

# Homework 1

PSTAT 115, Winter 2023

Due on Sunday, January 22 at 11:59 pm

**Note:** You may work with a partner but you must writeup and submit your own assignment. Submit your Rmarkdown (.Rmd) and the compiled pdf on Gradescope in a zip file. Include any addition files (e.g. scanned handwritten solutions) in zip file with the pdf.

## Text Analysis of JK Rowling's Harry Potter Series

### Question 1

You are interested in studying the writing style and tone used by JK Rowling (JKR for short), the author of the popular Harry Potter series. You select a random sample of chapters of size  $n$  from all of JKR's books. You are interested in the rate at which JKR uses the word *fire* in her writing, so you count how many times the word *fire* appears in each chapter in your sample,  $(y_1, \dots, y_n)$ . In this set-up,  $y_i$  is the number of times the word *fire* appeared in the  $i$ -th randomly sampled chapter. In this context, the population of interest is all chapters written by JKR and the population quantity of interest (the estimand) is the rate at which JKR uses the word *fire*. The sampling units are individual chapters. Note: this assignment is partially based on text analysis package known as [tidytext](#). You can read more about tidytext [here](#).

1a.

Model: let  $Y_i$  denote the quantity that captures the number of times the word *fire* appears in the  $i$ -th chapter. As a first approximation, it is reasonable to model the number of times *fire* appears in a given chapter using a Poisson distribution. *Reminder:* Poisson distributions are for integer outcomes and useful for events that occur independently and at a constant rate. Let's assume that the quantities  $Y_1, \dots, Y_n$  are independent and identically distributed (IID) according to a Poisson distribution with unknown parameter  $\lambda$ ,

$$p(Y_i = y_i \mid \lambda) = \text{Poisson}(y_i \mid \lambda) \quad \text{for } i = 1, \dots, n.$$

Write the likelihood  $L(\lambda)$  for a generic sample of  $n$  chapters,  $(y_1, \dots, y_n)$ . Simplify as much as possible (i.e. get rid of any multiplicative constants)

Likelihood Model

$$L(\mu) = \prod_{i=1}^n \frac{\lambda^{y_i} * e^{-\lambda}}{y_i!} \propto \lambda^{\sum y_i} * e^{-\lambda n}$$

## 1b.

Write the log-likelihood  $\ell(\lambda)$  for a generic sample of  $n$  articles,  $(y_1, \dots, y_n)$ . Simplify as much as possible. Use this to compute the maximum likelihood estimate for the rate parameter of the Poisson distribution.

Log-Likelihood

$$l(\lambda) = \log\left(\prod_{i=1}^n \frac{\lambda^{y_i} * e^{-\lambda}}{y_i!}\right) = \sum_{i=1}^n \log(\lambda^{y_i} * \frac{1}{y_i!} * e^{-\lambda}) =$$

$$\sum_{i=1}^n (\log(\lambda^{y_i}) - \log(y_i!) + \log(e^{-\lambda})) = \sum_{i=1}^n (y_i \log(\lambda) - \log(y_i!) - \lambda) = \sum y_i \log(\lambda) - \sum \log(y_i!) - n\lambda$$

MLE

$$\hat{\lambda}_{mle} = \frac{d}{d\lambda} l(\lambda) = 0 \Rightarrow \frac{d}{d\lambda} \sum y_i \log(\lambda) - \sum \log(y_i!) - n\lambda = \frac{\sum y_i}{\lambda} - n = 0 \Rightarrow \hat{\lambda}_{mle} = \frac{\sum y_i}{n}$$

