

Assigned Nov 24: Handout

Assigned Nov 26: S4.4 and S4.10 to S4.14

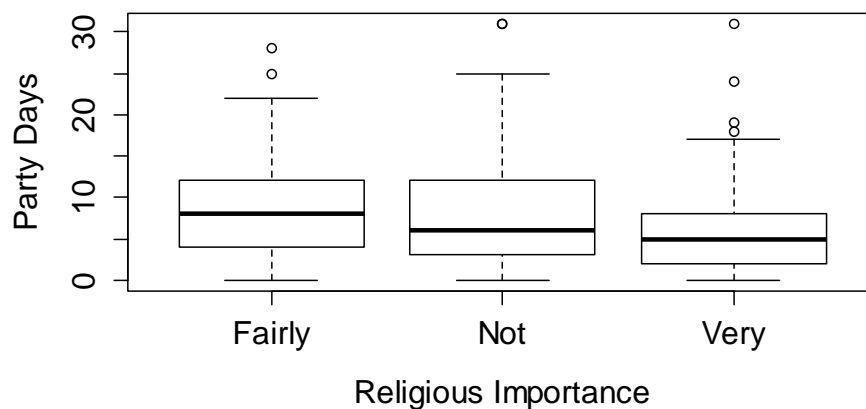
Assigned Dec 2: S4.5 to S4.9

Assigned Nov 24:

1. Create side by side boxplots comparing the PartyDays for the three Religious Importance groups.

Solution: See plot below. The R command used:

```
> boxplot(PartyDays~ReligImp,xlab="Religious Importance",ylab="PartyDays")
```



2. Find and report the sample means, standard deviations and sample sizes for the three Religious Importance groups.

Solution: The R commands (one version!) and results are as shown below.

```
> numSummary(Dataset[, "PartyDays"], groups=Dataset$ReligImp,
statistics=c("mean", "sd", "length"))
```

	mean	sd	data:n
Fairly	8.158228	5.181802	316
Not	7.716216	5.805724	222
Very	5.884354	5.190240	147

3. Using the boxplots in question 1 and the standard deviations in question 2, comment on whether or not it is appropriate to conduct an analysis of variance.

Solution: The purpose of this question is to examine the assumptions needed for ANOVA. The box plots show just a few outliers and some skewness, so you would be justified if you question the normality assumption. But the sample sizes are large so it is okay to proceed. The sample standard deviations are all very similar, so the assumption of equal standard deviations in the population is reasonable.

4. Write the population model and the null and alternative hypotheses for analysis of variance using both versions, as follows.

a. Write the version of the model using group means. Define the parameters. Specify the conditions that accompany the model.

Solution: $Y = \mu_i + \varepsilon$; conditions are on the error terms ε and state that (1) They have mean 0, (2) They have common standard deviation σ , (3) They come from a normal distribution, (4) They are independent.

The parameters in this case are μ_1, μ_2, μ_3 , where in each case μ_i is the population mean Party Days for all students in the population who would give the level i answer (Not, Fairly, Very) to the religious importance question if asked.

b. Write the “factor effects” version of the model.

Solution: $Y = \mu + \alpha_i + \varepsilon$ where $\alpha_i = \mu_i - \mu$ and $\sum \alpha_i = 0$.

c. Give the null and alternative hypotheses using the group means version, and then using the “factor effects” version.

Solution: Group means version: $H_0: \mu_1 = \mu_2 = \mu_3$ and H_a : Not all μ_i are equal.

Factor effects version: $H_0: \alpha_1 = \alpha_2 = \alpha_3 = 0$ and H_a : Not all α_i are 0.

5. Conduct a one-way analysis of variance to compare mean party days for the three religious importance groups, as follows:

a. State null and alternative hypotheses in symbols and words *in the context of this situation*. Define the parameters used in your hypotheses.

Solution: Hypotheses in symbols are given in 4c; you can state them here using either version.

