

# BCB744 (BioStatistics): Summative Task 2, 12 April 2023

AJ Smit

2024-03-25

## On this page

Honesty Pledge . . . . .	1
Instructions . . . . .	2
Question 1 . . . . .	3
Chromosomal effects of mercury-contaminated fish consumption . . . . .	3
Question 2 . . . . .	3
Malignant glioma pilot study . . . . .	3
Question 3 . . . . .	4
Risk factors associated with low infant birth weight . . . . .	4
Question 4 . . . . .	4
The LungCapData.csv data . . . . .	4
Question 5 . . . . .	5
The air quality data . . . . .	5
Question 6 . . . . .	5
The <b>shells.csv</b> data . . . . .	5
Question 7 . . . . .	5
The fertiliser_crop_data.csv data . . . . .	5
The end . . . . .	6

## Honesty Pledge

This assignment requires that you work as an individual and not share your code, results, or discussion with your peers. Penalties and disciplinary action will apply if you are found cheating.

### **i** Acknowledgement of the Pledge

Copy the statement, below, into your document and replace the underscores with your name acknowledging adherence to the UWC's Honesty Pledge.

**I, \_\_\_\_\_, hereby state that I have not communicated with or gained information in any way from my peers and that all work is my own.**

## **Instructions**

Please note the following instructions. Failing to comply with them in full will result in a loss of marks.

- **QUARTO → HTML** Submit your assessment answers as an .html file compiled from your Quarto document. Produce *fully annotated reports*, including the meta-information at the top (name, date, purpose, etc.). Provide ample commentary explaining the purpose of the various tests/sections as necessary.
- **TESTING OF ASSUMPTIONS** For all questions, make sure that when *formal inferential statistics are required*, each is preceded by the appropriate tests for the assumptions, i.e., state the assumptions, state the statistical procedure for testing the assumptions and mention their corresponding  $H_0$ . If a graphical approach is used to test assumptions, explain the principle behind the approach. Explain the findings emerging from the test of assumptions, and justify your selection of the appropriate inferential test (e.g. *t*-test, ANOVA, etc.) that you will use.
- **STATE HYPOTHESES** When inferential statistics are required, please provide the full  $H_0$  and  $H_A$ , and conclude the analysis with a statement of which is accepted or rejected.
- **GRAPHICAL SUPPORT** All descriptive and inferential statistics must be supported by the appropriate figures of the results.
- **STATEMENT OF RESULTS** Make sure that the textual statement of the final result is written exactly as required for it to be published in a journal article. Please consult a journal if you don't know how.
- **FORMATTING** Pay attention to formatting. Some marks will be allocated to the appearance of the script, including considerations of aspects of the tidiness of the file, the use of the appropriate headings, and adherence to code conventions (e.g. spacing etc.).
- **MARK ALLOCATION** Please see the [Introduction Page](#) for an explanation of the assessment approach that will be applied to these questions.

Submit the .html file wherein you provide answers to Questions 1–7 by no later than 19:00 today. Label the script as follows:

BCB744\_<Name>\_<Surname>\_Summative\_Task\_2.html, e.g.

BCB744\_AJ\_Smit\_Summative\_Task\_2.html.

Upload your .html files onto [Google Forms](#).

## Question 1

### Chromosomal effects of mercury-contaminated fish consumption

These data reside in package **coin**, dataset **mercuryfish**. The dataframe contains the mercury level in blood, the proportion of cells with abnormalities, and the proportion of cells with chromosome aberrations in consumers of mercury-contaminated fish and a control group. Please see the dataset's help file for more information.

Analyse the dataset and answer the following questions:

- Does the presence of methyl-mercury in a diet containing fish result in a higher proportion of cellular abnormalities?
- Does the concentration of mercury in the blood influence the proportion of cells with abnormalities, and does this differ between the **control** and **exposed** groups?
- Is there a relationship between the variables **abnormal** and **ccells**? This will have to be for the **control** and **exposed** groups, noting that an interaction effect *might* be present.

## Question 2

### Malignant glioma pilot study

Package **coin**, dataset **glioma**: A non-randomized pilot study on malignant glioma patients with pretargeted adjuvant radioimmunotherapy using yttrium-90-biotin.

- Do **sex** and **group** interact to affect survival time (**time**)?
- Do **age** and **histology** interact to affect survival time (**time**)?
- Show a full graphical exploration of the data. Are there any other remaining patterns visible in the data that should be explored statistically? Study your results, select the most promising and insightful question that remains, and do the analysis.

### Question 3

#### Risk factors associated with low infant birth weight

Package **MASS**, dataset **birthwt**: A dataset about the risk factors associated with low infant birth mass collected at Baystate Medical Center, Springfield, Mass. during 1986.

State three hypotheses and test them. Make sure one of the tests makes use of the 95% confidence interval approach rather than a formal inferential methodology.

### Question 4

#### The **LungCapData.csv** data

- a. Using the Lung Capacity data provided, please calculate the 95% CIs for the **LungCap** variable as a function of:
  - Gender
  - Smoke
  - Caesarean
- b. Create a graph of the mean  $\pm$  95% CIs and determine if there are statistical differences in **LungCap** between the levels of **Gender**, **Smoke**, and **Caesarean**. Do the same using inferential statistics. Are your findings the same using these two approaches?
- c. Produce all the associated tests for assumptions—i.e. the assumptions to be met when deciding whether to use your choice of inferential test or its non-parametric counterpart.
- d. Create a combined tidy dataframe (observe tidy principles) with the estimates for the 95% CI for the **LungCap** data (**LungCap** as a function of **Gender**), estimated using both the traditional and bootstrapping approaches. Create a plot comprising two panels (one for the traditional estimates, one for the bootstrapped estimates) of the mean, median, scatter of raw data points, and the upper and lower 95% CI.
- e. Undertake a statistical analysis that incorporates both the effect of **Age** *and* one of the categorical variables on **LungCap**. What new insight does this provide?

## Question 5

### The air quality data

Package **datasets**, dataset **airquality**. These are daily air quality measurements in New York, May to September 1973. See the help file for details.

- a. Which two of the four response variables are best correlated with each other?

## Question 6

### The **shells.csv** data

This dataset contains measurements of shell widths and lengths of the left and right valves of two species of mussels, *Aulacomya* sp. and *Choromytilus* sp. Length and width measurements are presented in mm.

Fully analyse this dataset.

## Question 7

### The **fertiliser\_crop\_data.csv** data

The data represent an experiment designed to test whether or not fertiliser type and the density of planting have an effect on the yield of wheat. The dataset contains the following variables:

- Final yield (kg per acre)—make sure to convert this to the most suitable SI unit before continuing with your analysis
- Type of fertiliser (fertiliser type A, B, or C)
- Planting density (1 = low density, 2 = high density)
- Block in the field (north, east, south, west)

Fully analyse this dataset.

## **The end**

Submit the .html file wherein you provide answers to Questions 1–7 by no later than 19:00 today. Label the script as follows:

BCB744\_<Name>\_<Surname>\_Summative\_Task\_2.html, e.g.

BCB744\_AJ\_Smit\_Summative\_Task\_2.html.

Upload your .html files onto [Google Forms](#).