From now on, we'll focus on JKR's writing style in the Harry Potter book, *The Goblet of Fire*. This book has 37 chapters. Below is the code for counting the number of times *fire* appears in each chapter of *The Goblet of Fire*. We use the `tidytext` R package which includes functions that parse large text files into word counts. The code below creates a vector of length 37 which has the number of times the word *fire* was used in that chapter (see [https://uc-r.github.io/tidy\\_text](https://uc-r.github.io/tidy_text) for more on parsing text with `tidytext`)

```
library(tidyverse)      # data manipulation & plotting
library(stringr)        # text cleaning and regular expressions
library(tidytext)       # provides additional text mining functions
library(harrypotter)    # text for the seven novels of the Harry Potter series

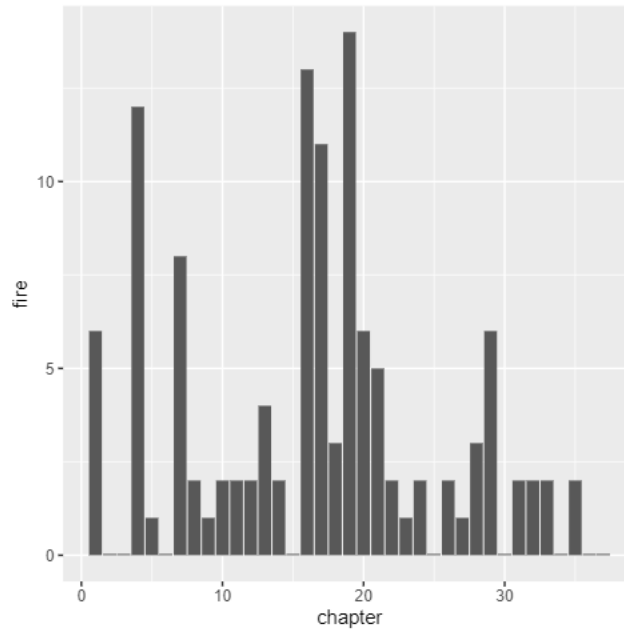
text_tb <- tibble(chapter = seq_along(goblet_of_fire),
                  text = goblet_of_fire)
tokens <- text_tb %>% unnest_tokens(word, text)
word_counts <- tokens %>% group_by(chapter) %>%
  count(word, sort = TRUE) %>% ungroup
word_counts_mat <- word_counts %>% spread(key=word, value=n, fill=0)

fire_counts <- word_counts_mat$fire
```

1c.

Make a bar plot where the heights are the counts of the word *fire* and the x-axis is the chapter.

```
ggplot(data = word_counts_mat, aes(x = chapter, y = fire)) + geom_bar(stat = 'identity')
```



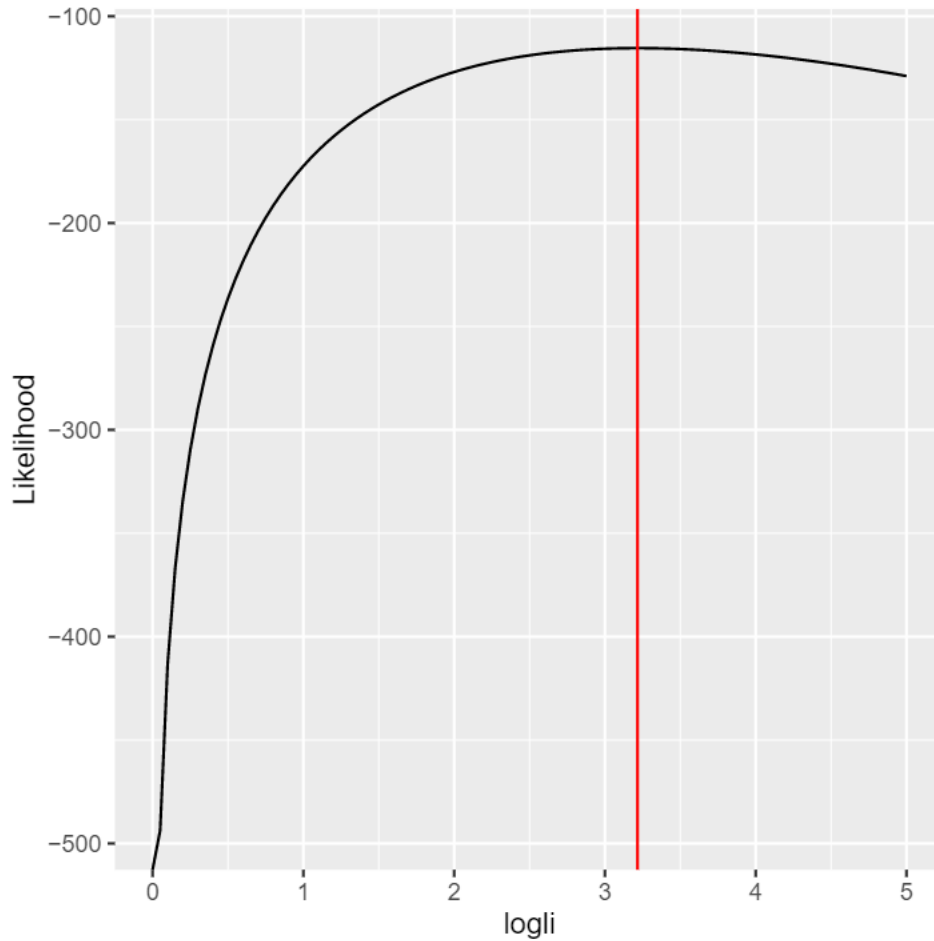
1d.

