

STAT W4201 001, Homework 3

Brian Weinstein (bmw2148)

Feb 17, 2016

Code is attached here and also posted at <https://github.com/BrianWeinstein/advanced-data-analysis>. Where relevant, code snippets and output are included in-line.

Problem 1: Ramsey 4.30

A boxplot of SPF estimates for each subject is shown in Figure 1, where the SPF estimate is calculated for each subject as the number minutes that a person was able to withstand sunlight after applying sunscreen divided by the number of minutes they were able to withstand before applying sunscreen. The dataset contains paired observations, and the quantity in question has a nearly normal distribution with no outliers, making it a strong candidate for the one-sample paired t-test.

	PreTreatment	Sunscreen	spfEstimate
1	30	120	4.000000
2	45	240	5.333333
3	180	480	2.666667
4	15	150	10.000000
5	200	480	2.400000
6	20	270	13.500000
7	15	300	20.000000
8	10	180	18.000000
9	20	300	15.000000
10	20	240	12.000000
11	60	480	8.000000
12	60	300	5.000000
13	120	480	4.000000

The SPF Estimate has a mean of 9.223, and the standard error on the mean is $SD(sp\hat{f}Estimate)/\sqrt{13} = 1.664$. Using a t distribution with $(13 - 1) = 12$ degrees of freedom, a 95% confidence interval for the mean is

$$9.223 \pm t_{12}(1 - (0.05/2)) \times 1.664$$

$$9.223 \pm (2.179)(1.664)$$

$$9.223 \pm 3.626$$

$$\Rightarrow 95\% \text{ CI for the mean SPF Estimate: } [5.598, 12.848]$$

A potentially confounding variable could be the strength of sunlight during the measurement period for each subjects' data points. For a given participant, if the sunlight strength wasn't constant for both measurements, it likely would have had an impact on the duration for which they were able to withstand the sunlight.

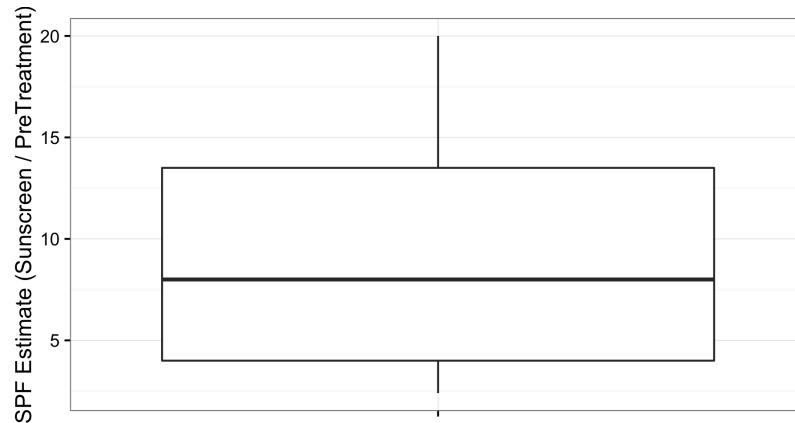


Figure 1: SPF estimates for 13 subjects.

Problem 2: Ramsey 4.32

See attached code. A boxplot of the difference in the number of retching episodes while on the marijuana treatment vs the placebo is show in Figure 2.

The data provides convincing evidence that marijuana reduced the frequency of retching episodes compared to the placebo (estimated one-sided p-value = 0.000156 from the sign test). The marijuana treatment reduced the number of retching episodes for a given patient by an estimated 25 episodes (95% confidence interval for an additive treatment effect: 7 to 43 episodes).

