# Final_Project

Mikayla

3/26/2022

```r
#read in the original dataset
jpDf = read.csv('JEOPARDY_CSV.csv')

#cleaning the columns in jpDf

#cleaning category
library(quanteda)
```

```
## Warning: package 'quanteda' was built under R version 4.1.3
```

```
## Package version: 3.2.1
## Unicode version: 13.0
## ICU version: 69.1
```

```
## Parallel computing: 16 of 16 threads used.
```

```
## See https://quanteda.io for tutorials and examples.
```

```r
jpDf$charCategory = as.character(jpDf$Category)
categoryCorp = corpus(jpDf$charCategory)
categoryTokens = tokens(categoryCorp, remove_punct = TRUE, remove_symbols =
TRUE)
categoryLower = tokens_tolower(categoryTokens)
categoryCleaned = dfm(categoryLower, remove = stopwords())
```

```
## Warning: 'remove' is deprecated; use dfm_remove() instead
```
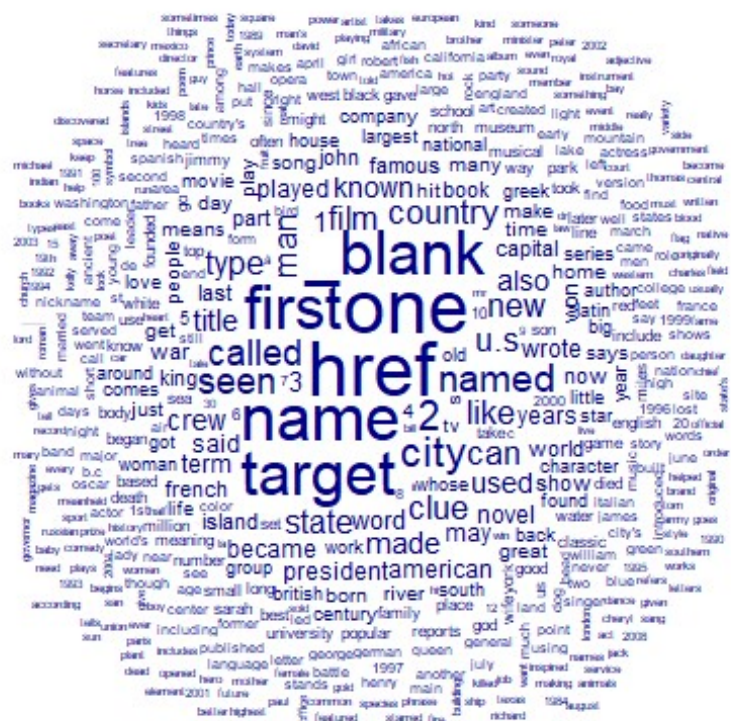
```r
#Cleaning Question
jpDf$charQuestion = as.character(jpDf$Question)
questionCorp = corpus(jpDf$charQuestion)
questionTokens = tokens(questionCorp, remove_punct = TRUE, remove_symbols =
TRUE,
                        remove_url = TRUE)
questionLower = tokens_tolower(questionTokens)
questionCleaned = dfm(questionLower, remove = stopwords())
```

```
## Warning: 'remove' is deprecated; use dfm_remove() instead
```

```r
library(quanteda.textplots)
```

```
## Warning: package 'quanteda.textplots' was built under R version 4.1.3
```

```r
#word cloud for overall top recurring words in Category
textplot_wordcloud(categoryCleaned, min_count =10)
```

```
textplot_wordcloud(questionCleaned, min_count =10)
```

```r
library(quanteda.textstats)

## Warning: package 'quanteda.textstats' was built under R version 4.1.3

#Frequency of Words in Category
fiftyFreqWordsCategory = textstat_frequency(categoryCleaned, n=50)

#Frequency of Words in Question
fiftyFreqWordsQuestion = textstat_frequency(questionCleaned, n=70)

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

#Grouping by: Round (Jeopardy!/Double Jeopardy!/Final Jeopardy!)