Plot the log-likelihood of the Poisson rate of *fire* usage in R using the data in `fire_counts`. Then use `fire_counts` to compute the maximum likelihood estimate of the rate of the usage of the word *fire* in The Goblet of Fire. Mark this maximum on the log-likelihood plot with a vertical line (use `abline` if you make the plot in base R or `geom_vline` if you prefer `ggplot`).

```
f = function(logli){
  (sum(fire_counts) * log(logli)) - (37 * logli) - sum(log(factorial(fire_counts)))
}

logli <- ggplot(data = data.frame(logli = 0), mapping = aes(x=logli))+stat_function(fun = f)+
  xlim(0,5)+scale_y_continuous(name = "Likelihood")

logli + geom_vline(xintercept = (sum(fire_counts) / 37), color = 'red')
```



`#fire_counts`

## Question 2

For the previous problem, when computing the rate of *fire* usage, we were implicitly assuming each chapter had the same length. Remember that for  $Y_i \sim \text{Poisson}(\lambda)$ ,  $E[Y_i] = \lambda$  for each chapter, that is, the average number of occurrences of *fire* is the same in each chapter. Obviously this isn't a great assumption, since the lengths of the chapters vary; longer chapters should be more likely to have more occurrences of the word. We can augment the model by considering properties of the Poisson distribution. The Poisson is often used to express the probability of a given number of events occurring for a fixed "exposure". As a useful example of the role of the exposure term, when counting then number of events that happen in a set length of time, we need to account for the total time that we are observing events. For this text example, the exposure is not time, but rather corresponds to the total length of the chapter.

We will again let  $(y_1, \dots, y_n)$  represent counts of the word *fire*. In addition, we now count the total number of words in each chapter  $(\nu_1, \dots, \nu_n)$  and use this as our exposure. Let  $Y_i$  denote the random variable for the counts of the word *fire* in a chapter with  $\nu_i$  words. Let's assume that the quantities  $Y_1, \dots, Y_n$  are independent and identically distributed (IID) according to a Poisson distribution with unknown parameter  $\lambda \cdot \frac{\nu_i}{1000}$ ,

$$p(Y_i = y_i \mid \nu_i, 1000) = \text{Poisson}(y_i \mid \lambda \cdot \frac{\nu_i}{1000}) \quad \text{for } i = 1, \dots, n.$$

In the code below, `chapter_lengths` is a vector storing the length of each chapter in words.

```
chapter_lengths <- word_counts %>% group_by(chapter) %>%
  summarize(chapter_length = sum(n)) %>%
  ungroup %>% select(chapter_length) %>% unlist %>% as.numeric
```

## 2a.

What is the interpretation of the quantity  $\frac{\nu_i}{1000}$  in this model? What is the interpretation of  $\lambda$  in this model? State the units for these quantities in both of your answers.