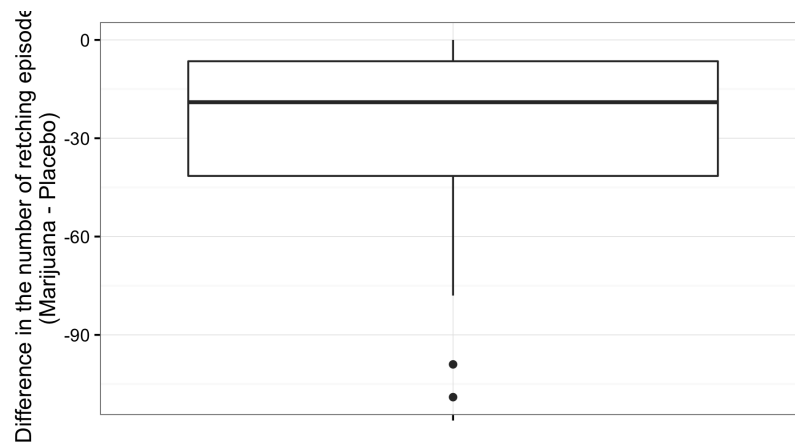


Figure 2: The difference in the number of retching episodes for 15 patients while on a marijuana treatment vs a placebo.

Problem 3: Ramsey 5.19

- (a) The pooled estimate of the variance s_p^2 is given by

$$\begin{aligned}
s_p^2 &= \frac{\sum_{i=1}^I (n_i - 1) s_i^2}{\sum_{i=1}^I (n_i - 1)} \\
&= \frac{(127 - 1)(0.4979)^2 + (44 - 1)(0.4235)^2 + (24 - 1)(0.3955)^2 + \cdots}{(127 - 1) + (44 - 1) + (24 - 1) + \cdots} \\
&\quad \frac{\cdots + (41 - 1)(0.3183)^2 + (18 - 1)(0.3111)^2 + (16 - 1)(0.4649)^2 + \cdots}{\cdots + (41 - 1) + (18 - 1) + (16 - 1) + \cdots} \\
&\quad \frac{\cdots + (11 - 1)(0.2963)^2 + (7 - 1)(0.3242)^2 + (6 - 1)(0.5842)^2}{\cdots + (11 - 1) + (7 - 1) + (6 - 1)} \\
&= 0.1919322.
\end{aligned}$$

- (b) To construct an ANOVA table to test for species differences, we first compute SS_W , the sum of squared residuals of the “within” group.

$$\begin{aligned}
SS_W &= \sum_{i=1}^9 \left[\sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2 \right] = \sum_{i=1}^9 [(n_i - 1) \cdot \text{SE}(Y_i)^2] \\
&= (127 - 1)(0.4979)^2 + (44 - 1)(0.4235)^2 + \cdots + (6 - 1)(0.5842)^2 \\
&= 54.70068.
\end{aligned}$$

Given that the sample standard deviation of all 294 observations as one group is $\text{SD}_{Total} = 0.4962$, the total sum of squared residuals SS_{Total} is

$$\begin{aligned}
SS_{Total} &= \sum_{i=1}^9 \left[\sum_{j=1}^{n_i} (Y_{ij} - \bar{Y})^2 \right] = (n_1 + n_2 + \cdots + n_9 - 1)(\text{SD}_{Total})^2 \\
&= (127 + 44 + \cdots + 6 - 1)(0.4962)^2 = (293)(0.4962)^2 \\
&= 72.14083.
\end{aligned}$$

Therefore, the sum of squared residuals of the “between” group, SS_B is:

$$SS_B = SS_{Total} - SS_W = 72.14083 - 54.70068 = 17.44015.$$

The Total, Within, and Between degrees of freedom are $\text{df}_{Total} = (n - 1) = 293$, $\text{df}_W = (n - I) = (294 - 9) = 285$, and $\text{df}_B = \text{df}_{Total} - \text{df}_W = 8$, respectively.

The Mean Square for the Within and Between groups is the $[(\text{Sum of Squares}) / \text{df}]$, and the F-statistic is defined as the $[(\text{Between Mean Square}) / (\text{Within Mean Square})]$.