grpRoundJ <- jpDf[jpDf$Round == 'Jeopardy!',]
grpRoundDJ <- jpDf[jpDf$Round == 'Double Jeopardy!',]
grpRoundFJ <- jpDf[jpDf$Round == 'Final Jeopardy!',]

#Grouping by: Values -  Values (400/600/800/1000/1200/1600/2000/None)
# Where None is represented as final Jeopardy
grpValue400 <- jpDf[jpDf$Value == '$400',]
grpValue600 <- jpDf[jpDf$Value == '$600',]
grpValue800 <- jpDf[jpDf$Value == '$800',]
grpValue1000 <- jpDf[jpDf$Value == '$1000',]
grpValue1200 <- jpDf[jpDf$Value == '$1200',]
grpValue1600 <- jpDf[jpDf$Value == '$1600',]
grpValue2000 <- jpDf[jpDf$Value == '$2000',]
grpValueNone <- jpDf[jpDf$Value == 'None',]

#cleaning the grouped by Round dataframes
#Jeopardy! (Question)
grpRoundJ$charQuestion = as.character(grpRoundJ$Question)
questionCorp = corpus(grpRoundJ$charQuestion)
questionTokens = tokens(questionCorp, remove_punct = TRUE, remove_symbols =
TRUE,
                        remove_url = TRUE)
questionLower = tokens_tolower(questionTokens)
questionCleanedJ = dfm(questionLower, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead
```

```r
#Double Jeopardy!
grpRoundDJ$charQuestion = as.character(grpRoundDJ$Question)
questionCorpDJ = corpus(grpRoundDJ$charQuestion)
questionTokensDJ = tokens(questionCorpDJ, remove_punct = TRUE, remove_symbols
= TRUE,
                          remove_url = TRUE)
questionLowerDJ = tokens_tolower(questionTokensDJ)
questionCleanedDJ = dfm(questionLowerDJ, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead

#Final Jeopardy!
grpRoundFJ$charQuestion = as.character(grpRoundFJ$Question)
questionCorpFJ = corpus(grpRoundFJ$charQuestion)
questionTokensFJ = tokens(questionCorpFJ, remove_punct = TRUE, remove_symbols
= TRUE,
                          remove_url = TRUE)
questionLowerFJ = tokens_tolower(questionTokensFJ)
questionCleanedFJ = dfm(questionLowerFJ, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead

#Jeopardy! (Category)
grpRoundJ$charCategory = as.character(grpRoundJ$Category)
categoryCorpJ = corpus(grpRoundJ$charCategory)
categoryTokensJ = tokens(categoryCorpJ, remove_punct = TRUE, remove_symbols =
TRUE)
categoryLowerJ = tokens_tolower(categoryTokensJ)
categoryCleanedJ = dfm(categoryLowerJ, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead

#Double Jeopardy! (Category)
grpRoundDJ$charCategory = as.character(grpRoundDJ$Category)
categoryCorpDJ = corpus(grpRoundDJ$charCategory)
categoryTokensDJ = tokens(categoryCorpDJ, remove_punct = TRUE, remove_symbols
= TRUE)
categoryLowerDJ = tokens_tolower(categoryTokensDJ)
categoryCleanedDJ = dfm(categoryLowerDJ, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead

#Double Jeopardy! (Category)
grpRoundFJ$charCategory = as.character(grpRoundFJ$Category)
categoryCorpFJ = corpus(grpRoundFJ$charCategory)
categoryTokensFJ = tokens(categoryCorpFJ, remove_punct = TRUE, remove_symbols
= TRUE)
categoryLowerFJ = tokens_tolower(categoryTokensFJ)
categoryCleanedFJ = dfm(categoryLowerFJ, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead
```

```r
#cleaning the grouped by Value dataframes
#400 (Question)
grpValue400$charQuestion = as.character(grpValue400$Question)
questionCorp400 = corpus(grpValue400$charQuestion)
questionTokens400 = tokens(questionCorp400, remove_punct = TRUE,
remove_symbols = TRUE,
                          remove_url = TRUE)
questionLower400 = tokens_tolower(questionTokens400)
questionCleaned400 = dfm(questionLower400, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead

#600
grpValue600$charQuestion = as.character(grpValue600$Question)
questionCorp600 = corpus(grpValue600$charQuestion)
questionTokens600 = tokens(questionCorp600, remove_punct = TRUE,
remove_symbols = TRUE,
                          remove_url = TRUE)
questionLower600 = tokens_tolower(questionTokens600)
questionCleaned600 = dfm(questionLower600, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead

#800
grpValue800$charQuestion = as.character(grpValue800$Question)
questionCorp800 = corpus(grpValue800$charQuestion)
questionTokens800 = tokens(questionCorp800, remove_punct = TRUE,
remove_symbols = TRUE,
                          remove_url = TRUE)
questionLower800 = tokens_tolower(questionTokens800)
questionCleaned800 = dfm(questionLower800, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead

#1000
grpValue1000$charQuestion = as.character(grpValue1000$Question)
questionCorp1000 = corpus(grpValue1000$charQuestion)
questionTokens1000 = tokens(questionCorp1000, remove_punct = TRUE,
remove_symbols = TRUE,
                          remove_url = TRUE)
