

DPH112 Biostatistics 2

Assignment 1

Elmer V Villanueva
Elmer.Villanueva@xjtlu.edu.cn

1 Background

You will make use of the Excel file called *HINTS 4 Cycle 4 Data.csv*. The file contains data from the Health Information National Trends Survey (HINTS) conducted by the US National Cancer Institute. The survey collects nationally representative data on the American public's use of cancer-related information. The survey has been conducted nine times since 2003. The data you will be analysing comes from the 2014 iteration. Codebooks and methodology papers may be found on the HINTS website, which you may access as you wish.

You will be exploring the relationship between the respondent's body mass index and the average daily amount of time the respondent spent on sedentary leisure activities in the past month. You must perform a *simple linear regression* of BMI on sedentary leisure time and report your findings. The data analysis must be conducted in R. You must use the variables `bmi` and `averagedailytvgames` and only these two variables. You must clean the data set, ensuring that only valid data are analysed. You must perform the analysis and produce a regression equation and a suitable graphical summary. You must also test the assumptions of the model. Finally, you must present your R command file. In a written report, you must describe your methods and findings.

2 Instructions

Prepare a written report of no more than 250 words, at least one and no more than five figures and at least one and no more than three tables. The report must show that you are able to

1. describe the methods you've employed (10 marks);
 2. summarise the results of the regression analysis (30 marks);
 3. summarise the results of the diagnostic tests (10 marks);
 4. interpret the findings in the context of the research question (10 marks);
- and

5. use technical English appropriately (10 marks)

The R file must show that you are able to

1. import data from Excel to R (5 points);
2. identify and exclude invalid data from analysis (5 points);
3. perform a simple linear regression (5 points);
4. produce a graphical summary of the data (5 points); and
5. perform tests of the validity of the model (10 points)

3 Submission

You must submit *two* files.

- File 1 is the research report in PDF format. The file must be named < StudentID > .pdf. For example, 123456789.pdf.
- File 2 is the R command file. The file must be named < StudentID > .R. For example, 123456789.R.

You must submit the files via ICE. All penalties for late or incomplete submissions apply.

4 Marking Rubric

Criterion	Level 1	Level 2	Level 3	Level 4
Report: Methods	Not done (0 marks)	Major mistakes or omissions (3 marks)	Minor mistakes or omissions (6 marks)	No mistakes or omissions (10 marks)
Report: Results - Regression equation	Not done (0 marks)	Major mistakes or omissions (3 marks)	Minor mistakes or omissions (6 marks)	No mistakes or omissions (10 marks)
Report: Results - confidence intervals and p-values	Not done (0 marks)	Major mistakes or omissions (3 marks)	Minor mistakes or omissions (6 marks)	No mistakes or omissions (10 marks)
Report: Results - R^2	Not done (0 marks)	Major mistakes or omissions (3 marks)	Minor mistakes or omissions (6 marks)	No mistakes or omissions (10 marks)
Report: Diagnostic tests	Not done (0 marks)	Major mistakes or omissions (3 marks)	Minor mistakes or omissions (6 marks)	No mistakes or omissions (10 marks)
Report: Interpretation	Not done (0 marks)	Major mistakes or omissions (3 marks)	Minor mistakes or omissions (6 marks)	No mistakes or omissions (10 marks)
Report: Technical English	Major mistakes (0 marks)	Minor mistakes (2 marks)	Appropriate (10 marks)	
R: Data import	Not done (0 marks)	Done (5 marks)		
R: Invalid data	Not done (0 marks)	Done (5 marks)		
R: Regression	Not done (0 marks)	Inappropriate (2 marks)	Appropriate (5 marks)	
R: Graphs	Not done (0 marks)	Inappropriate (2 marks)	Appropriate (5 marks)	
R: Model validity	Not done (0 marks)	Inappropriate (4 marks)	Appropriate (10 marks)	