# Appendix(A) - R script code for data cleanse up

November 17, 2016

```r
rm(list = ls())
getexpr = function(s,g)substring(s,g,g+attr(g,"match.length")-1)
thepage <- readLines('https://en.wikipedia.org/wiki/Historical_polling_for_U.S._
    Presidential_elections')

tbregexp <- '<table class="wikitable">'
tbBlock <- grep(tbregexp, thepage)

regExp <- c('<caption><b><a href=.* title=.*>([0-9]*)</a>', #year
            '<th>([^<]*)<?.* \\(.*\\).*</th>',  #candidate name
            '<td.*>[^>]*>*([A-Z][a-z]*).*</td>', #month name
            '<td.*>[^0-9]*([0-9]+)%.*</td>')

totalYr <- c()
totalMn <- c()
totalCan <- c()
totalRt <- c()
for(i in 1:length(tbBlock)) {
  result <- list()
  for (j in 1:3) {
    ofrom <- tbBlock[i]
    if (i == length(tbBlock))
      oto = length(thepage)-1
    else oto = tbBlock[i+1]
    datalines <- grep(regExp[j],thepage[ofrom:oto],value=TRUE)
    gg <- gregexpr(regExp[j],datalines)
    matches <- c()
    for(k in 1:length(gg)) {
      matches <- c(matches,getexpr(datalines[k],gg[[k]]))
    }

    result[[j]] <- gsub(regExp[j],'\\1',matches)
  }
  #datalines at this time should store line# of different month
  datalines <- grep(regExp[3],thepage[ofrom:oto])
  repMon <- c()
  wmatches <- c()
  cand <- c()

  for(j in 1:(length(datalines))) {
    from = ofrom+datalines[j]-1
    if (j == length(datalines)) to = oto
    else to = ofrom+datalines[j+1]-1

    rate <- grep(regExp[4], thepage[from:to],value=TRUE)
    gg <- gregexpr(regExp[4],rate)
```

```r
46
47    for(k in 1:length(gg)) {
48      matches <- getexpr(rate[k], gg[[k]])
49      wmatches <- c(wmatches,matches)
50      repMon <- c(repMon, rep(result[[3]][j], time=length(matches)))
51    }
52    wmatches <- gsub(regExp[4],'\\1',wmatches)
53  }
54
55  for(k in 1:(length(repMon)/length(result[[2]])))
56    cand <- c(cand, result[[2]])
57
58
59 }
60 totalYr <- c(totalYr, rep(result[[1]], length(wmatches)))
61 totalMn <- c(totalMn, repMon)
62 totalCan <- c(totalCan, cand)
63 totalRt <- c(totalRt, wmatches)
```

```r
1 rm(list = ls())
2
3 library(rvest)
4 library(stringr)
5
6 tb <- "https://www.archives.gov/federal-register/electoral-college/scores.html" %>%
7     read_html() %>%
8     html_nodes(xpath = '//tr/td/table') %>%
9     .[[1]]
10
11 hPath <- c('//tr[1]/th[1]',
12          '//tr[2]/th',
13          '//tr[3]/th',
14          '//tr[4]/th',
15          '//tr[4]/th',
16          '//tr[5]/th',
17          '//tr[5]/th',
18          '//tr[6]/th',
19          '//tr[7]/th',
20          '//tr[8]/th')
21 tPath <- c('//tr[1]/th[2]',
22          '//tr[2]/td',
23          '//tr[3]/td[1]',
24          '//tr[4]/td[1]',
25          '//tr[4]/td[2]',
26          '//tr[5]/td[1]',
27          '//tr[5]/td[2]',
28          '//tr[6]/td',
29          '//tr[7]/td',
30          '//tr[8]/td')
31 dtset <- data.frame(index=1:53)
32
33 #Election Column
34 head <- tb  %>%
35   html_nodes(xpath = hPath[1]) %>%
36   html_text(trim = TRUE)
37 text <- tb  %>%
38   html_nodes(xpath = tPath[1]) %>%
39   html_text(trim = TRUE)
40 colname <- head[1]
```

```r
41  dtset[colname] <- text
42
43  #President Column
44  head <- tb  %>%
45    html_nodes(xpath = hPath[2]) %>%
46    html_text(trim = TRUE)
47  text <- tb  %>%
48    html_nodes(xpath = tPath[2]) %>%
49    html_text(trim = TRUE)
50
51  remove <- c("") #Used to Eliminate the extra "" col
52  head <- head[!head %in% remove]
53  text <- text[!text %in% remove]
54  colname <- head[1]
55  dtset[colname] <- text
56
57  #Main Opponent Column
58  head <- tb  %>%
59    html_nodes(xpath = hPath[3]) %>%
60    html_text(trim = TRUE)
61  text <- tb  %>%
62    html_nodes(xpath = tPath[3]) %>%
63    html_text(trim = TRUE)
64
65  pat = '([^[:digit:]])*\\]$'
66  head <- head[!head %in% remove]
67  text <- text[grep(pat, text)]
68  colname <- head[1]
69  dtset[colname] <- text
70
71  #Winner Electoral Column
72  head <- tb  %>%
73    html_nodes(xpath = hPath[4]) %>%
74    html_text(trim = TRUE)
75  text <- tb  %>%
76    html_nodes(xpath = tPath[4]) %>%
77    html_text(trim = TRUE)
78
79  head <- head[!head %in% remove]
80  text <- text[!text %in% remove] #Used to Eliminate the extra "" col
