# Scientific Programming: Assignment 1

University of Cambridge

Henrik Åhl

November 1, 2016

# 1  Introduction

This is an assignment report in connection to the *Scientific Programming with R* module in the Computational Biology course at the University of Cambridge, Michelmas term 2016. All related code is as of November 1, 2016 available on `https://github.com/supersubscript/compbio/tree/master/src/sp_assignments/assignment_1`, or available per request by contacting hpa22@cam.ac.uk. Likewise, the corresponding assignment can be found on `https://github.com/sje30/rpc2016/tree/master/assigns`.

# 2  Solutions

The assignment consists of three sections, related to various parts of scientific programming. The first section deals with geometric plotting, the second with developing an algorithm for solving cryptarithms, and the third with data manipulation and visualisation.
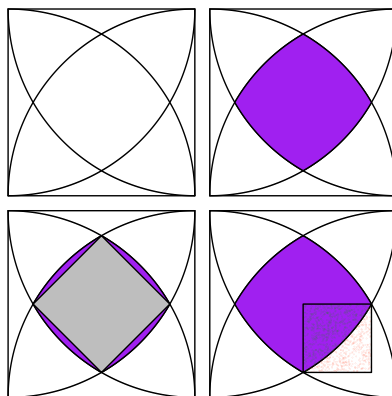
## 2.1  Curved Squares



Figure 1: The arcs and corresponding geometries, as well as the area estimate.

The analytic area can be derived by tracing the two upper quarter circles to their intersection and realising that a triangle based in the two bottom corners and this point must be equilateral. We can then use two sides of the triangle to form a circle sector, which must have area $\pi/6$. Likewise, we can use the triangle to extract a circle segment, which correspondingly will have area $\pi/6 - \sin(\pi/3)/2$. Knowing that a full quarter circle is here of area $pi/4$, we can easily deduce that the area surrounding the shape must be four times the area of the circle quartile minus the segment and sector. *Summa sumarum*, we find that the enclosed area is

$$1 - (\pi/4 - \pi/6 - (\pi/6 - \sin(\pi/3)/2)) \approx 0.315.$$

```
1  The area of the square contained is:        0.2678984
2  The estimated area of the shape is:  0.3150979
```

<div align="center">curved_squares_out.txt</div>

## 2.2 Cryptarithms

This exercise was solved using matrix manipulations, inserting character columns from a letter-value permutation matrix, parsing and evaluating. It has an approximate runtime on the order of minutes.

```
1  A B C D E F G H I
2  1 7 2 3 4 5 6 9 0
3  1 7 4 6 8 2 5 9 3
4  1 8 2 3 6 5 4 9 0
5  3 8 2 7 6 1 4 9 0
6  ###################
7  s e n d m o r y
8  9 5 6 7 1 0 8 2
9  ###################
10 s n o w r a i l e t
11 1 2 3 6 9 5 4 0 7 8
12 1 2 4 6 9 5 3 0 7 8
13 1 3 2 4 9 5 6 0 8 7
14 1 3 6 4 9 5 2 0 8 7
15 1 4 5 9 8 7 6 0 2 3
16 1 4 6 9 8 7 5 0 2 3
17 1 9 2 5 8 7 3 0 6 4
18 1 9 3 5 8 7 2 0 6 4
19 ###################
20 o n e t w h r l v
21 3 9 1 8 0 4 6 7 2
22 ###################
```

<div align="center">cryptarithms_out.txt</div>

## 2.3 Word Processing

In the word processing part, $y$ is treated as a vowel (as it should be). The dictionary does not distinguish between capitalised letters or not, so in particular in the palindrome exercise, the palindromes will be replicated if variants with different capitalisations exist.

```
Number of words:     99171
Number of ascii words:    244
Number of non-v words:    254
Number of words with apostrophes:  26141
The most frequent bigram is 's ' with a frequency of 0.03208276.
```

