

Problem Set 1: Written Part

Your Name Goes Here

Details

How to Write Up

The written part of this assignment can be either typeset using latex or hand written.

Grading

5% of your grade on this assignment is for turning in something legible. This means it should be organized, and any Rmd files should knit to pdf without issue.

An additional 15% of your grade is for completion. A quick pass will be made to ensure that you've made a reasonable attempt at all problems.

Across both the written part and the R part, in the range of 1 to 3 problems will be graded more carefully for correctness. In grading these problems, an emphasis will be placed on full explanations of your thought process. You don't need to write more than a few sentences for any given problem, but you should write complete sentences! Understanding and explaining the reasons behind what you are doing is at least as important as solving the problems correctly.

Solutions to all problems will be provided.

Collaboration

You are allowed to work with others on this assignment, but you must complete and submit your own write up. You should not copy large blocks of code or written text from another student.

Sources

You may refer to our text, Wikipedia, and other online sources. All sources you refer to must be cited in the space I have provided at the end of this problem set.

Problem I: χ^2 , t , and F distributions

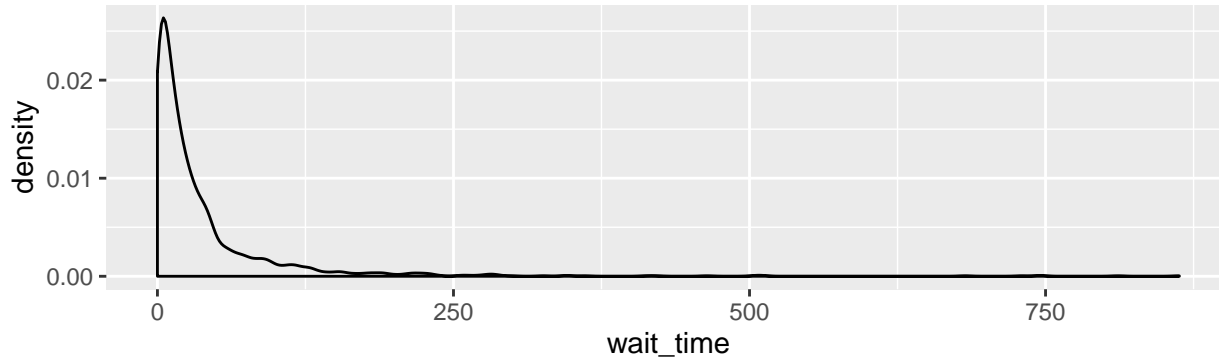
(1) Suppose that $Z_1, Z_2, Z_3, Z_4 \stackrel{\text{i.i.d.}}{\sim} \text{Normal}(0, 1)$, and $Y = \frac{Z_1}{\sqrt{(Z_2^2 + Z_3^2 + Z_4^2)/3}}$. What is the distribution of Y ? What is the distribution of Y^2 ? Justify your answers briefly.

(2) Let $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \text{Normal}(\mu_X, \sigma^2)$ and $Y_1, \dots, Y_m \stackrel{\text{i.i.d.}}{\sim} \text{Normal}(\mu_Y, \sigma^2)$, with all X 's and Y 's independent. Note that the two distributions have different means but the same variance. Show how you can use the F distribution to find $P(S_X^2/S_Y^2 > 2)$, where S_X^2 is the sample variance of X_1, \dots, X_n and S_Y^2 is the sample variance of Y_1, \dots, Y_m .

Problem II: Emergency Department Waiting Times

The National Center for Health Statistics, a division within the U.S. Centers for Disease Control, conducts a nationally representative survey of hospitals each year to track the waiting times for emergency room visits (that is, how much time passed between when a patient arrived at the hospital and when they were seen by a doctor or registered nurse). In this problem, we will model the distribution of waiting times for 1874 emergency department visits from 2012. The output below shows a plot of the data.

```
ggplot(data = er_visits, mapping = aes(x = wait_time)) +
  geom_density()
```



An exponential distribution is often used to model waiting times. Suppose we adopt the data model

$X_i \stackrel{\text{i.i.d.}}{\sim} \text{Exponential}(\theta), i = 1, \dots, n,$

where X_i is the waiting time for visit number i . In our data set we have observed values x_1, \dots, x_n , where $n = 1874$.

Here are some facts about the exponential distribution (note that there are two common parameterizations of the exponential distribution; please work with the definitions and properties stated below for this problem):

parameter	$\theta > 0$: scale parameter
p.f.	$f_{X \Theta}(x \theta) = \theta^{-1}e^{-x/\theta}$ on the support $x \geq 0$
Mean	θ
Variance	θ^2

(1) Find the maximum likelihood estimator of θ .

(2) Is your result from part (1) a random variable or a number? If it is a random variable, explain why it is random. If it is a number, explain why it is not random. (Your answer should be more detailed than “because that’s the definition of an estimator.”)