$\frac{\nu_i}{1000}$  is the number of words per chapter divided by 1000. In this model  $\lambda$  is the number of times “fire” appears per 1000 words.

## 2b.

List the known and unknown variables and constants, as described in lecture 2. Make sure you include  $Y_1, \dots, Y_n, y_1, \dots, y_n, n, \lambda$ , and  $\nu_i$ .

The known variables with variance  $> 0$  are  $Y_i$ . The known variables with variance  $= 0$  are  $y_i, \nu_i$ , and  $n$ . There are no unknown variables with variance  $> 0$ . The unknown variable with variance  $= 0$  is  $\lambda$ .

## 2c.

Write down the likelihood in this new model. Use this to calculate maximum likelihood estimator for  $\lambda$ . Your answer should include the  $\nu_i$ 's.

$$L(\lambda) = \prod_{i=1}^n \frac{\frac{\lambda \nu_i}{1000}^{y_i} * e^{-\frac{\lambda \nu_i}{1000}}}{y_i!} \propto \prod_{i=1}^n \frac{\lambda \nu_i^{y_i}}{1000^{y_i}} * e^{-\frac{\lambda \nu_i}{1000}}$$

The maximum likelihood estimator would be

$$\begin{aligned} \frac{d}{d\lambda} l(\lambda) = 0 &\Rightarrow \log\left(\prod_{i=1}^n \frac{\frac{\lambda \nu_i}{1000}^{y_i} * e^{-\frac{\lambda \nu_i}{1000}}}{y_i!}\right) = \Sigma(\log(\frac{\lambda \nu_i}{1000}^{y_i} * e^{-\frac{\lambda \nu_i}{1000}})) = \\ &\Sigma(\log(\frac{\lambda \nu_i}{1000}^{y_i}) - \log(y_i!) + \log(e^{-\frac{\lambda \nu_i}{1000}})) = \Sigma(y_i \log \frac{\lambda \nu_i}{1000}) - \log(y_i!) - \frac{\lambda \nu_i}{1000} = \\ &\Sigma(y_i \log \frac{\lambda \nu_i}{1000}) - \Sigma(y_i!) - \Sigma(\frac{\lambda \nu_i}{1000}) \Rightarrow \frac{d}{d\lambda} \Sigma(y_i \log \frac{\lambda \nu_i}{1000}) - \log(y_i!) - \frac{\lambda \nu_i}{1000} = \\ &\frac{\Sigma y_i}{\lambda} - \frac{\Sigma \nu_i}{1000} = 0 \Rightarrow \lambda = \frac{1000 * \Sigma y_i}{\Sigma \nu_i} \end{aligned}$$

## 2d.

Compute the maximum likelihood estimate and save it in the variable `lambda_mle`. In 1-2 sentences interpret its meaning (make sure you include units in your answers!).

```
lambda_mle <- (1000 * sum(fire_counts)) / sum(chapter_lengths)
```

