

# R Notebook

First, we import our necessary libraries.

```
library(corpora)
library(HMM)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
```

The main library we care about is from corpora and called “BrownBigrams.” It is a data frame with 24167 rows and the following columns:

id: unique ID of the bigram entry

word1: the first word form in the bigram (character)

pos1: part-of-speech category of the first word (factor)

word2: the second word form in the bigram (character)

pos2: part-of-speech category of the second word (factor)

O11: co-occurrence frequency of the bigram (numeric)

O12: occurrences of the first word without the second (numeric)

O21: occurrences of the second word without the first (numeric)

O22: number of bigram tokens containing neither the first nor the second word (numeric)

It looks like this:

```
head(BrownBigrams,5)
```

```
##   id word1 pos1 word2 pos2 O11 O12   O21   O22
## 1  1     %    N     .    .    8  67 45465 864228
## 2  2     %    N   in    I    6  69 17981 891712
## 3  3     %    N   of    I   18  57 35481 874212
## 4  4     %    S     .    .    5  49 45468 864246
## 5  5     %    S   of    I   15  39 35484 874230
```

Our states are the possible parts of speech tags. We construct the list of states like this:

```
states = c("C", "D", "E", "F", "G", "I", "J", "L", "M", "N", "P", "R", "S", "T", "U", "V", "W", "Y", ".")
states

## [1] "C" "D" "E" "F" "G" "I" "J" "L" "M" "N" "P" "R" "S" "T" "U" "V" "W"
## [18] "Y" "."
```

Our Symbols are the words. We construct the list of symbols like this:

```
symbols = unique(BrownBigrams$word1)
```

Next, we make our transition matrix. This is done by finding the counts of all the bigrams of parts of speech and dividing by the occurrence of the first one.

```
transition_matrix = matrix(0L, nrow = length(states), ncol = length(states), dimnames = list(states, states))

for(i in 1:24167){
  transition_matrix[BrownBigrams$pos1[i], BrownBigrams$pos2[i]] = transition_matrix[BrownBigrams$pos1[i], BrownBigrams$pos2[i]] + 1
}

transition_matrix
```

```
##      C    D    E  F G    I    J    L  M    N    P    R    S    T U    V    W    Y .
## C   43   67   47   3 0   28  103 136 19  322  30 111   4  17 0  198 24  71 0
## D   10   42   15   1 3   54   37 961 14 2850   9  58   2   5 0   41   6 368 0
## E    0    0    0   0 0   0    0   0  9   0   0   1   0   0 0   9   0   0 0
## F    1    0    0   0 6   0    0   0  0   0   0   0   0   0 0   0   0   0 0
## G   88   15 158   0 0   58  531   7  0   21  28  43   0 174 0   2   2   0 0
## I   22 163 345 11 0 124  136 328   2  808 263  97 15   21 0   13 48 296 0
## J  108 135   16   0 0   0  176  28   0  628   4   7   0  92 0   5   1   3 0
## L    7    0    0   0 0   0    0   0  0   0   0   0   0   0 0   0   0   0 0
## M    4    3    9   0 0   0    1   0  0   0  21  61   0   1 0 225   0   0 0
## N 1299 587   30   1 0   13 2143 14 36  311  48  43   7 284 0 435 95   5 0
## P   30   21  46   0 0   4  125   7 79   2  14 133   0 13 0 655 15   0 0
## R  118   50 125   3 0   54  299 149   6   5  62 174   2  56 0 175 17   0 0
## S    2    2    0   0 0   0    8   0  2   10   0   0   0   1 0   4   2   0 0
## T    2   21  16   0 0   1    6  26   0  38  13   5   1   0 0 518   5  33 0
## U    4    0    0   0 0   0    0   0  0   0   1   0   0   0 1   0   0   0 0
## V  111   34 667   4 0 722  599 198   0  55 296 430   5 205 0   25 43   7 0
## W    5    4   33   3 0   0   11   6 39   12  53  17   2   5 0 133   0   4 0
## Y  130   92   2   0 0   0   74   1  5   47   1   0   0   4 0   72   1 419 0
## .    0    0    0   0 0   0    0   0  0   0   0   0   0   0 0   0   0   0 0
```