questionLower1000 = tokens_tolower(questionTokens1000)
questionCleaned1000 = dfm(questionLower1000, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead

#1200
grpValue1200$charQuestion = as.character(grpValue1200$Question)
questionCorp1200 = corpus(grpValue1200$charQuestion)
questionTokens1200 = tokens(questionCorp1200, remove_punct = TRUE,
remove_symbols = TRUE,
                          remove_url = TRUE)
```

```
questionLower1200 = tokens_tolower(questionTokens1200)
questionCleaned1200 = dfm(questionLower1200, re200move = stopwords())

## Warning: re200move argument is not used.

#1600
grpValue1600$charQuestion = as.character(grpValue1600$Question)
questionCorp1600 = corpus(grpValue1600$charQuestion)
questionTokens1600 = tokens(questionCorp1600, remove_punct = TRUE,
remove_symbols = TRUE,
                          remove_url = TRUE)
questionLower1600 = tokens_tolower(questionTokens1600)
questionCleaned1600 = dfm(questionLower1600, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead

#2000
grpValue2000$charQuestion = as.character(grpValue2000$Question)
questionCorp2000 = corpus(grpValue2000$charQuestion)
questionTokens2000 = tokens(questionCorp2000, remove_punct = TRUE,
remove_symbols = TRUE,
                          remove_url = TRUE)
questionLower2000 = tokens_tolower(questionTokens2000)
questionCleaned2000 = dfm(questionLower2000, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead

#None
grpValueNone$charQuestion = as.character(grpValueNone$Question)
questionCorpNone = corpus(grpValueNone$charQuestion)
questionTokensNone = tokens(questionCorpNone, remove_punct = TRUE,
remove_symbols = TRUE,
                          remove_url = TRUE)
questionLowerNone = tokens_tolower(questionTokensNone)
questionCleanedNone = dfm(questionLowerNone, remove = stopwords())

## Warning: 'remove' is deprecated; use dfm_remove() instead

#freq for grouped by Round
twentyFreqWordsQuestionJ = textstat_frequency(questionCleanedJ, n=50)
twentyFreqWordsQuestionJ

##      feature frequency rank docfreq group
## 1       one      6684    1    6406   all
## 2      href      5976    2    4962   all
## 3      name      5419    3    5277   all
## 4     first      4877    4    4744   all
## 5    target      4798    5    3894   all
## 6    _blank      4762    6    3860   all
## 7      city      3181    7    3086   all
## 8         2      3082    8    2952   all
## 9       can      2731    9    2641   all
```

```
## 10      named     2676   10   2644   all
## 11       seen     2629   11   2620   all
## 12     called     2589   12   2560   all
## 13       like     2568   13   2476   all
## 14      state     2535   14   2409   all
## 15        u.s     2523   15   2496   all
## 16        new     2421   16   2341   all
## 17       type     2399   17   2343   all
## 18        man     2280   18   2232   all
## 19    country     2226   19   2209   all
## 20       clue     2131   20   1951   all
## 21       made     2118   21   2086   all
## 22       film     2048   22   2018   all
## 23       used     1969   23   1942   all
## 24       also     1810   24   1808   all
## 25      known     1776   25   1760   all
## 26          1     1725   26   1622   all
## 27      title     1716   27   1708   all
## 28        may     1692   28   1668   all
## 29       crew     1654   29   1651   all
## 30          3     1599   30   1546   all
## 31       said     1505   31   1493   all
## 32       word     1473   32   1412   all
## 33       term     1460   33   1456   all
## 34     became     1421   34   1406   all
## 35      years     1416   35   1396   all
## 36     played     1409   36   1380   all
## 37  president     1372   37   1329   all
## 38      world     1336   38   1320   all
## 39   american     1330   39   1311   all
## 40       show     1287   40   1256   all
## 41       part     1278   41   1249   all
## 42    capital     1245   42   1209   all
## 43       last     1226   43   1206   all
## 44       song     1225   44   1211   all
## 45        hit     1222   45   1194   all
## 46       home     1212   46   1183   all
## 47        now     1183   47   1175   all
## 48        won     1170   48   1133   all
## 49       make     1163   49   1143   all
## 50        get     1159   50   1130   all

twentyFreqWordsQuestionDJ = textstat_frequency(questionCleanedDJ, n=50)
twentyFreqWordsQuestionDJ

##       feature frequency rank docfreq group
## 1        href      6954    1    5526   all
## 2         one      5958    2    5696   all
## 3      target      5882    3    4570   all
## 4      _blank      5852    4    4541   all
```