Using the version with μ_i the parameter definition is that they are the population mean Party Days for all students in the population who would give the level i answer (Not, Fairly, Very) to the religious importance question if asked.

The hypotheses in words:

H_0 : The *population* mean Party Days are the same for all three religious importance groups.

H_a : At least one of the population means differs from the others.

b. Show the ANOVA table produced by R.

Solution:

```
> Model<-aov(PartyDays~ReligImp)
> summary(Model)
              Df Sum Sq Mean Sq F value    Pr(>F)
ReligImp      2     531   265.25    9.118 0.000124 ***
Residuals   682   19840    29.09
```

c. Give the test statistic and p-value for the test.

Solution: The test statistic is $F = 9.118$, $p\text{-value} = 0.000124$.

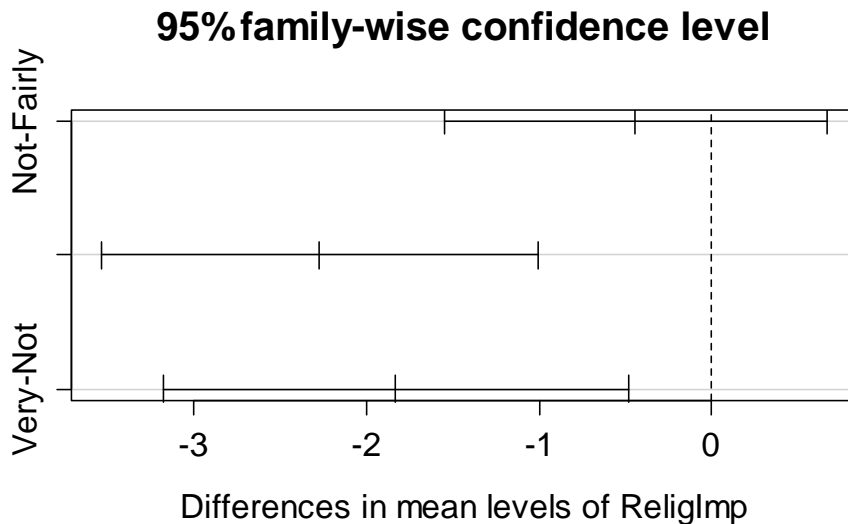
d. State a conclusion in statistical terms *and* in the context of this situation.

Solution: The $p\text{-value}$ of $0.000124 < 0.05$; do reject the null hypothesis. We conclude that the *population* mean number of days students go to parties in a month differ for students in the different religious importance groups.

6. Use the Tukey multiple comparisons procedure to determine which population means differ (if any). Show both numerical and graphical results. Write a summary stating your conclusions.