<div align="center">word_processing_out.txt</div>

```
B  B's  BB's   BC's   BLT's  BM's
BMW BMW's BS's   Bk   Bk's   Blvd
Br   Br's   C  C's  CD's   CFC's
CPR's  CRT's  CST's  CT's   Cd   Cd's
Cf   Cf's   Cl   Cl's   Cm   Cm's
Cr   Cr's   Cs   D  D's  DC's
DD's   DDS's  DP's   Dr   F  F's
FDR FDR's  FM's   Fm   Fm's   Fr
Fr's   G  G's  GHQ's  GMT's  GNP's
GP's   Gd   Gd's   H  H's  HF's
HP's   HQ's   HSBC  HSBC's   HTML's   Hf
Hf's   Hg   Hg's   Hz   Hz's   J
J's  JFK  JFK's  Jr   Jr's   K
K's  KFC  KFC's  KKK's  Kr   Kr's
L  L's  LBJ  LBJ's  LCD's  LPN's
LSD's  Ln   Lr   Lt   Ltd  M
M's  MB's   MD's   MGM MGM's  MHz
MP's   MS's   MSG's  MST's  MT's   MVP's
Mb   Md   Md's   Mg   Mg's   Mn
Mn's   Mr   Mr's   Mrs Ms   Mt
N  N's  NW's   Nb   Nb's   Nd
Nd's   Np   Np's   P  P's  PBS's
PC's   PM's   PMS's  PS's   PST's  PVC's
Pb   Pb's   Pd   Pd's   PhD  PhD's
Pl   Pm   Pm's   Pt   Pt's   Q
R  R's  RN's   RV's   Rb   Rb's
Rd   Rh   Rh's   Rn   Rn's   Rx
S  S's  SC's   SW's   Sb   Sb's
Sc   Sc's   Sgt  Sm   Sm's   Sn
Sn's   Sq   Sr   Sr's   St   T
T's  TB's   TLC's  TNT's  TV's   Tb
Tb's   Tc   Tc's   Th   Th's   Tl
Tl's   Tm   Tm's   V  V's  VCR's
VD's   VHF's  VLF's  W  W's  WWW's
Wm   Wm's   X  X's  XL's   Z
Z's  Zn   Zn's   Zr   Zr's   b
brr  c  cs   d  dB   f
g  gs   h  h'm  j  k
kHz kW   kc   ks   l  ls
m  ms   n  nth  p  pH
pj's   psst   q  r  rs   s
sh   t  ts   v  vs   w
x  z
```

<div align="center">word_processing_noVowelWords.txt</div>

<div align="center">4</div>

```
1  A Ada Ala Ana Anna   Ara
2  Ava B Bib Bob C D
3  E Eve F G Gog H
4  Hannah  I J K L Laval
5  M MGM Malayalam N Nan O
6  Ono Otto  P Q R S
7  S's Salas T Tet Tevet Tut
8  U V W X Y Z
9  a aha b bib bob boob
10 c civic d dad deed  deified
11 did dud e eke ere eve
12 ewe eye f g gag gig
13 h hah huh i j k
14 kayak kook  l level m ma'am
15 madam minim mom mum n non
16 noon  nun o oho p pap
17 peep  pep pip poop  pop pup
18 q r radar redder  refer rotor
19 s sagas sees  sexes shahs sis
20 solos sos stats t tat tenet
21 tit toot  tot u v w
22 wow x y z
```
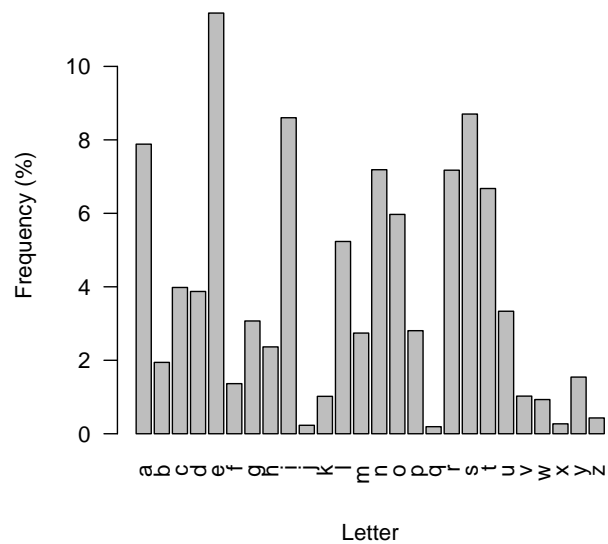
word_processing_palindromes.txt



Figure 2: Letter frequencies for the Ubuntu English (UK) dictionary.

Figure 3: Bigram for the letters contained in the dictionary. Square areas are proportional to the logarithm of the number of occurences. The horizontal row marks the first character, the vertical the subsequent. (You said frequency in the exercise, but I think log(occ) makes more sense for visualisation.)