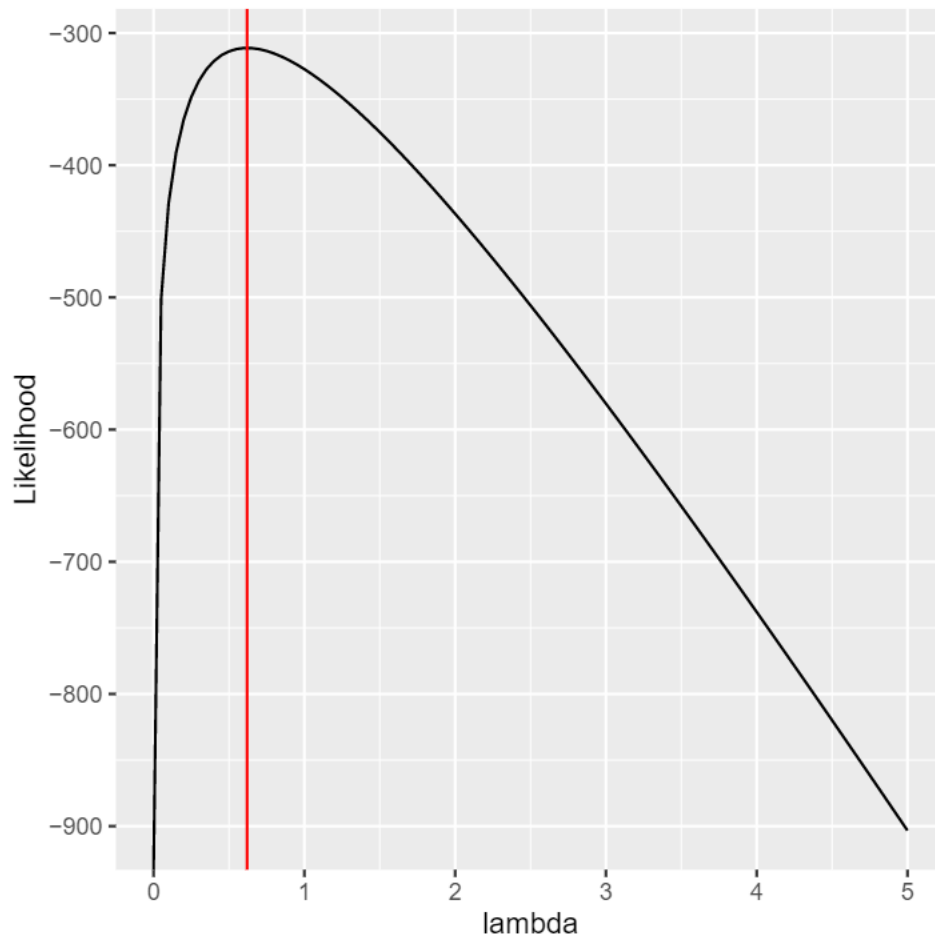
```
. = ottr::check("tests/q2d.R")
```

The mle calculated here is 0.62 which means that a  $\lambda$  value of 0.62 makes the data most likely. It means that for every 1000 pages the word *fire* will come up 0.62 times.

2e.

Plot the log-likelihood from the previous question in R using the data from on the frequency of *fire* and the chapter lengths. Add a vertical line at the value of `lambda_mle` to indicate the maximum likelihood.

```
f2 = function(lambda){  
  sum(fire_counts) * log(lambda) - sum(log(factorial(fire_counts))) - sum(chapter_lengths) * (lambda /  
}  
  
g <- ggplot(data = data.frame(lambda = 0.5), mapping = aes(x=lambda))+stat_function(fun = f2)+  
  xlim(0,5)+scale_y_continuous(name = "Likelihood")  
  
g + geom_vline(xintercept = (lambda_mle), color = 'red')
```



### Question 3

Correcting for chapter lengths is clearly an improvement, but we're still assuming that JKR uses the word *fire* at the same rate in all chapters. In this problem we'll explore this assumption in more detail.

3a.

Why might it be unreasonable to assume that the rate of *fire* usage is the same in all chapters? Comment in a few sentences.

It might be unreasonable to assume that the rate of *fire* usage is the same in all chapters because different chapters will be about different topics. For example, the introduction can be vastly different from the ending which could lead to different rates of *fire*. In addition certain chapters could have more dialogue which could lead to a lower rate of *fire*.