```
## 5          name      5581      5      5468      all
## 6         first      4980      6      4862      all
## 7          city      3287      7      3190      all
## 8             2      2906      8      2798      all
## 9          clue      2848      9      2609      all
## 10        named      2833     10      2790      all
## 11       called      2772     11      2737      all
## 12         film      2597     12      2550      all
## 13         seen      2593     13      2590      all
## 14         like      2519     14      2442      all
## 15      country      2519     14      2491      all
## 16          new      2421     16      2348      all
## 17          man      2418     17      2354      all
## 18         crew      2328     18      2321      all
## 19         type      2270     19      2229      all
## 20        title      2269     20      2246      all
## 21          u.s      2241     21      2213      all
## 22          can      2178     22      2131      all
## 23        state      2060     23      1956      all
## 24        known      2013     24      1984      all
## 25         used      1973     25      1943      all
## 26         made      1923     26      1894      all
## 27         also      1877     27      1876      all
## 28        novel      1857     28      1850      all
## 29        wrote      1838     29      1804      all
## 30         word      1744     30      1671      all
## 31       became      1652     31      1635      all
## 32       played      1567     32      1514      all
## 33         king      1548     33      1484      all
## 34         said      1548     33      1532      all
## 35        years      1538     35      1520      all
## 36            3      1523     36      1490      all
## 37    president      1502     37      1455      all
## 38      capital      1467     38      1415      all
## 39            1      1455     39      1347      all
## 40     american      1443     40      1423      all
## 41         last      1436     41      1414      all
## 42          may      1410     42      1394      all
## 43         term      1404     43      1396      all
## 44          war      1401     44      1331      all
## 45       french      1400     45      1380      all
## 46       island      1391     46      1334      all
## 47         part      1375     47      1339      all
## 48         book      1342     48      1316      all
## 49         play      1302     49      1278      all
## 50        world      1301     50      1282      all
```

```r
twentyFreqWordsQuestionFJ = textstat_frequency(questionCleanedFJ, n=50)
twentyFreqWordsQuestionFJ
```

```
##         feature frequency rank docfreq group
## 1            2       360    1     344   all
## 2        first       351    2     338   all
## 3         name       311    3     295   all
## 4          one       308    4     285   all
## 5          u.s       228    5     226   all
## 6          man       207    6     200   all
## 7        named       169    7     164   all
## 8        state       155    8     141   all
## 9    president       153    9     139   all
## 10           1       144   10     133   all
## 11        last       137   11     136   all
## 12     country       132   12     123   all
## 13        city       131   13     125   all
## 14           3       124   14     120   all
## 15       years       119   15     114   all
## 16        said       117   16     117   all
## 17       title       111   17     110   all
## 18        film       111   17     109   all
## 19         new       104   19     101   all
## 20      called       101   20      99   all
## 21        word       100   21      95   all
## 22       wrote        96   22      96   all
## 23       novel        93   23      92   all
## 24      became        92   24      92   all
## 25       whose        92   24      91   all
## 26       world        87   26      85   all
## 27     century        86   27      85   all
## 28   character        83   28      82   all
## 29        made        82   29      81   all
## 30        book        76   30      74   all
## 31      famous        73   31      72   all
## 32        used        73   31      70   all
## 33       known        71   33      71   all
## 34        time        70   34      69   all
## 35        born        69   35      66   all
## 36           4        69   35      68   all
## 37        work        69   35      69   all
## 38    american        68   38      67   all
## 39      people        67   39      66   all
## 40       names        67   39      63   all
## 41        best        66   41      61   all
## 42         won        65   42      62   all
## 43         war        64   43      60   all
## 44        died        62   44      61   all
## 45         may        62   44      59   all
## 46         men        62   44      62   all
## 47       since        61   47      60   all
## 48     capital        61   47      56   all
```

