

# STAT S4201 001, Homework 2

Brian Weinstein (bmw2148)

Feb 10, 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:

Suppose that  $\{X_1, X_2, \dots, X_n\}$  is a random sample from  $N(\mu, \sigma^2)$ . Construct a 95% confidence interval for  $\sigma^2$  under the following scenarios: (a)  $\mu$  is known to be 0. (b)  $\mu$  is unknown.

(a) asdf

(b) asdf

Problem  
1a

Type up  
Problem  
1b

Fix  $n = 10$  and  $\sigma = 1$ . Run a Monte Carlo simulation to confirm that the confidence interval you constructed under the scenario (a) produces a coverage of 95%. Report how many random samples were drawn in your simulation and how close your coverage was to 95%.

See attached code. In my simulation I drew 1000 random samples, with a 95% confidence interval capturing the true variance in 948 of the trials (94.8% coverage).

## Problem 2: Ramsey 3.22

```
# Data input
time26 <- c(5.79, 1579.52, 2323.70)
time28 <- c(68.8, 108.29, 110.29, 426.07, 1067.60)
```

(a) Form two new variables by taking the logarithms of the breakdown times.

```
> Y1 <- log(time26) ; Y1
[1] 1.756132 7.364876 7.750916
> Y2 <- log(time28) ; Y2
[1] 4.231204 4.684813 4.703113 6.054604 6.973168
```

(b)  $\bar{Y}_1 - \bar{Y}_2 = 5.6240 - 5.3294 = 0.2946$

(c)  $\exp(\bar{Y}_1 - \bar{Y}_2) = \exp(0.2946) = 1.3426$ , where 1.3426 is the multiplicative treatment effect, indicating that the breakdown time at 26 kV is estimated to be 1.3426 times larger than the breakdown time at 28 kV.

(d) Compute a 95% confidence interval for the difference in mean log breakdown times. Take the antilogarithms of the endpoints and express the result in a sentence.

The pooled standard deviation is given by

$$s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{(n_1 + n_2 - 2)}} = \sqrt{\frac{(3 - 1)(11.257) + (5 - 1)(1.310)}{(3 + 5 - 2)}} = 2.1508,$$

where  $n_1 = 3$  and  $n_2 = 5$  are the sample sizes and  $s_1^2 = 11.2574$  and  $s_2^2 = 1.3104$  are the sample variances of the log-transformed measurements.

The standard error of  $(\bar{Y}_1 - \bar{Y}_2)$  is given by

$$SE(\bar{Y}_1 - \bar{Y}_2) = s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} = (2.1508) \sqrt{\frac{1}{3} + \frac{1}{5}} = 1.5707.$$

Therefore a 95% confidence interval for the difference in mean log breakdown times is

$$\begin{aligned} & (\bar{Y}_1 - \bar{Y}_2) \pm t_{(3+5-2)} \left(1 - \frac{\alpha}{2}\right) SE(\bar{Y}_1 - \bar{Y}_2) \\ & 0.2946 \pm t_6 \left(1 - \frac{0.05}{2}\right) 1.5707 \\ & 0.2946 \pm (2.4469)(1.5707) \\ & 0.2946 \pm 3.8435 \\ & \Rightarrow -3.5489 \leq (\bar{Y}_1 - \bar{Y}_2) \leq 4.1381. \end{aligned}$$

Taking the anilogarithms of the confidence interval endpoints,

$$\text{Lower confidence limit} = e^{-3.5489} = 0.0288$$

$$\text{Upper confidence limit} = e^{4.1381} = 62.682,$$

we find that the breakdown time at 26 kV is estimated to be 1.3426 (from part c) times longer than the breakdown time at 28 kV, (95% confidence: 0.0288 to 62.682 times).

### Problem 3: Ramsey 3.25

See attached code. When excluding either (1) no observations, (2) observation 646, or (3) observations 646 and 645, there is no evidence that the mean dioxin level in Vietnam veterans is greater than the mean dioxin level in non-Vietnam veterans. The difference in means non-(veteran minus veteran) (parts per trillion), one-sided p-values, and 95% confidence intervals for the difference in means (parts per trillion) are:

| Case | Difference in means | one-sided p-value | 95% CI            |
|------|---------------------|-------------------|-------------------|
| (1)  | -0.0745             | 0.3963            | -0.6305 to 0.4815 |
| (2)  | -0.0113             | 0.4805            | -0.6305 to 0.4815 |
| (3)  | 0.0210              | 0.5386            | -0.6305 to 0.4815 |

### Problem 4: Ramsey 3.28

When including all of the observations, the t-test generates a two-sided p-value of 0.0809:

```
> t.test(formula=Humerus~Status, data=sparrowData,
+ var.equal=TRUE, conf.level=0.95)

Two Sample t-test

data:  Humerus by Status
t = -1.777, df = 57, p-value = 0.0809
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -0.021446053  0.001279386
sample estimates:
mean in group Perished mean in group Survived
      0.7279167           0.7380000
```

And when excluding the smallest length in the Perished group, we have a two-sided p-value of 0.18:

```
> t.test(formula=Humerus~Status, data=sparrowData,
+         subset=!(Humerus==min(sparrowData$Humerus) & Status=="Perished")),
+         var.equal=TRUE, conf.level=0.95)
```

Two Sample t-test

```
data: Humerus by Status
t = -1.3578, df = 56, p-value = 0.18
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -0.017542698  0.003368785
sample estimates:
mean in group Perished mean in group Survived
      0.730913          0.738000
```

With the full dataset, the p-value indicates that the data is suggestive, but isn't conclusive, in the hypothesis that the humerus length differs in the two populations. When excluding the smallest length in the Perished group, the p-value indicates that there's little evidence of a difference between the groups.

Since the conclusion of the study depends on which dataset is used, we must either use resistant analysis (Chapter 4) or report the results of both analysis, as we've done here.

**Problem 5:** [Ramsey 3.32](#)

**Problem 6:** [Ramsey 4.19](#)

## Todo list

|                              |   |
|------------------------------|---|
| Problem 1a . . . . .         | 1 |
| Type up Problem 1b . . . . . | 1 |