81  colname <- paste(head[1],"-Winner")
82  dtset[colname] <- text
83
84  #Opponent Electoral Column
85  head <- tb  %>%
86    html_nodes(xpath = hPath[5]) %>%
87    html_text(trim = TRUE)
88  text <- tb  %>%
89    html_nodes(xpath = tPath[5]) %>%
90    html_text(trim = TRUE)
91
92  head <- head[!head %in% remove]
93  text <- text[!text %in% remove] #Used to Eliminate the extra "" col
94  colname <- paste(head[1],"-Opponent")
95  dtset[colname] <- text
96
97  #Popular Vote Winner Column
98  head <- tb  %>%
99    html_nodes(xpath = hPath[6]) %>%
```

```r
   html_text(trim = TRUE)
text <- tb   %>%
  html_nodes(xpath = tPath[6]) %>%
  html_text(trim = TRUE)

pat = '\r'
head <- head[!head %in% remove]
text <- text[!grepl(pat,text)]
colname <- paste(head[1],"-Winner")
dtset[colname] <- text

#Popular Vote Opponent Column
head <- tb   %>%
  html_nodes(xpath = hPath[7]) %>%
  html_text(trim = TRUE)
text <- tb   %>%
  html_nodes(xpath = tPath[7]) %>%
  html_text(trim = TRUE)

pat = 'Return to Index'
head <- head[!head %in% remove]
text <- c(rep("no record", time=9), text[!grepl(pat,text)])
colname <- paste(head[1],"-Opponent")
dtset[colname] <- text

#Vote for Others Column
head <- tb   %>%
  html_nodes(xpath = hPath[8]) %>%
  html_text(trim = TRUE)
text <- tb   %>%
  html_nodes(xpath = tPath[8]) %>%
  html_text(trim = TRUE)

pat = 'Votes for Others'

head <- head[!head %in% remove]
text <- text[!text %in% remove]
condition <- !grepl(pat,head)
copy <-text
text[condition] <- 'NA'

colname <- head[1]
dtset[colname] <- text

#Vice President Column
head <- tb   %>%
  html_nodes(xpath = hPath[9]) %>%
  html_text(trim = TRUE)
text <- tb   %>%
  html_nodes(xpath = tPath[9]) %>%
  html_text(trim = TRUE)

pat = 'Vice President'
rmpat = 'Notes|(Return to Index)'
head <- head[!head %in% remove]
text <- text[!text %in% remove]
text <- text[!grepl(rmpat,text)]
head[condition] <- pat
prev <- text
```

```r
159 text[condition] <- copy[condition]
160
161 colname <- head[1]
162 dtset[colname] <- text
163
164
165 #Further Cleanse up
166 dtset <- dtset[!names(dtset) %in% 'index']
167 partyPattern = '\\[([^[:digit:]])*\\]'
168 newCol1 = unlist(str_extract_all(dtset[2][,],partyPattern))
169 newCol2 = unlist(str_extract_all(dtset[3][,],partyPattern))
170
171 newCol1 <- gsub('(\\[)|(\\])','',newCol1)
172 newCol2 <- gsub('(\\[)|(\\])','',newCol2)
173
174 dtset[2][,] <-gsub(' \\[([^[:digit:]])*\\]','',dtset[2][,])
175 dtset[3][,] <-gsub(' \\[([^[:digit:]])*\\]','',dtset[3][,])
176
177 dtset$WinnerParty <- newCol1
178 dtset$OpponentParty <- newCol2
179
180 for (i in 4:7) {
181   dtset[i][,] <- gsub('[^[:digit:]]*', '', dtset[i][,])
182 }
183
184 vpevpat <- ' (\\([0-9]*\\))'
185 vpev <- unlist(str_extract_all(dtset[9][,],vpevpat)) %>%
186         gsub(pattern='(\\()|(\\))',replacement = '')
187 vpev <- c(rep('',time = 4),vpev)
188 dtset[9][,] <- gsub(' \\([0-9]*\\)', '', dtset[9][,])
189 dtset$VicePresidentEV <- vpev
190
191 #Candidate Pool Dataset
192 data.frame(dtset$Election ,dtset$`Votes for Others`)
193 ele <- dtset$Election
194 cpYr <- c()
195 cpLst <- c()
196 cpVote <- c()
197
198 for (i in 1:length(ele)) {
199   sigYrVote <- dtset$`Votes for Others`[i]
200   pat <- '[:alpha:]([:alpha:]| |\\.)*[:alpha:] \\([0-9]*\\)'
201   if (grepl('^NA$', sigYrVote)) {
202     cpYr <- c(cpYr, ele[i])
203     cpLst <- c(cpLst, '')
204     cpVote <- c(cpVote, '')
205   } else {
206     temp <- unlist(str_extract_all(sigYrVote, pat))
207     sigYrVote <- temp %>% gsub(pattern=vpevpat, replacement='')
208     vote <- temp %>% gsub(pattern='[^0-9]', replacement='')
209     cpYr <- c(cpYr, rep(ele[i], length(sigYrVote)))
210     cpLst <- c(cpLst, sigYrVote)
211     cpVote <- c(cpVote, vote)
212   }
213 }
214
215 dtset <- dtset[!names(dtset) %in% 'Votes for Others']
216 data.frame(Election=cpYr, 'Candidate List'=cpLst, 'Candidate Vote'=cpVote) %>% write.
     csv(row.names=FALSE, file='C:/Users/Yufan/Desktop/CMSC424/CandidatePool.csv')
```

```
217  dtset %>% write.csv(row.names=FALSE, file='C:/Users/Yufan/Desktop/CMSC424/
         ElectoralVoteDataSet.csv')
```

1