# Code

```
1  sink("curved_squares_out.txt")
2
3  ### Take off axes
4  pdf("curved_squares.pdf", width = 3, height = 3)
5  par(bty = 'n',xaxt = 'n',yaxt = 'n')
6  par(mfrow = c(2,2), mar = c(0,0,0,0) + 0.0)
7
8  ### Define the functions, upper-to-lower, lower-to-upper
9  x = seq(0, 1, 10 ^ -4)
10 lowerUTL = function(x){-sqrt(1 - (x-1)^2) + 1}
11 upperUTL = function(x){sqrt(1 - x^2)}
12 lowerLTU = function(x){-sqrt(1 - x^2) + 1}
13 upperLTU = function(x){sqrt(1 - (x-1)^2)}
14
15 ####### FIGURE 1
16 ### Draw the arcs
17 plot(x, lowerUTL(x), type='l', ylab='', xlab='')
18 lines(x, lowerLTU(x))
19 lines(x, upperLTU(x))
20 lines(x, upperUTL(x))
21 rect(0,0,1,1)
22
23 ### Indices we need
24 intersectionIndex = which(abs(lowerUTL(x) - upperLTU(x)) < 10 ^ -4)
25 midpointIndex = length(x)/2
26
27 ####### FIGURE 2
28 ### Draw the arcs
29 plot(x, lowerUTL(x), type='l', ylab='', xlab='')
30 lines(x, lowerLTU(x))
31 lines(x, upperLTU(x))
32 lines(x, upperUTL(x))
33 rect(0,0,1,1)
34
35 ### Le polygon
36 polygon(c(x[seq(intersectionIndex, length(x) - intersectionIndex)],
37   rev(x[seq(intersectionIndex,length(x) - intersectionIndex)])),
38   c(upperLTU(x[seq(intersectionIndex, length(x) - midpointIndex)]),
39   upperUTL(x[seq(midpointIndex,length(x) - intersectionIndex)]),
40   rev(lowerLTU(x[seq(midpointIndex,length(x) - intersectionIndex)])),
41   rev(lowerUTL(x[seq(intersectionIndex, midpointIndex)])))), col='purple')
42
43 ####### FIGURE 3
44 ### Draw the arcs
45 plot(x, lowerUTL(x), type = 'l', ylab = '', xlab = '')
46 lines(x, lowerLTU(x))
47 lines(x, upperLTU(x))
48 lines(x, upperUTL(x))
49 rect(0,0,1,1)
```

```r
50
51 ### Le polygon
52 polygon(c(x[seq(intersectionIndex, length(x) − intersectionIndex)],
53    rev(x[seq(intersectionIndex,length(x) − intersectionIndex)])),
54    c(upperLTU(x[seq(intersectionIndex, length(x) − midpointIndex)]),
55    upperUTL(x[seq(midpointIndex,length(x) − intersectionIndex)]),
56    rev(lowerLTU(x[seq(midpointIndex,length(x) − intersectionIndex)])),
57    rev(lowerUTL(x[seq(intersectionIndex, midpointIndex)]))), col='purple')
58
59 ### Le cube
60 y1 =  x + upperLTU(x[intersectionIndex]) − x[intersectionIndex]
61 y2 = −x + upperLTU(x[intersectionIndex]) − x[intersectionIndex] + 1
62 y3 =  x + lowerLTU(.5) − .5
63 y4 = −x + lowerLTU(.5) + .5
64
65 area.square = (.5 − x[intersectionIndex])^2 +
66    (upperLTU(.5) − upperLTU(x[intersectionIndex]))^2
67 cat("The area of the square contained is:\t", area.square, "\n")
68
69 polygon(c(x[seq(intersectionIndex, length(x) − intersectionIndex)],
70    rev(x[seq(intersectionIndex,length(x) − intersectionIndex)])),
71    c(y1[seq(intersectionIndex, length(x) − midpointIndex)],
72    y2[seq(midpointIndex, length(x) − intersectionIndex)],
73    rev(y3[seq(midpointIndex,length(x) − intersectionIndex)]),
74    rev(y4[seq(intersectionIndex,midpointIndex)])), col='gray')
75
76 ###### FIGURE 4
77 ### Draw the arcs
78 plot(x, lowerUTL(x), type='l', ylab='', xlab='')
79 lines(x, lowerLTU(x))
80 lines(x, upperLTU(x))
81 lines(x, upperUTL(x))
82 rect(0,0,1,1)
83
84 ### Le polygon
85 polygon(c(x[seq(intersectionIndex, length(x) − intersectionIndex)],
86    rev(x[seq(intersectionIndex,length(x) − intersectionIndex)])),
87    c(upperLTU(x[seq(intersectionIndex, length(x) − midpointIndex)]),
88    upperUTL(x[seq(midpointIndex,length(x) − intersectionIndex)]),
89    rev(lowerLTU(x[seq(midpointIndex,length(x) − intersectionIndex)])),
90    rev(lowerUTL(x[seq(intersectionIndex, midpointIndex)]))), col='purple')
91
92 ### Naive Monte−Carlo (box around purple shape)
93 ### Note: Our problem is symmetric, so we can limit ourselves to a quartile.
94 nrTries = 1e7
95 xHit = runif(nrTries, .5, x[length(x) − intersectionIndex])
96 yHit = runif(nrTries, x[intersectionIndex], .5)
97 hits = which(yHit > lowerLTU(xHit))
98 misses = which(yHit <= lowerLTU(xHit))
99
100 ### Plot every 100th point.
```

```r
101  points ( xHit [ hits [ seq (1 , length ( hits ) ,10000) ] ] , yHit [ hits [ seq (1 , length ( hits )
          ,10000) ] ] ,
102    type = 'p', cex = .01 , col = 'forestgreen ')
103  points ( xHit [ misses [ seq (1 , length ( misses ) ,10000) ] ] , yHit [ misses [ seq (1 , length (
          misses ) , 10000) ] ] ,
104    type = 'p', cex = .01 , col = 'tomato ')
105  rect ( .5 , x [ intersectionIndex ] , x [ length ( x ) − intersectionIndex ] , .5 )
106
107  area . estimate = length ( hits ) / nrTries *
108    ( x [ length ( x ) − intersectionIndex ] − x [ intersectionIndex ] ) ^2
109  cat ("The estimated area of the shape is :\ t", area . estimate ,"\n")
110
111  ### Kill off graphical output
112  dev . off ()
```