Thus the ANOVA table to test for species differences is:

Source of Variation	Sum of Squares	d.f.	Mean Square	F-Statistic	p-Value
Between Groups	17.44015	8	2.180019	11.35828	5.6683×10^{-14}
Within Groups	54.70068	285	0.1919322		
Total	72.14083	293			

Where the p-value for the F-Statistic is given by the CDF of an F distribution with 8 and 285 degrees of freedom.

```
> # p-value of F-statistic
> fstat
[1] 11.35828
> pf(q=fstat, df1=8, df2=285, lower.tail=FALSE)
[1] 5.6683e-14
```

- (c) The ANOVA method for calculating SS_B yields the same answer as the formula:

$$SS_B = \sum_{i=1}^I \left(n_i \bar{Y}_i^2 \right) - n \bar{Y}^2 = \sum_{i=1}^9 \left(n_i \bar{Y}_i^2 \right) - n \bar{Y}^2.$$

The overall mean $\bar{Y} = 7.4746$ is computed as a weighted average of the group means, and $n = 294$, thus

$$\begin{aligned} SS_B &= \sum_{i=1}^9 \left(n_i \bar{Y}_i^2 \right) - n \bar{Y}^2 \\ &= (127)(7.347)^2 + (44)(7.368)^2 + (24)(7.418)^2 + (41)(7.487)^2 + (18)(7.563)^2 + \cdots \\ &\quad \cdots + (16)(7.568)^2 + (11)(8.214)^2 + (7)(8.272)^2 + (6)(8.297)^2 - (294)(7.4746)^2 \\ &= 17.45402, \end{aligned}$$

which is the same, to within rounding errors, as the SS_B calculated using the ANOVA method (17.44015).

- (d) If the first 6 species (“Group A”) have one common mean and the last 3 species (“Group B”) have another common mean, our full and reduced models are:

Species Number:	1	2	3	4	5	6	7	8	9
Full model (separate-means):	\bar{Y}_1	\bar{Y}_2	\bar{Y}_3	\bar{Y}_4	\bar{Y}_5	\bar{Y}_6	\bar{Y}_7	\bar{Y}_8	\bar{Y}_9
Reduced model (two-group):	\bar{Y}_A	\bar{Y}_A	\bar{Y}_A	\bar{Y}_A	\bar{Y}_A	\bar{Y}_A	\bar{Y}_B	\bar{Y}_B	\bar{Y}_B

where $\bar{Y}_1, \bar{Y}_2, \dots, \bar{Y}_9$ are the group means given in Display 5.23, and \bar{Y}_A and \bar{Y}_B are weighted averages of $\bar{Y}_1, \dots, \bar{Y}_6$ and $\bar{Y}_7, \dots, \bar{Y}_9$, respectively, as shown below:

$$\begin{aligned} \bar{Y}_A &= \frac{\sum_{i=1}^6 n_i \bar{Y}_i}{\sum_{i=1}^6 n_i} = \frac{(127)(7.347) + \cdots + (16)(7.568)}{127 + \cdots + 16} = 7.405489 \\ \bar{Y}_B &= \frac{\sum_{i=7}^9 n_i \bar{Y}_i}{\sum_{i=7}^9 n_i} = \frac{(11)(8.214) + \cdots + (6)(8.297)}{11 + \cdots + 6} = 8.251667. \end{aligned}$$

The group Sample Standard Deviations are computed using the pooled standard deviations (from the given Sample Standard Deviation values in Display 5.23):

$$SD(Y_A) = \sqrt{\frac{\sum_{i=1}^6 (n_i - 1) (SD(Y_i))^2}{\sum_{i=1}^6 (n_i - 1)}} = 0.4416123$$

$$SD(Y_B) = \sqrt{\frac{\sum_{i=7}^9 (n_i - 1) (SD(Y_i))^2}{\sum_{i=7}^9 (n_i - 1)}} = 0.3912750.$$

Therefore the RSS for Groups A and B are

$$RSS_A = (n_A - 1) (SD(Y_A))^2 = (270 - 1)(0.4416123)^2 = 52.460766$$

$$RSS_B = (n_B - 1) (SD(Y_B))^2 = (24 - 1)(0.3912750)^2 = 3.521211.$$

These values are summarized in the table below:

	group	n	mean	pooled_sd	rss
1	A	270	7.405489	0.4416123	52.460766
2	B	24	8.251667	0.3912750	3.521211

The Sum of Squares for the two-group model is just the sum of the RSS for Group A and B (55.98198), with 293 degree of freedom.