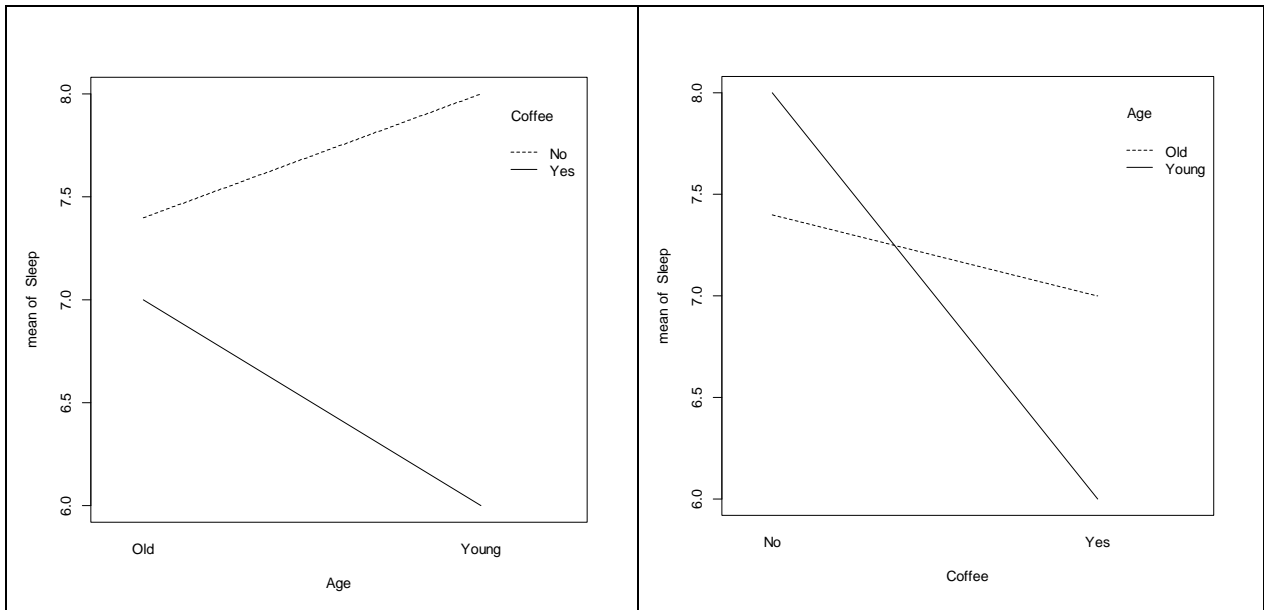
Solution: Results are shown below. Both the numerical results and the plot show the population means for party days are significantly different for the “Very” and “Fairly” groups and the “Very” and “Not” groups, but are not significantly different for the “Fairly” and “Not” groups.

```
> TukeyHSD(Model)
  Tukey multiple comparisons of means
    95% family-wise confidence level
Fit: aov(formula = PartyDays ~ ReligImp)
$ReligImp
      diff      lwr      upr    p adj
Not-Fairly -0.4420116 -1.551456  0.6674325 0.6178056
Very-Fairly -2.2738741 -3.538676 -1.0090720 0.0000812
Very-Not    -1.8318625 -3.179001 -0.4847237 0.0041855
> plot(TukeyHSD(Model))
```



Assigned Nov 26: S4.4 and S4.10 to S4.14

S4.4 The two versions of the interaction plot are shown below. You only need to show one of them. (Yours might look different if you reversed the factor levels Old, Young and No, Yes on the horizontal axis. I did these in R, so the factor levels are shown in alphabetical order.) There may be a minor Factor A effect (Age) but it looks like the average across the two coffee groups is quite similar for the Old and Young. There is a larger Factor B effect (Coffee) as illustrated by the gap between the two lines in the plot on the left. The mean for the “No” coffee group is considerably higher than for the “Yes” coffee group. There is also an interaction. The difference between mean sleep hours for the older group is much smaller than it is for the younger group.



- S4.10**
- a. Yes. Mean DDT appears to be much higher in the Arctic region than in the U.S. or Canada.
 - b. Yes, Mean DDT increases with age.
 - c. No. The change over age is similar for all three regions (the lines are almost parallel), and the differences among regions are similar at all ages. (You might think there is an interaction effect because two of the lines cross, but that is very minor – they are almost on top of each other.)

- S4.11**
- a. H_0 : The population mean DDT levels in falcons are the same in the Artic, the U.S. and Canada.
 H_a : For at least one of the 3 sites the population mean DDT level is different.
 - b. H_0 : The population mean DDT levels in falcons are the same for the three age groups.
 H_a : For at least one of the 3 age categories the population mean DDT level is different.
 - c. Here is one way to write them.
 H_0 : The differences in population mean DDT levels in falcons at various sites do not depend on the ages of the falcons.
 H_a : The differences in population mean DDT levels in falcons at various sites depend on the ages of the falcons.