../../../src/sp_assignments/assignment_1/curved_squares.R

```r
1   sink ("cryptarithms _out . txt")
2
3   #####################
4   ### FUNCTIONS
5   #####################
6
7   ### Returns possible permutations of input
8   permutations <− function (n) {
9     if (n == 1)
10    { return ( matrix (1) ) } else
11    {
12      sp <− permutations (n − 1)
13      p <− nrow ( sp )
14      A <− matrix ( nrow = n * p , ncol = n )
15      for (i in 1:n)
16      {
17        A[ ( i − 1) * p + 1:p ,] <− cbind (i , sp + ( sp >= i ) )
18      }
19      return (A)
20    }
21  }
22  getUnique = function (x)
23  {
24    if (x == 0) {
25      return (1)
26    }
27    x * getUnique (x − 1)
28  }
29
30  ### Which letters can 't be 0?
31  getFirstCharacters = function ( string )
32  {
33    string = gsub ("[[: punct :]]", " ", string ) # remove some symbols
34    string = gsub ("\\s+"," ", string ) # remove extra whitespaces
35    string = strsplit ( string ," ") [[1]]
```

9

```r
36    unique(sapply(string, function(x)
37      substring(x,1,1)))
38 }
39
40 ### Which are the unique letters?
41 uniqueLetters = function(inString)
42 {
43   inString = strsplit(inString, "")[[1]]
44   inString = inString[!inString %in% c(" ", "+", "*", "-", "=", "/", "&")]
45   inString = unique(inString)
46 }
47
48 #####################
49 ### MAIN
50 #####################
51
52 ### Takes in string and evaluates the corresponding cryptarithm.
53 ### Note: Several statements are separated by '&'.
54 ### Note: Only handles base 10 at the moment.
55 ### Note: Very inefficient; struggles with larger strings (=> larger
       matrices)
56 crypta = function(string.in) {
57
58   ### Get the number of row duplicates we will have
59   m = permutations(10) - 1
60   inString = string.in
61   letters = uniqueLetters(inString)
62   firstCharacters = getFirstCharacters(inString)
63   inString = strsplit(inString, "")[[1]]
64   colnames(m) = c(letters, rep("NA", 10 - length(letters)))
65   if (length(letters) != 10)
66   {
67     m = m[, -grep("NA", colnames(m))] # Take off edge.
68   }
69   m = m[seq(1,nrow(m),getUnique(10 - length(letters))),] # Remove duplicates
         .
70   for (ii in firstCharacters) {
71     m = m[-which(m[, ii] == 0),] # Kill off those w 0's at beginning.
72   }
73   inString = inString[!inString %in% c(" ")] # Take out essentials
74   strings = matrix(
75     inString, ncol = length(inString), nrow = nrow(m), byrow = TRUE
76   )
77   colnames(strings) = strings[1,]
78
79   ### Adjust our columns correspondingly.
80   count = 1
81   for (ii in colnames(strings)) {
82     if (ii %in% colnames(m)) {
83       strings[, count] = m[,ii]
84     }
```

```r
85     else if (ii == "=")
86     {
87       strings[, count] = rep("==", nrow(m))
88     }
89     count = count + 1
90   }
91
92   ### Add strings back together
93   strings = matrix(do.call(paste0, as.data.frame(strings)))
94   output = sapply(strings, function(x)
95     eval(parse(text = x)))
96   print(m[which(output == TRUE),])
97   cat("################################","\n")
98 }
99
100 ### Run program with our different setups.
101 crypta("AB * C = DE & DE + FG = HI")
102 crypta("send + more = money")
103 crypta("snow + rain = sleet")
104 crypta("one + two + two + three + three = eleven")
```