Now that we have our counts, we perform operations to get probabilities.

```
tag_count = BrownBigrams %>% group_by(pos1) %>% tally()

for(i in 1:18){
  for(j in 1:19){
    transition_matrix[states[i],states[j]] = transition_matrix[states[i],states[j]] / tag_count$n[i]
  }
}

transition_matrix
```

```
##      C      D      E      F      G
## C 0.035159444 0.054783320 0.038430090 0.0024529845 0.0000000000
## D 0.002234138 0.009383378 0.003351206 0.0002234138 0.0006702413
## E 0.000000000 0.000000000 0.000000000 0.0000000000 0.0000000000
## F 0.142857143 0.000000000 0.000000000 0.0000000000 0.8571428571
## G 0.078083407 0.013309672 0.140195209 0.0000000000 0.0000000000
## I 0.008172363 0.060549777 0.128157504 0.0040861813 0.0000000000
## J 0.089775561 0.112219451 0.013300083 0.0000000000 0.0000000000
## L 1.000000000 0.000000000 0.000000000 0.0000000000 0.0000000000
## M 0.012307692 0.009230769 0.027692308 0.0000000000 0.0000000000
## N 0.242758363 0.109699122 0.005606429 0.0001868810 0.0000000000
## P 0.026223776 0.018356643 0.040209790 0.0000000000 0.0000000000
```

```

## R 0.091119691 0.038610039 0.096525097 0.0023166023 0.0000000000
## S 0.064516129 0.064516129 0.000000000 0.000000000 0.000000000
## T 0.002919708 0.030656934 0.023357664 0.000000000 0.000000000
## U 0.666666667 0.000000000 0.000000000 0.000000000 0.000000000
## V 0.032637460 0.009997060 0.196118789 0.0011761247 0.000000000
## W 0.015290520 0.012232416 0.100917431 0.0091743119 0.000000000
## Y 0.153301887 0.108490566 0.002358491 0.000000000 0.000000000
## . 0.000000000 0.000000000 0.000000000 0.000000000 0.000000000
##           I           J           L           M           N           P
## C 0.022894522 0.084219133 0.111201962 0.0155355683 0.263286999 0.024529845
## D 0.012064343 0.008266309 0.214700626 0.0031277927 0.636729223 0.002010724
## E 0.000000000 0.000000000 0.000000000 0.4736842105 0.000000000 0.000000000
## F 0.000000000 0.000000000 0.000000000 0.000000000 0.000000000 0.000000000
## G 0.051464064 0.471162378 0.006211180 0.000000000 0.018633540 0.024844720
## I 0.046062407 0.050520059 0.121842496 0.0007429421 0.300148588 0.097696880
## J 0.000000000 0.146300914 0.023275145 0.000000000 0.522028263 0.003325021
## L 0.000000000 0.000000000 0.000000000 0.000000000 0.000000000 0.000000000
## M 0.000000000 0.003076923 0.000000000 0.000000000 0.000000000 0.064615385
## N 0.002429452 0.400485890 0.002616333 0.0067277144 0.058119978 0.008970286
## P 0.003496503 0.109265734 0.006118881 0.0690559441 0.001748252 0.012237762
## R 0.041698842 0.230888031 0.115057915 0.0046332046 0.003861004 0.047876448
## S 0.000000000 0.258064516 0.000000000 0.0645161290 0.322580645 0.000000000
## T 0.001459854 0.008759124 0.037956204 0.000000000 0.055474453 0.018978102
## U 0.000000000 0.000000000 0.000000000 0.000000000 0.000000000 0.166666667
## V 0.212290503 0.176124669 0.058218171 0.000000000 0.016171714 0.087033226
## W 0.000000000 0.033639144 0.018348624 0.1192660550 0.036697248 0.162079511
## Y 0.000000000 0.087264151 0.001179245 0.0058962264 0.055424528 0.001179245
## . 0.000000000 0.000000000 0.000000000 0.000000000 0.000000000 0.000000000
##           R           S           T           U           V           W
## C 0.090760425 0.0032706460 0.013900245 0.00000000 0.161896975 0.0196238757
## D 0.012957998 0.0004468275 0.001117069 0.00000000 0.009159964 0.0013404826
## E 0.052631579 0.0000000000 0.000000000 0.00000000 0.473684211 0.000000000
## F 0.000000000 0.0000000000 0.000000000 0.00000000 0.000000000 0.000000000
## G 0.038154392 0.0000000000 0.154392192 0.00000000 0.001774623 0.0017746229
## I 0.036032689 0.0055720654 0.007800892 0.00000000 0.004829123 0.0178306092
## J 0.005818786 0.0000000000 0.076475478 0.00000000 0.004156276 0.0008312552
## L 0.000000000 0.0000000000 0.000000000 0.00000000 0.000000000 0.000000000
## M 0.187692308 0.0000000000 0.003076923 0.00000000 0.692307692 0.000000000
## N 0.008035881 0.0013081667 0.053074192 0.00000000 0.081293216 0.0177536909
## P 0.116258741 0.0000000000 0.011363636 0.00000000 0.572552448 0.0131118881
## R 0.134362934 0.0015444015 0.043243243 0.00000000 0.135135135 0.0131274131
## S 0.000000000 0.0000000000 0.032258065 0.00000000 0.129032258 0.0645161290
## T 0.007299270 0.0014598540 0.000000000 0.00000000 0.756204380 0.0072992701
## U 0.000000000 0.0000000000 0.000000000 0.1666667 0.000000000 0.000000000
## V 0.126433402 0.0014701558 0.060276389 0.00000000 0.007350779 0.0126433402
## W 0.051987768 0.0061162080 0.015290520 0.00000000 0.406727829 0.000000000
## Y 0.000000000 0.0000000000 0.004716981 0.00000000 0.084905660 0.0011792453
## . 0.000000000 0.0000000000 0.000000000 0.00000000 0.000000000 0.000000000
##           Y .
## C 0.0580539657 0
## D 0.0822162645 0
## E 0.0000000000 0
## F 0.0000000000 0
## G 0.0000000000 0

