

Practical 1: A Markov Text Model

New translations of the bible into modern English always arouse controversy. Proponents of the new translation argue that it makes the text more understandable, while opponents argue that the new text merely removes the rhythm and poetry of the language while making it no more accessible. There are various ways of testing claims about understandability. One test is to see how readily readers can distinguish genuine text from the book, with computer generated text designed to match word patterns seen in the book, but generated randomly so that it contains no actual meaning.

This practical is about creating such computer generated text. The idea is to use a Markov model — that is a model in which we generate words sequentially, with the each word being drawn with a probability dependent on the word preceding it. The probabilities are obtained by training the model on the actual text. That is by simply tabulating the frequency with which each word follows any other word.

To make this work requires simplification. The model will not cover every word used in the text. Rather the model's 'vocabulary' will be limited to the m most common words. $m \approx 1000$ is sensible. Suppose that the m most common words are in a vector \mathbf{b} . We will construct an $m \times m$ matrix \mathbf{A} , such that

$$P(b_j|b_i) = A_{ij}.$$

Given \mathbf{A} , \mathbf{b} and a starting word from \mathbf{b} , you can then iterate to generate text from the model. That is, given a word b_i the i th row of \mathbf{A} gives the probabilities of each of the other words following it, so we choose the next word from \mathbf{b} randomly according to these probabilities. To estimate \mathbf{A} you need to go through the text of the bible, counting up the number of times b_j follows b_i for all words in \mathbf{b} .

Because this is the first practical, the instructions for how to produce code will be unusually detailed: the task has been broken down for you. Obviously the process of breaking down a task into constituent parts before coding is part of programming, so in future practicals you should expect less of this detailed specification.

As a group of 3 you should aim to produce well commented, clear and efficient code for training the model and simulating short sections of text using it. The code should be written in a plain text file and is what you will submit. Your solutions should use only the functions available in base R and the recommended package `mgcv` (which is installed with base R by default). The work must be completed in your work group of 3, which you must have arranged and registered on Learn. The first comment in your code should list the names and university user names of each member of the group.

1. Create a repo for this project on github, and clone to your local machines.
2. Download the text as plain text from <https://www.gutenberg.org/ebooks/1581>.
3. The following code will read the file into R. You will need to change the path in the `setwd` call to point to your local repo. Please use the given file name for the bible text file, to facilitate marking.

```
setwd("put/your/local/repo/location/here")
a <- scan("1581-0.txt", what="character", skip=156)
n <- length(a)
a <- a[-((n-2909):n)] ## strip license
```

Check the help file for any function whose purpose you are unclear of. The read in code gets rid of text you don't want at the start and end of the file. Check out what is in `a`.

4. Some pre-processing of the data file is needed. Write a function, called `split_punct`, which takes a vector of words as input along with a punctuation mark (e.g. " , " , " . " etc.). The function should search for each word containing the punctuation mark, remove it from the word, and add the mark as a new entry in the vector of words, after the word it came from. The updated vector should be what the function returns. For example, if looking for commas, then input vector

```
"An" "omnishambles," "in" "a" "headless" "chicken" "factory"
```

should become output vector

```
"An" "omnishambles" "," "in" "a" "headless" "chicken" "factory"
```

Functions `grep`, `rep` and `gsub` are the ones to use for this task. Beware that some punctuation marks are special characters in regular expressions, which `grep` and `gsub` can interpret. The notes tell you how to deal with this.

5. Use your `split_punct` function to separate the punctuation marks, " , " , " . " , " ; " , " ! " , " : " and " ? " from words they are attached to in the bible text.
6. `library(mgcv)` will load the recommended package `mgcv` in to R. `mgcv` provides a function `uniquecombs` that can be used to find the unique words in your text vector, and to return an index vector of the same length as the original text vector, indicating which entry in the unique word list each original word corresponds to. See `?uniquecombs`
 - (a) Use `uniquecombs` to find the vector of unique words. Do this having replaced the capital letters in words with lower case letters using function `tolower`.
 - (b) Making use of the index attribute returned by `uniquecombs`, and the `tabulate` function, count up how many time each unique word occurs in the text.
 - (c) You need to decide on a threshold number of occurrences at which a word should be included in the set of $m \approx 1000$ most common words. Write code to search for the threshold required to retain ≈ 1000 words.
 - (d) Hence create a vector, `b`, of the m most commonly occurring words.
7. Now you need to make the **A** matrix.
 - (a) The `match` function allows you to create a vector giving which element of your most common word vector, `b`, each element of the full text vector corresponds to. If a word is not in `b`, then an NA is given. e.g.


```
match(c("tum", "tee", "tum", "tee", "tumpy", "tum", "wibble", "wobble"), c("tum", "tee"))
[1]  1  2  1  2 NA  1 NA NA
```

Use `match` to get a vector indicating which common word each bible word is (don't forget to work on the lower case bible text).
 - (b) Now create a two column matrix (e.g. using `cbind`), in which the first column is the index of common words, and the next column is the index for the following word. i.e. the index vector created by `match` followed by that vector shifted by one place (you need to remove first and last entries as appropriate). Each row of your matrix is a pair of subsequent words in the text. When a row has no NAs then we have a pair of common words, which will contribute to our **A** matrix.
 - (c) Using `rowSums` and `is.na` identify the common word pairs, and drop the other word pairs (those that contain an NA).
 - (d) Now loop through the common word pairs adding a 1 to `A[i, j]` every time the j th common word follows the i th common word. Make sure `A` is initialized to the right thing before you start counting.
 - (e) Finally standardize the rows of `A`, so that each row's entries sum to 1, and we can interpret `A[i, j]` as the probability that `b[j]` will follow `b[i]`.
8. Finally write code to simulate 50-word sections from your model, starting from a randomly selected entry in `b`. Do this by simulating integers indexing words from the model, and then printing out the corresponding text with `cat`. The `sample` function will be useful.
9. If you get everything working and have time to go for the last 3 marks, then modify your code so that words that most often start with a capital letter in the main text, also start with a capital letter in your simulation. But do make sure that you achieve this in a way that does not mess up the word frequencies!

One piece of work - the text file containing your commented R code - is to be submitted for each group of 3 on Learn by 16:00 Friday 8th October 2021. You may be asked to supply us with an invitation to your github repo, so ensure this is in good order.

No extensions are available on this course, because of the frequency with which work has to be submitted. So late work will automatically attract a penalty. Technology failures will not be accepted as a reason for lateness (unless it is provably the case that Learn was unavailable for an extended period), so aim to submit ahead of time.

Marking Scheme: Full marks will be obtained for code that:

1. does what it is supposed to do, and has been coded in R approximately as indicated (that is marks will be lost for simply finding a package or online code that simplifies the task for you).
2. is carefully commented, so that someone reviewing the code can easily tell exactly what it is for, what it is doing and how it is doing it without having read this sheet, or knowing anything else about the code. Note that *easily tell* implies that the comments must also be as clear and *concise* as possible. You should assume that the reader knows basic R, but not that they know exactly what every function in R does.
3. is well structured and laid out, so that the code itself, and its underlying logic, are easy to follow.
4. is reasonably efficient. As a rough guide the whole code should take less than a minute to run - much longer than that and something is probably wrong.
5. includes the final part - but this is only worth 3 marks out of 18.
6. was prepared collaboratively using git and github in a group of 3.