../../../src/sp_assignments/assignment_1/cryptarithms.R

```r
1 rm(list = ls())
2 sink("word_processing_out.txt")
3 ### 1 Number of words
4 words = read.table(
5   "../../../src/sp_assignments/assignment_1/usr-share-dict-words")
6 words = unlist(words)
7 #words = tolower(words)
8 #words = unique(words)
9 cat("Number of words:\t", length(words), "\n")
10
11 ### 2 Ascii
12 words.ascii = words[-grep("[^\x21-\x7e]", words)]
13 words.ascii = as.character(words.ascii)
14 cat("Number of ascii words:\t", length(words) - length(words.ascii),
15   "\n")
16
17 ### 3 Vowels
18 noVowelWords = words.ascii[-grep("[aeiouy]", words.ascii,
19   ignore.case = TRUE)]
20 cat("Number of non-v words:\t", length(noVowelWords), "\n")
21 write(noVowelWords, "word_processing_noVowelWords.txt", sep="\t", ncolumns =
22   6)
23 ### 4 Palindromes
24 reverseString = function(x) paste(substring(x, nchar(x):1,
25   nchar(x):1), collapse = "")
26 palindromes = words.ascii[which(sapply(words.ascii,
27   function(x) tolower(x) == tolower(reverseString(x))) == TRUE)]
```

11

```r
28  write(palindromes, "word_processing_palindromes.txt", sep="\t", ncolumns =
         6)
29
30  ### 5 Apostrophes
31  no.apostrophes = grep("[']", words.ascii, value = TRUE,
32    invert = TRUE)
33  cat("Number of words with apostrophes:\t", length(words.ascii)
34    - length(no.apostrophes), "\n")
35
36  ### 6 Probability distribution
37  letter.frequency = paste(no.apostrophes, collapse = "")
38  letter.frequency = summary(factor(tolower(strsplit(
39    letter.frequency, "")[[1]])))
40  letter.frequency = letter.frequency / sum(letter.frequency)
41  pdf("word_processing_letter_frequency.pdf", width = 5, height = 5)
42  par(mfrow = c(1, 1), las = 2)
43  barplot(letter.frequency * 100, xlab = 'Letter',
44    names.arg = letters[1:26], ylab="Frequency (%)")
45  x = dev.off()
46
47  ### 7 Bigram
48  bigram.matrix = matrix(0, length(letters) + 1, length(letters) + 1)
49  bigram.data   = tolower(paste(c("", no.apostrophes, ""),
50    collapse = " "))
51  colnames(bigram.matrix) = c(letters, " ")
52  rownames(bigram.matrix) = c(letters, " ")
53  pairs = substring(bigram.data, 1:(nchar(bigram.data) - 1),
54    2:nchar(bigram.data))
55  pairs.data = table(pairs)
56
57  cat("The most frequent bigram is \'",
58    names(pairs.data)[which(pairs.data == max(pairs.data))],
59    "\' with a frequency of ",
60    pairs.data[which(pairs.data == max(pairs.data))]
61    / sum(pairs.data), ".", sep = "")
62
63  for(ii in names(pairs.data))
64  {
65    chars = strsplit(ii,"")[[1]]
66    bigram.matrix[chars[1], chars[2]] =
67      bigram.matrix[chars[1], chars[2]]  + pairs.data[ii]
68  }
69
70  pdf("word_processing_bigram.pdf", width = 10, height = 10)
71  par(mfrow=c(1,1))
72  nrs = as.numeric(factor(rownames(bigram.matrix)))
73  plot(0, 0, type='n', xlim = c(1,27), ylim = c(1,27),
74      xaxt = 'n', yaxt = 'n', xlab = '', ylab = '')
75  for(ii in nrs)
76  {
77    for(jj in nrs)
```

```R
78    {
79        points ( ii , jj , cex = 0.22 * max(0 , log (bigram . matrix [ ii , jj ])) ,
80                 pch = 15 , col='forestgreen ')
81    }
82 }
83 axis (1 , at =1:27 , labels=c( letters [1:26] , '_'), tck= 0, cex . axis = .7)
84 axis (2 , at =1:27 , labels=c( letters [1:26] , '_'), tck= 0, cex . axis = .7)
85 x = dev . off ()
```

../../../src/sp_assignments/assignment_1/word_processing.R