### 3b.

We can use simulation to check our Poisson model, and in particular the assumption that the rate of *fire* usage is the same in all chapters. Generate simulated counts of the word *fire* by sampling counts from a Poisson distribution with the rate  $(\hat{\lambda}_{MLE}\nu_i)/1000$  for each chapter  $i$ .  $\hat{\lambda}_{MLE}$  is the maximum likelihood estimate computing in 2d. Store the vector of these values for each chapter in a variable of length 37 called `lambda_chapter`. Make a side by side plot of the observed counts and simulated counts and note any similarities or differences (we've already created the observed histogram for you). Are there any outliers in the observed data that don't seem to be reflected in the data simulated under our model?

```
observed_histogram <- ggplot(word_counts_mat) + geom_histogram(aes(x=fire)) +
  xlim(c(0, 25)) + ylim(c(0,7.5)) + ggtitle("Observed")

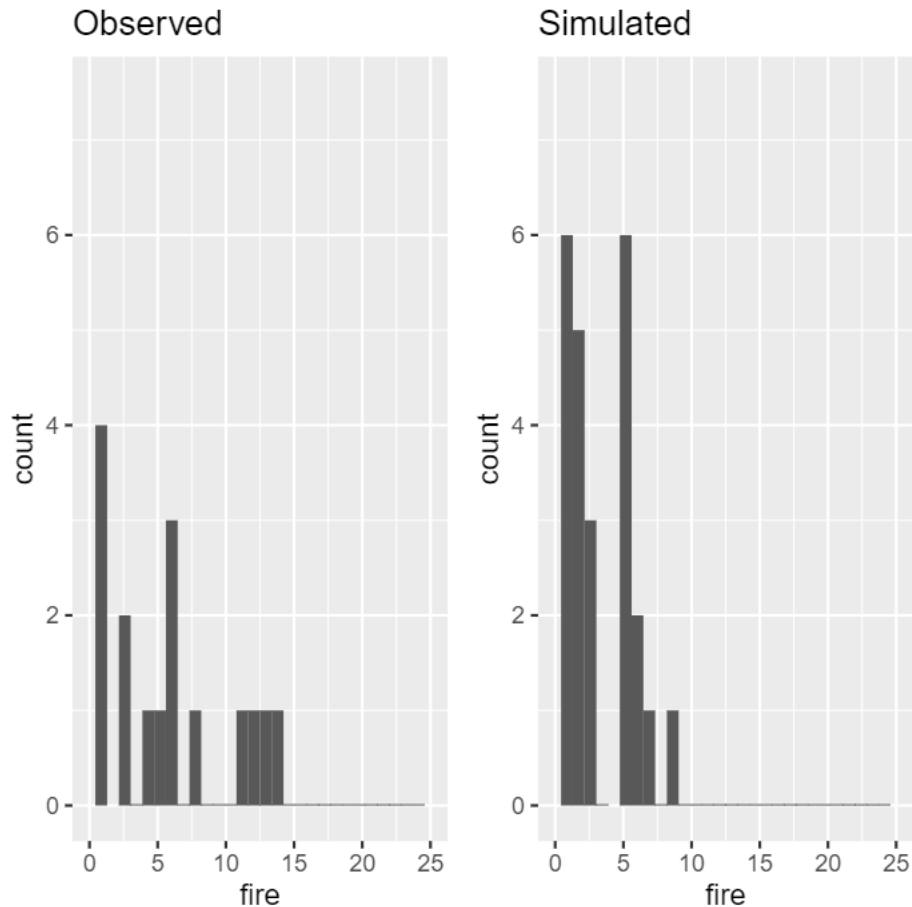
lambda_chapter <- lambda_mle * chapter_lengths / 1000

simulated_counts <- tibble(fire = rpois(37, lambda_chapter))

simulated_histogram <- ggplot(simulated_counts) + geom_histogram(aes(x=fire)) +
  xlim(c(0, 25)) + ylim(c(0,7.5)) + ggtitle("Simulated")

## This uses the patchwork library to put the two plots side by side
observed_histogram + simulated_histogram

## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## Warning: Removed 3 rows containing missing values (`geom_bar()`).
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## Warning: Removed 3 rows containing missing values (`geom_bar()`).
```



```
. = ottr::check("tests/q3b.R")
```

```
##
```

```
## All tests passed!
```

The simulated graph appears to have more counts more frequently. In addition, the observed graph appears to have values when  $\text{fire} > 10$ , but on our simulated graph there are none. It is possible that the observed graph had these outliers which are not in the simulated one, but it could also be just by chance this time that it didn't.