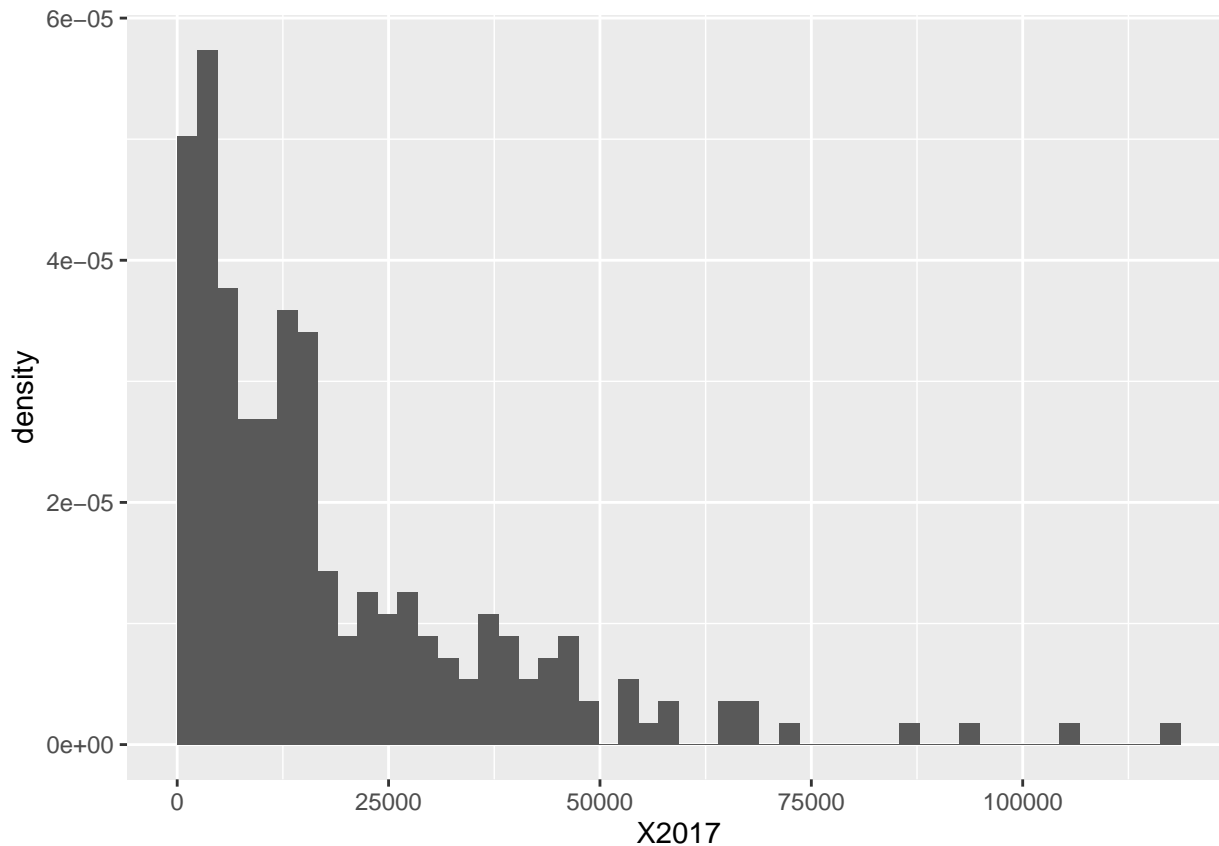
Problem III: Per Capita GDP

The code below reads in and plots a data set with measurements of per capita GDP at purchasing power parity as of 2017 for 235 countries, measured in inflation-adjusted 2011 international dollars; these data are from the World Bank, here: <https://data.worldbank.org/indicator/NY.GDP.PCAP.PP.KD>. Per capita GDP can be roughly interpreted as the amount of income generated in a country in one year divided by the number of people living in that country. The purchasing power parity adjustment attempts to adjust GDP to account for differences in cost of living in different countries.

```
gdp <- read.csv("http://www.evanlray.com/data/worldbank/worldbank_percap_gdp_ppp.csv")
```

```
gdp <- gdp %>%
  filter(!is.na(X2017))
```

```
ggplot(data = gdp, mapping = aes(x = X2017)) +
  geom_histogram(mapping = aes(y = ..density..), boundary = 0, bins = 50)
```



A lognormal distribution is often used to model non-negative variables that are skewed right, like incomes. In the written part of this assignment you will find the maximum likelihood estimator for the parameters of a lognormal distribution, and in the R part of the assignment you will fit the model to this data set.

For the purpose of this assignment, let's assume that the per capita GDP of different countries in a given year can be modelled as independent, identically distributed random variables (this is not actually reasonable, but may be good enough if our goal is to describe the distribution of values for per capita GDP across different countries).

Let's adopt the model $X_i \stackrel{i.i.d.}{\sim} \text{lognormal}(\mu, \sigma)$, $i = 1, \dots, n$.

The pdf of a lognormal distribution is given by $f(x|\mu, \sigma) = x^{-1}(2\pi\sigma^2)^{-\frac{1}{2}} \exp\left[-\frac{1}{2} \frac{\{\log(x)-\mu\}^2}{\sigma^2}\right]$

(1) Find the maximum likelihood estimators of μ and σ . For this problem, you do not have to check second-order conditions to verify that you have found a global maximum of the log-likelihood function.