

Statistical Computing: Midterm Exam

2023-10-18

Name & email:

Instructions:

1. The exam can be found in a pdf and Rmarkdown form in Github.
2. Use the Rmarkdown form to answer the questions of the exam, using R.
3. You will need to submit your Rmarkdown file and the knit PDF file as your solutions to this exam, in D2L->Assignments-> Midterm Exam.
4. The exam is prepared so that you can complete it by 4:30pm. Notice that you will have until 5 pm to upload your submission on D2L.
5. There are 3 questions in the exam, each in a different page. Each question has sub-questions. The exam is with open notes. Points for each subquestion can be found in a parenthesis.

Question 1: Maximum Likelihood Estimation for the parameters α, β of the Gamma distribution.

Let X be a random variable that follows the $\text{Gamma}(\alpha, \beta)$, where α is the shape parameter and β is the rate parameter. Then the density function of X is:

$$f(X = x | \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}.$$

If $X_1, X_2, X_3, \dots, X_n$ are n independent and identically distributed random variables that follow the same $\text{Gamma}(\alpha, \beta)$, then the likelihood function of $X_1 = x_1, X_2 = x_2, X_3 = x_3, \dots, X_n = x_n$ is:

$$L(\alpha, \beta | x_1, x_2, x_3, \dots, x_n) = \prod_{i=1}^n \frac{\beta^\alpha}{\Gamma(\alpha)} x_i^{\alpha-1} e^{-\beta x_i} \text{ and the loglikelihood is}$$

$$\ell(\alpha, \beta | x_1, x_2, x_3, \dots, x_n) = \sum_{i=1}^n \log\left(\frac{\beta^\alpha}{\Gamma(\alpha)} x_i^{\alpha-1} e^{-\beta x_i}\right) = n \log\left(\frac{\beta^\alpha}{\Gamma(\alpha)}\right) + (\alpha - 1) \sum_{i=1}^n \log(x_i) - \beta \sum_{i=1}^n x_i.$$

In R:

- You can use the R function `gamma(x)` to calculate values of $\Gamma(x)$.
- For your convenience you can find the R functions for the Gamma distribution here.
- Notice that for this problem we define a Gamma distribution only by the shape and rate parameter and we ignore the scale parameter.

1.1) (18 points) Write a function that evaluates the negative log-likelihood of a random sample of IID $\text{Gamma}(\alpha, \beta)$ distributed random variables.

1.2) (18 points) Using your function from 1.1) and the following data find the maximum likelihood estimators of α and β .

- Use $\alpha = 1, \beta = 1$ as starting parameters.
- Report your parameter estimates.

```
x = as.matrix(read.table('https://raw.githubusercontent.com/gdlc/STAT_COMP/master/DATA/Gamma_data.txt'))  
head(x)
```

1.3) (18 points) Provide Standard Errors, and 95% Confidence Intervals for each of the parameters.

Question 2: Logistic regression

Below you are provided with a simulated data set containing information on ten thousand customers. Specifically, for each customer we are given information about the following variables:

1. *default* : A factor with levels No and Yes indicating whether the customer defaulted on their debt. Default usually happens after six months in a row of not making at least the minimum payment due.
2. *student* : A factor with levels No and Yes indicating whether the customer is a student
3. *balance*: The average balance that the customer has remaining on their credit card after making their monthly payment
4. *income* : Income of customer

```
library(ISLR2)
Default_data = Default
head(Default_data)
```

2.1) (18 points) We model the probability of defaulting based on the variables *student*, *balance* and *income*. Using the logistic regression model $\text{default} \sim \text{student} + \text{balance} + \text{income}$, estimate the probability that a student with median balance and median income will default. To do that follow the steps below:

- Fit a logistic regression model and
- Find the median of the variable balance and the median of the variable income.
- Find the predicted probability of defaulting for a student with balance equal with the median balance of the data and income equal with the median income of the data.
- Construct 95% Confidence interval for the prediction of the probability of defaulting for a student with median balance and median income.
- Based on your results would you expect a student with median income and balance to default?

Question 3: Linear regression

A statistics professor regularly wears a smart watch to keep track of their steps and calories burnt within a day. The smart watch records not only the number of steps taken each day, but also the number of minutes walked at a moderate pace, and the number of miles total that they walked.

The professor created a data set of the above information for each day and added the variable *Rain* to record whether it was a rainy, sunny, or cold day.

The following dataset has 68 observations on the following 8 variables.

1. *Steps*: Total number of steps for the day
2. *Moderate*: Number of steps at a moderate walking speed
3. *Min*: Number of minutes walking at a moderate speed
4. *kcal*: Number of calories burned walking at a moderate speed
5. *Mile*: Total number of miles walked
6. *Rain*: Type of weather (rain or shine)
7. *Day*: Day of the week (U=Sunday, M=Monday, T=Tuesday, W=Wednesday, R=Thursday, F=Friday, S=Saturday)
8. *DayType*: Coded as Weekday or Weekend

```
path='https://vincentarelbundock.github.io/Rdatasets/csv/Stat2Data/Pedometer.csv'

walk_data = read.csv(path, row.names = 1)
```

3.1) (10 points) Use the *walk_data* to fit linear regression model that predicts calories *kcal* burnt on a day from walking as function of *Steps* and *Min*. Summarize in no more than 2 sentence your conclusions regarding the fit of this model.

3.2) (18 points) The professor would typically bike on sunny days, so they would walk less on those days. At the same time, if it was raining they would choose to walk instead of biking. Based on that, the variable *Rain* could possibly affect the variables *Steps* and the variable *kcal*.

-Test if adding the variable *Rain* and the interaction of *Rain* with *Steps* to the above model, would help the prediction of *kcal*. What is your conclusion?