

FIT2086 Assignment 2

Due Date: 11:55PM, Tuesday, 17/9/2024

1 Introduction

There are total of three questions worth $11 + 11 + 8 = 30$ marks in this assignment. This assignment is worth a total of 20% of your final mark, subject to hurdles and any other matters (e.g., late penalties, special consideration, etc.) as specified in the FIT2086 Unit Guide or elsewhere in the FIT2086 Moodle site (including Faculty of I.T. and Monash University policies).

Students are reminded of the Academic Integrity Awareness Training Tutorial Activity and, in particular, of Monash University's policies on academic integrity. In submitting this assignment, you acknowledge your awareness of Monash University's policies on academic integrity and that work is done and submitted in accordance with these policies.

Submission Instructions: Please follow these submission instructions:

1. No files are to be submitted via e-mail. Submissions are to be made via Moodle.
2. You **may not** use generative A.I. in any capacity when answering this assignment.
3. Please provide a **single** file containing your report, i.e., your answers to these questions. Provide code/code fragments as required in your report, and make sure the code is written in a **fixed width font** such as **Courier New** (or a screen shot is taken and inserted – please make sure this is neat and readable), or similar, and is grouped with the question the code is answering. You can submit hand-written answers, but if you do, please make sure they are clear and legible. **Do not submit multiple files** – all your files should be combined into a single PDF file as required. Please ensure that your assignment answers the questions in the **order specified** in the assignment. Multiple files and questions out of order make the life of the tutors marking your assignment much more difficult than it needs to be, and may attract penalties, so please **ensure your assignment follows these requirements**.

Question 1 (11 marks)

The fuel efficiency of cars is usually measured in the number of kilometers (on average) that a car can travel on one litre of fuel, under “typical” conditions. Higher fuel efficiency is obviously desirable. The file `fuel.ass2.2024.csv` contains records on a subset of actual vehicles measured for fuel efficiency by the US government in the period 2017-2020. The data has fuel efficiency recordings on a number of vehicles along with information indicating whether they are either all-wheel-drive (coded as A) or part-time four-wheel-drive (coded as P). Please use this file to answer the following questions.

Important: you may use R to determine the means and variances of the data, as required, and the R functions `qt()` and `pnorm()` but you must perform all the remaining steps by hand. Please provide appropriate R code fragments and all working out.

1. Calculate an estimate of the average fuel efficiency of vehicles that are all-wheel drive. Calculate a 95% confidence interval for this estimate using the t -distribution, and summarise/describe your results appropriately. Show working as required. [4 marks]
2. An obvious and important question is: is there a difference in fuel efficiency between all-wheel-drive vehicles and part-time four-wheel drive vehicles? Using the provided data and the approximate method for difference in means with unknown variances presented in Lecture 4, calculate the estimated mean difference in fuel efficiency between all-wheel-drive vehicles and part-time four-wheel-drive vehicles, and a 95% confidence interval for this difference. Summarise/describe your results appropriately. Show working as required. [3 marks]
3. Given that in an all-wheel-drive vehicle all four wheels are continuously powered, while in a part-time four-wheel-drive vehicle all wheels are powered only at certain times, it seems plausible that all-wheel-drive cars may have worse fuel efficiency. Using the provided data, test the null hypothesis that all-wheel-drives are less efficient than part-time four-wheel-drive vehicles. Write down explicitly the hypothesis you are testing, and then calculate a p -value using the approximate hypothesis test for differences in means with unknown variances presented in Lecture 5. What does this p -value suggest about the difference between vehicles with all-wheel-drive and part-time four-wheel-drive transmissions? Show working as required. [3 marks]
4. Can you identify any possible problems with your conclusions based on the available data? Could there be an alternative explanation for the results you obtained other than their difference in drive-systems (all-wheel-drive *vs* part-time four-wheel-drive)? [1 mark]

Question 2 (11 marks)