3c. Assume the word usage rate varies by chapter, that is,

$$p(Y_i = y_i \mid \lambda, \nu_i, 1000) = \text{Poisson}(y_i \mid \lambda_i \cdot \frac{\nu_i}{1000}) \quad \text{for } i = 1, \dots, n.$$

Compute a separate maximum likelihood estimate of the rate of *fire* usage (per 1000 words) in each chapter,  $\hat{\lambda}_i$ . Make a bar plot of  $\hat{\lambda}_i$  by chapter. Save the chapter-specific MLE in a vector of length 37 called `lambda_hats`. Which chapter has the highest rate of usage of the word *fire*? Save the chapter number in a variable called `most_fiery_chapter`.

```
# Maximum likelihood estimate
lambda_hats <- 1000 * fire_counts / chapter_lengths

most_fiery_chapter <- which.max(lambda_hats)

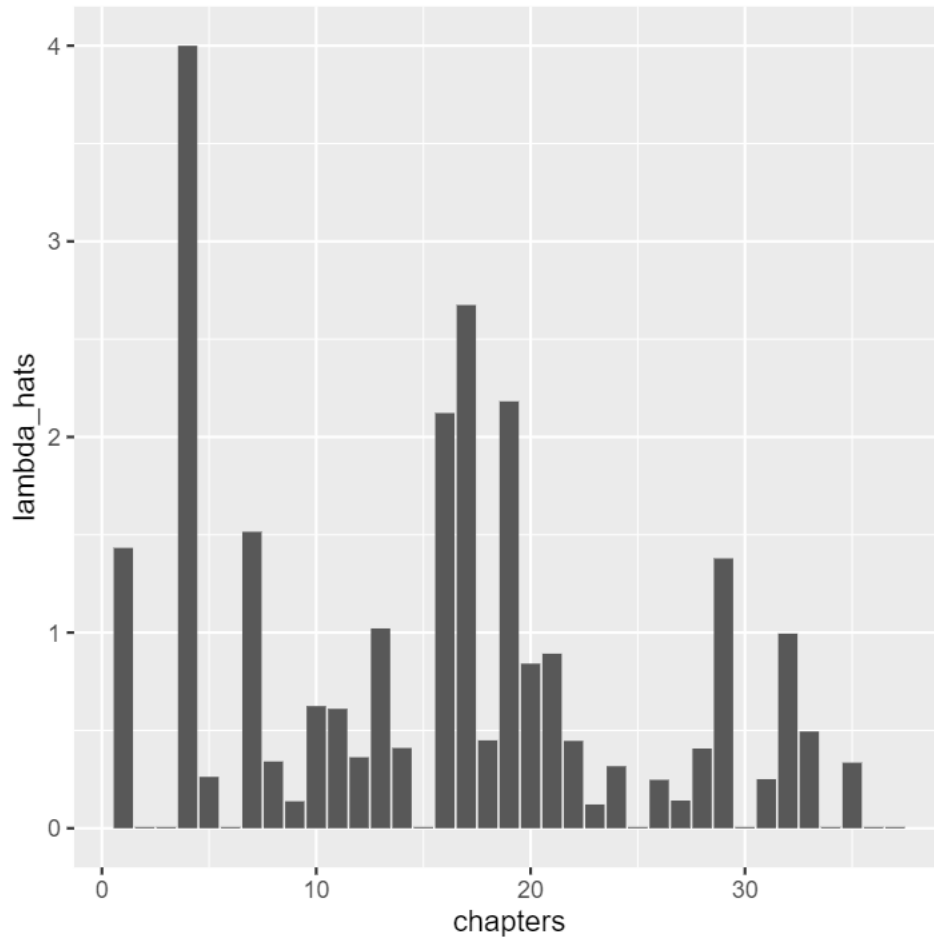
# Make a bar plot of the MLEs, lambda_hats
# YOUR CODE HERE
```



```

chapters <- seq(1:37)
df <- data.frame(lambda_hats, chapters)
plot <- ggplot(data = df, aes(x = chapters, y = lambda_hats)) + geom_bar(stat = 'identity')
plot