```
## 49      states       61   47      60   all
## 50       year        61   47      57   all
```

```
twentyFreqWordsCategoryJ = textstat_frequency(categoryCleanedJ, n=50)
twentyFreqWordsCategoryJ
```

```
##           feature frequency rank docfreq group
## 1           words      2093    1    2068   all
## 2              tv      1735    2    1695   all
## 3           world      1681    3    1681   all
## 4         history      1566    4    1566   all
## 5             u.s      1344    5    1344   all
## 6           state      1266    6    1241   all
## 7            time      1098    7    1098   all
## 8          sports      1074    8    1074   all
## 9            food      1030    9    1005   all
## 10              s      1009   10     956   all
## 11          names      1003   11    1003   all
## 12       american       936   12     936   all
## 13            new       839   13     824   all
## 14        science       834   14     834   all
## 15           name       822   15     822   all
## 16          music       807   16     807   all
## 17          movie       799   17     799   all
## 18        century       792   18     792   all
## 19          first       740   19     725   all
## 20         movies       731   20     726   all
## 21          songs       687   21     682   all
## 22          rhyme       678   22     678   all
## 23         people       656   23     651   all
## 24         cities       607   24     607   all
## 25      geography       599   25     599   all
## 26           rock       586   26     586   all
## 27       business       538   27     523   all
## 28       crossword       537   28     537   all
## 29          clues       537   28     537   all
## 30         famous       514   30     514   all
## 31            old       511   31     511   all
## 32   presidential       495   32     495   all
## 33            pop       489   33     489   all
## 34          women       488   34     488   all
## 35         quotes       479   35     479   all
## 36        classic       473   36     473   all
## 37        america       469   37     469   all
## 38        animals       465   38     465   all
## 39            war       454   39     454   all
## 40           word       452   40     452   all
## 41           ends       449   41     449   all
## 42       capitals       448   42     448   all
## 43        literary       440   43     440   all
```

```
## 44        category         432    44       432    all
## 45          around         430    45       430    all
## 46         college         427    46       427    all
## 47         phrases         425    47       425    all
## 48            york         425    47       410    all
## 49             get         424    49       424    all
## 50      literature         423    50       423    all
```

```r
twentyFreqWordsCategoryDJ = textstat_frequency(categoryCleanedDJ, n=50)
twentyFreqWordsCategoryDJ
```

```
##            feature frequency rank docfreq group
## 1            world      2290    1    2290    all
## 2            words      2154    2    2121    all
## 3          history      2029    3    2029    all
## 4          century      1332    4    1332    all
## 5          science      1325    5    1325    all
## 6         american      1306    6    1306    all
## 7               tv      1257    7    1252    all
## 8                s      1155    8    1075    all
## 9            names      1103    9    1103    all
## 10           movie      1066   10    1066    all
## 11             u.s      1053   11    1053    all
## 12      literature      1051   12    1051    all
## 13             art      1043   13    1043    all
## 14           music      1022   14    1022    all
## 15          movies      1003   15     993    all
## 16         authors       950   16     925    all
## 17         literary       935   17     930    all
## 18            time       928   18     923    all
## 19           women       920   19     915    all
## 20          sports       807   20     807    