The geometric distribution is a probability distribution for non-negative integers. It models the number of tails observed in a sequence of (weighted) coin tosses until the first head is observed. As such it is used widely throughout data science to model the number of times until some specific binary event occurs, i.e, the number of years between major natural disasters, etc. The version that we will look at has a probability mass function of the form

$$p(y | \phi) = (e^\phi + 1)^{-y-1} e^{y\phi} \quad (1)$$

where $y \in \mathbb{Z}_+$, i.e., y can take on the values of non-negative integers. In this form it has one parameter: ϕ , the log-odds of seeing a failure (tail) when the coin is tossed. If a random variable follows a geometric distribution with log-odds ϕ we say that $Y \sim \text{Exp}(\phi)$. If $Y \sim \text{Exp}(\phi)$, then $\mathbb{E}[Y] = e^\phi$ and $\mathbb{V}[Y] = e^\phi(e^\phi + 1)$.

1. Produce a plot of the geometric probability mass function (1) for the values $y \in \{0, 1, \dots, 20\}$, for $\phi = 0$, $\phi = 1$ and $\phi = 2$. Ensure that the graph is readable, the axis are labelled appropriately and a legend is included. [2 marks]
2. Imagine we are given a sample of n observations $\mathbf{y} = (y_1, \dots, y_n)$. Write down the joint probability of this sample of data, under the assumption that it came from a geometric distribution with log-odds parameter ϕ (i.e., write down the likelihood of this data). Make sure to simplify your expression, and provide working. (*hint: remember that these samples are independent and identically distributed.*) [2 marks]
3. Take the negative logarithm of your likelihood expression and write down the negative log-likelihood of the data \mathbf{y} under the geometric model with log-odds ϕ . Simplify this expression. [1 mark]
4. Derive the maximum likelihood estimator $\hat{\phi}$ for ϕ . That is, find the value of ϕ that minimises the negative log-likelihood. You must provide working. [3 marks]
5. Determine the approximate bias and variance of the maximum likelihood estimator $\hat{\phi}$ of ϕ for the geometric distribution. (*hints: utilise techniques from Lecture 2, Slide 21 and the mean/variance of the sample mean*) [3 marks]

Question 3 (8 marks)

It is frequent in nature that animals express certain asymmetries in their behaviour patterns. It has been suggested that this might be nature's way of "breaking gridlocks" that might occur if we were to act purely rationally (think: why does a beetle decide to move one way over another when put in a featureless bowl?).

An interesting study regarding preferences was undertaken by Irish researchers in 2006. In the experiment, 240 volunteer students from Stanmillis University College in Belfast were asked to stand directly in front of a symmetrical doll's face and asked to kiss the doll on the cheek or lips; researchers then recorded whether the student tilted their head to the right or left when kissing the doll. Of the 240 students, 168 turned their head to the right and 72 turned their head to the left. You must analyse this data to see if there is an inbuilt preference in humans for the direction of head tilt when kissing. Provide working, reasoning or explanations and R commands that you have used, as appropriate.

1. Calculate an estimate of the preference for humans turning their heads to the right when kissing using the above data, and provide an approximate 95% confidence interval for this estimate. Summarise/describe your results appropriately. **[3 marks]**
2. Test the hypothesis that there is no preference in humans for tilting their head to one particular side when kissing. Write down explicitly the hypothesis you are testing, and then calculate a p -value using the approximate approach for testing a Bernoulli population discussed in Lecture 5. What does this p -value suggest? **[2 marks]**
3. Using R, calculate an exact p -value to test the above hypothesis. What does this p -value suggest? Please provide the appropriate R command that you used to calculate your p -value. **[1 mark]**
4. It is entirely possible that any preference for head turning to the right/left could be simply a product of right/left-handedness. To test this the handedness of the 240 volunteers was also recorded. It was found that 213 of the participants were right-handed and 27 were left handed. Using the approximate hypothesis testing procedure for testing two Bernoulli populations from Lecture 5, test the hypothesis that the rate of right-handedness in the population from which the participants was drawn is the same as the preference for turning heads to the right when kissing. Summarise your findings. What does the p -value suggest? **[2 marks]**