```



*#Chapter 4 has the highest rate of usage of the word fire*

```
. = ottr::check("tests/q3c.R")
```

```
##
```

```
## All tests passed!
```

## Question 4

Let's go back to our original model for usage rates of the word *fire*. You collect a random sample of book chapters penned by JKR and count how many times she uses the word *fire* in each of the chapter in your sample,  $(y_1, \dots, y_n)$ . In this set-up,  $y_i$  is the number of times the word *fire* appeared in the  $i$ -th chapter, as before. However, we will no longer assume that the rate of use of the word *fire* is the same in every chapter. Rather, we'll assume JKR uses the word *fire* at different rates  $\lambda_i$  in each chapter. Naturally, this makes sense, since different chapters have different themes and tone. To do this, we'll further assume that the rate of word usage  $\lambda_i$  itself, is distributed according to a  $\text{Gamma}(\alpha, \beta)$  with known parameters  $\alpha$  and  $\beta$ ,

$$f(\Lambda = \lambda_i \mid \alpha, \beta) = \text{Gamma}(\lambda_i \mid \alpha, \beta).$$

and that  $Y_i \sim \text{Pois}(\lambda_i)$  as in problem 1. For now we will ignore any exposure parameters,  $\nu_i$ . Note: this is a “warm up” to Bayesian inference, where it is standard to treat parameters as random variables and specify distributions for those parameters.

4a.

Write out the the data generating process for the above model. REDO REDO REDO

The data generating process for the above model would include randomly sampling from a Gamma distribution with known  $\alpha$  and  $\beta$  1000 times and seeing the results. It could look something like:

```
#for i in range(1:1000):
#  sample.append(rgamma($\alpha$, $\beta$))

#return sample
```

4b.

In R simulate 1000 values from the above data generating process, assume  $\alpha = 10$  (shape parameter of `rgamma`) and  $\beta = 1$  (rate parameter of `rgamma`). Store the value in a vector of length 1000 called `counts`. Compute the empirical mean and variance of values you generated. For a Poisson distribution, the mean and the variance are the same. In the following distribution is the variance greater than the mean (called **overdispersed**''') or is the variance less than the mean (underdispersed'')? Intuitively, why does this make sense?

```
## Store simulated data in a vector of length 1000
```

```
counts <- rgamma(1000, shape = 10, rate = 1)
```

```
print(mean(counts))
```

```
## [1] 10.00413
```

```
print(var(counts))
```

```
## [1] 10.00831
```

```
. = ottr::check("tests/q4b.R")
```

```
##
```

```
## All tests passed!
```

In our case the variance is slightly greater than the mean so our distribution is overdispersed. This makes sense as we did take random samples and there could have been a single relatively small outlier which slightly skews the variance to be greater than the mean.

4c.

List the known and unknown variables and constants as described in lecture 2. Make sure your table includes  $Y_1, \dots, Y_n, y_1, \dots, y_n, n, \lambda, \alpha$ , and  $\beta$ .

The known variables with variance  $> 0$  are  $Y_i$ . The known variables with variance  $= 0$  are  $y_i, \alpha, \beta$ , and  $n$ . The unknown variable with variance  $> 0$  is  $\lambda$ . There are no unknown variables with variance  $= 0$ .