```

```
## I 0.1099554235 0
## J 0.0024937656 0
## L 0.0000000000 0
## M 0.0000000000 0
## N 0.0009344048 0
## P 0.0000000000 0
## R 0.0000000000 0
## S 0.0000000000 0
## T 0.0481751825 0
## U 0.0000000000 0
## V 0.0020582182 0
## W 0.0122324159 0
## Y 0.4941037736 0
## . 0.0000000000 0
```

Finally, we make our emission probabilities

```
emissions = matrix(OL, nrow = length(states), ncol = length(symbols), dimnames = list(states, symbols))

for(i in 1:24167){
  emissions[BrownBigrams$pos1[i], BrownBigrams$word1[i]] = emissions[BrownBigrams$pos1[i], BrownBigrams$word1[i]]
}

symbol_count = BrownBigrams %>% group_by(word1) %>% tally()
for(a in 1:length(states)){
  for(b in 1:length(symbols)){
    emissions[states[a],symbols[b]] = emissions[states[a],symbols[b]] / symbol_count$n[b]
  }
}
```

(We will not print out this matrix because it is far too long).

Now, we can make our HMM and run the viterbi algorithm

```
model = initHMM(States=states, Symbols=symbols, startProbs=NULL, transProbs=transition_matrix, emissions=emissions)
viterbi(model, c("how", "are", "you"))

## [1] "W" "V" "P"

viterbi(model, c("i", "like", "this", "class"))

## [1] "P" "V" "D" "N"
```