Therefore we can construct a new ANOVA table and perform an F-test:

is this correct?

Source of Variation	Sum of Squares	d.f.	Mean Square	F-Statistic	p-Value
Between Groups A and B	16.15885	1	16.15885	84.28401	1.36×10^{-70}
Within Groups A and B	55.98198	292	0.1917191		
Total	72.14083	293			

Problem 4:

Consider the Bumpus's data in Chapter 2, compute the power of the two-sided two sample t -test of size 0.05 (i.e., reject the null hypothesis if the absolute value the t -statistic is greater than or equal to 2), under the alternative that $\mu_x - \mu_y = \bar{x} - \bar{y} = 0.01$ and $\sigma = s_p = 0.0214$.

Problem 5:

Show that the two-sided two sample t -test is equivalent to the anova F -test, if the number of groups is two.

For $I = 2$ groups, the F -statistic is given by

$$F\text{-statistic} = \frac{SS_B / [(n - 1) - (n - I)]}{SS_W / (n - I)},$$

where n_1 and n_2 are the sizes of samples 1 and 2, respectively, $n = n_1 + n_2$ is the total sample size, SS_B is the "between groups" sum of squared residuals, and SS_W is the "within groups" sum of squared residuals.

Simplifying, we find

$$F\text{-statistic} = \frac{SS_B / (I - 1)}{SS_W / (n - I)} = \frac{SS_B / (2 - 1)}{SS_W / (n - 2)} = \frac{SS_B / 1}{SS_W / (n - 2)}.$$

If the observations from group 1 are $\sim N(\mu_1, \sigma^2)$ and the observations from group 2 are $\sim N(\mu_2, \sigma^2)$, we know that

$$F\text{-statistic} \sim F_{1, n-2}, \text{ which is equivalent to } t_{n-2}^2.$$

i.e., an F distribution with a numerator degrees of freedom of 1 and a denominator degrees of freedom of $n - 2$ is equivalent to the square of a t distribution with $n - 2$ degrees of freedom.

Problem 6:

Consider X_1, \dots, X_{10} are i.i.d. $N(0, \sigma^2)$, Y_1, \dots, Y_{10} are i.i.d. $N(\mu, \sigma^2)$ and hypothesis testing:

$$H_0 : \mu = 0$$

$$H_A : \mu \neq 0.$$

Compute the power of a two sided two sample t-test of size 0.05 when $\sigma^2 = 1$ and $\mu = 0.1, 0.5, 1$, and 2. Plot the power as a function of μ . Then, increase the sample size in each group to 20 and draw the power function in the same plot as that of the sample size 10.

See attached code. A plot of the power for a two-sided, two-sample t-test of size 0.05, for $\sigma^2 = 1$, $\mu = \{0.1, 0.5, 1, 2\}$, and $(n_1, n_2) = \{(10, 10), (20, 20)\}$ is shown in Figure 3.

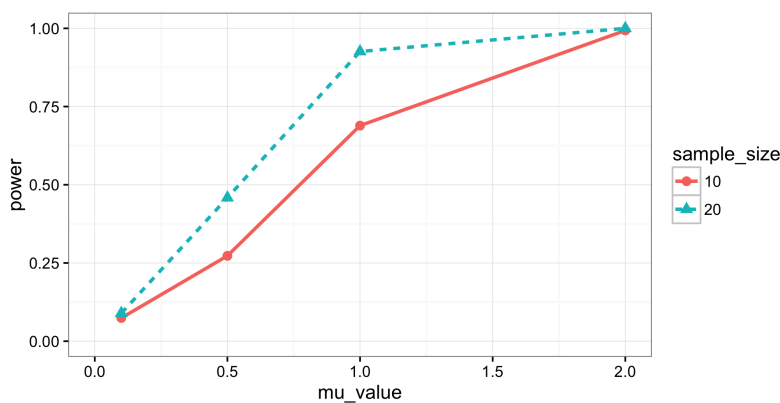


Figure 3: The power of a two-sided, two-sample t-test of size 0.05 at various sample sizes and μ values.

