Table 15: Coefficients

Model		Unstandardized	Standard Error	Standardized	t	p	95% CI		Collinearity Statistics	
							Lower	Upper	Tolerance	VIF
0	(Intercept)	134.140	7.537		17.799	< .001	119.278	149.002		
	adverts	0.096	0.010	0.578	9.979	< .001	0.077	0.115	1.000	1.000
1	(Intercept)	-26.613	17.350		-1.534	0.127	-60.830	7.604		
	adverts	0.085	0.007	0.511	12.261	< .001	0.071	0.099	0.986	1.015
	airplay	3.367	0.278	0.512	12.123	< .001	2.820	3.915	0.959	1.043
	attract	11.086	2.438	0.192	4.548	< .001	6.279	15.894	0.963	1.038

Table 16: ANOVA

Model		Sum of Squares	df	Mean Square	F	p
0	Regression	433687.833	1	433687.833	99.587	< .001
	Residual	862264.167	198	4354.870		
	Total	1.296e+6	199			
1	Regression	861377.418	3	287125.806	129.498	< .001
	Residual	434574.582	196	2217.217		
	Total	1.296e+6	199			

- What is the population regression equations for Model 0 and 1?
- Describe the regression equations you wrote above, in words (1-2 sentences each model). What is the point of the comparison?
- Summarise the findings of table 13 and compare with fig. 3.
- Write the null and alternative hypothesis based on the regression equations you wrote in part (a). Now, describe these hypotheses in words (do not refer to beta coefficients, just use plain language - like you're informing a friend). What does table 14 inform you about your hypotheses?
- Under Model 0, what is the expected output if the explanatory variable input has value 600? Compare this with output with the output from Model 1 under the same conditions. Explain the difference in your results.
- Explain the fourth and seventh columns of table 14.
- Table 15 provides you with the VIF for both models. Interpret the results without making too many references to the exact value of the VIF, i.e. what do these values mean?
- Provide the standardised regression equations for both models.
- Use table 16 to make your decision about your hypotheses. Explain your reasoning.

*Solution.*

- Model 0:  $\text{sales} = \beta_0 + \beta_1 \text{adverts} + \varepsilon$ .  
Model 1:  $\text{sales} = \beta_0 + \beta_1 \text{adverts} + \beta_2 \text{airplay} + \beta_3 \text{attract} + \varepsilon$ .
- In Model 0, we propose that the album sales depends only on the number of advertisements for the given album.  
In Model 1, we propose that the album sales depends not only on the number of advertisements, but also on the radio airplay time