- S4.12** a. p -value = 0.000, reject the null hypothesis, conclude that mean DDT levels differ across sites.
 b. p -value = 0.000, reject the null hypothesis, conclude that mean DDT levels differ across ages.
 c. p -value = 0.313, do not reject the null hypothesis. We cannot conclude that the differences in population mean DDT levels at various sites depend on the ages of the falcons.
- S4.13** a. The mean DDT level in falcons is much higher in the Arctic than in the U.S. or Canada.
 b. The mean DDT level increases with age.
 c. No. There is little interaction between the factors, so the interpretation of one factor does not depend on the levels of the other factor.
- S4.14** For Factor A, $F = 8892.70/3.44 = 2585.1$; for Factor B, $F = 860.59/3.44 = 250.2$, for the interaction, $F = 4.43/3.44 = 1.29$. Note that these differ slightly from the F values given by the computer (Minitab) because the mean squares were rounded to 2 decimal places when they were printed, but were used with many more decimal places in the computer calculations of the F values.

Assigned Dec 2: S4.5 to S4.9

- S4.5** a. μ_{12} = mean for 18-21 non-coffee drinkers = 8.0 hours.
 b. $\mu_{..}$ = mean for all groups = $28.4/4 = 7.1$ hours.
 c. $\mu_{1.}$ = mean for 18-21 year olds = 7.0 hours, $\mu_{2.}$ = mean for over 21 year olds = 7.2 hours
 d. $\mu_{.1}$ = mean for coffee drinkers = 6.5 hours; $\mu_{.2}$ = mean for non-coffee drinkers = 7.7 hours
- S4.6** a. $\alpha_1 = \mu_{1.} - \mu_{..} = 7.0 - 7.1 = -0.1$ hours; $\alpha_2 = \mu_{2.} - \mu_{..} = 7.2 - 7.1 = 0.1$ hours
 b. $\beta_1 = \mu_{.1} - \mu_{..} = 6.5 - 7.1 = -0.6$ hours; $\beta_2 = \mu_{.2} - \mu_{..} = 7.7 - 7.1 = 0.6$ hours
 c. $\alpha\beta_{11} = \mu_{11} - (\mu_{..} + \alpha_1 + \beta_1) = 6.0 - (7.1 - 0.1 - 0.6) = 6.0 - 6.4 = -0.4$ hours. Interactions must sum to 0 over each row and column, so $\alpha\beta_{12} = \alpha\beta_{21} = 0.4$ hours and $\alpha\beta_{22} = -0.4$ hours.
- S4.7** $\mu_{11} = \mu_{..} + \alpha_1 + \beta_1 = 7.1 - 0.1 - 0.6 = 6.4$ hours, $\mu_{12} = \mu_{..} + \alpha_1 + \beta_2 = 7.6$ hours, $\mu_{21} = \mu_{..} + \alpha_2 + \beta_1 = 6.6$ hours, $\mu_{22} = \mu_{..} + \alpha_2 + \beta_2 = 7.8$ hours. These are each off by 0.4 hours one way or the other, so the additive model is not adequate.
- S4.8** $\mu_{11} = \mu_{..} + \alpha_1 + \beta_1 + \alpha\beta_{11} = 7.1 - 0.1 - 0.6 - 0.4 = 6.0$ hours, $\mu_{12} = \mu_{..} + \alpha_1 + \beta_2 + \alpha\beta_{12} = 7.1 - 0.1 + 0.6 + 0.4 = 8.0$ hours, $\mu_{21} = \mu_{..} + \alpha_2 + \beta_1 + \alpha\beta_{21} = 7.1 + 0.1 - 0.6 + 0.4 = 7.0$ hours, $\mu_{22} = \mu_{..} + \alpha_2 + \beta_2 + \alpha\beta_{22} = 7.1 + 0.1 + 0.6 - 0.4 = 7.4$ hours.
- S4.9** Moving from the first row to the second row averaged over the two columns shows that the average difference in sleep hours for the two age groups is small, 7.0 for the younger group compared with 7.2 for the older group. Moving from the first column to the second column averaged over the two rows shows that the average difference in sleep hours for coffee drinkers is lower than for non-coffee drinkers, 6.5 hours compared with 7.7 hours. The interaction effect is clear because for coffee drinkers the older age group gets more sleep, while for non-coffee drinkers the younger age group gets more sleep.