Problem 7:

Under the setting of the previous problem, show that, under the null hypothesis, the p-value follows the uniform distribution on the interval $[0, 1]$ and perform simulations to confirm it.

Under the null hypothesis that $\mu = 0$, our t-statistic $t = \frac{\bar{X} - \bar{Y}}{\text{SE}(\bar{X} - \bar{Y})}$ follows a t distribution with $n_x + n_y - 2$ degrees of freedom (since X and Y are normally distributed).

The p-value is defined as

$$p = P(T \geq t) = 1 - P(T < t) = 1 - F(t),$$

where T is the general t distribution from which our test-statistic t is a realization, and $F(t)$ is the cumulative distribution function of the t distribution at $T = t$.

Since $F(t)$ is a monotonically increasing function between $[0, 1]$, we can write

$$p = P(T \geq t) = P(F(T) \geq F(t)) = 1 - P(F(T) < F(t)).$$

Equating these, we see that

$$\begin{aligned} 1 - F(t) &= 1 - P(F(T) < F(t)) \\ &\Rightarrow F(t) = P(F(T) < F(t)). \end{aligned}$$

Therefore $F(T)$, and equivalently $1 - F(T)$, must follow a uniform distribution on $[0, 1]$. Since $p = 1 - F(T)$, we see that the p-value is also uniformly distributed on $[0, 1]$.

The code below generates 20,000 two-sided p-values, simulated under the conditions specified in Problem 6, and the histogram in Figure 4 shows the (nearly uniform) distribution of the 20,000 two-sided p-values.

```
set.seed(1)
sigma <- 1

# initialize an empty list to store results
pvalList <- list()

for(i in 1:20000){

  # generate random samples
  nx <- 10
  x.vals <- rnorm(n=nx, mean=0, sd=sigma)
  ny <- 10
  y.vals <- rnorm(n=10, mean=0, sd=sigma)

  # calculate difference in means
  diffMean <- mean(x.vals) - mean(y.vals)

  # calculate the pooled sd
  sp <- sqrt( ((nx-1)*(sd(x.vals)^2) + (ny-1)*(sd(y.vals)^2)) / (nx + ny - 2) )

  # calculate the t statistic and 2-sided p value
  tStat <- diffMean / (sp * sqrt((1/nx) + (1/ny)))
  pval <- 2 * pt(q=-1 * abs(tStat), df=(nx + ny - 2)) # 2-sided pvalue

  pvalList[[i]] <- data.frame(pval)
}

pvalList <- rbindlist(pvalList)
```

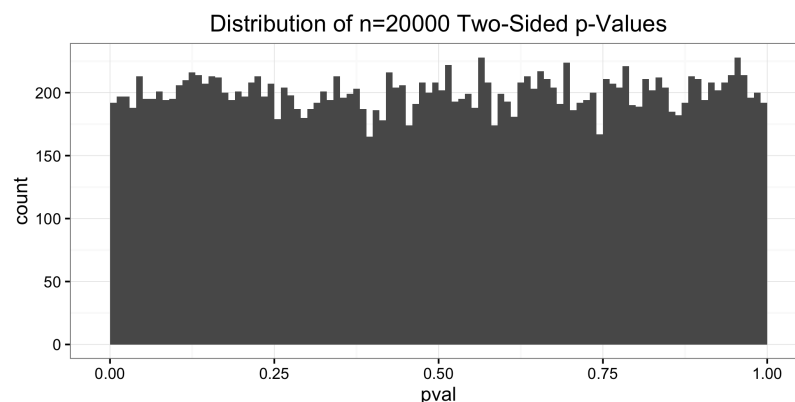


Figure 4: A histogram of 20,000 two-sided p-values, simulated under the conditions specified in Problem 6.