all
## 21       geography       787   21     787    all
## 22          famous       738   22     738    all
## 23            name       691   23     691    all
## 24          cities       666   24     666    all
## 25             war       663   25     663    all
## 26          people       661   26     656    all
## 27      characters       654   27     654    all
## 28           state       638   28     623    all
## 29           clues       627   29     627    all
## 30           books       609   30     609    all
## 31             new       603   31     598    all
## 32            film       594   32     594    all
## 33         historic       594   32     594    all
## 34        crossword       574   34     574    all
## 35         artists       562   35     562    all
## 36             lit       561   36     561    all
## 37           opera       561   36     561    all
## 38           films       545   38     540    all
```

```
## 39      america       540    39      540    all
## 40        first       539    40      515    all
## 41      musical       492    41      492    all
## 42         word       475    42      475    all
## 43        rhyme       475    42      475    all
## 44         20th       468    44      468    all
## 45          old       465    45      446    all
## 46  shakespeare       462    46      462    all
## 47        title       450    47      445    all
## 48       quotes       446    48      446    all
## 49         book       444    49      444    all
## 50     capitals       442    50      442    all
```

```
twentyFreqWordsCategoryFJ = textstat_frequency(categoryCleanedFJ, n=50)
twentyFreqWordsCategoryFJ
```

```
##            feature frequency rank docfreq group
## 1              u.s       203    1     203    all
## 2          century       194    2     194    all
## 3          history       178    3     178    all
## 4            world       160    4     160    all
## 5           famous       131    5     131    all
## 6         american       124    6     124    all
## 7            names       117    7     117    all
## 8         historic       111    8     111    all
## 9             20th        94    9      94    all
## 10       presidents        94    9      93    all
## 11          authors        92   11      91    all
## 12       literature        82   12      82    all
## 13             19th        67   13      67    all
## 14       characters        65   14      65    all
## 15            women        63   15      63    all
## 16         business        62   16      62    all
## 17           sports        59   17      59    all
## 18        geography        58   18      58    all
## 19         capitals        58   18      58    all
## 20        americans        58   18      58    all
## 21          origins        56   21      56    all
## 22            state        55   22      55    all
## 23             word        49   23      49    all
## 24           cities        49   23      49    all
## 25           movies        45   25      45    all
## 26          literary        44   26      44    all
## 27          british        43   27      43    all
## 28     presidential        40   28      40    all
## 29          leaders        40   28      40    all
## 30            music        37   30      37    all
## 31            words        37   30      37    all
## 32           people        35   32      35    all
## 33          science        35   32      35    all
```

```
## 34          novels         33    34        33    all
## 35           news          33    34        33    all
## 36          states         32    36        32    all
## 37          movie          31    37        31    all
## 38          films          31    37        31    all
## 39          books          29    39        28    all
## 40       government        29    39        29    all
## 41       television        28    41        28    all
## 42        landmarks        28    41        28    all
## 43         royalty         27    43        27    all
## 44         ancient         27    43        27    all
## 45        countries        26    45        26    all
## 46         industry        26    45        26    all
## 47         phrases         25    47        25    all
## 48          awards         25    47        25    all
## 49      organizations      25    47        25    all
## 50         european        24    50        24    all
```

*#freq for grouped by Round*
```
twentyFreqWordsQuestion400 = textstat_frequency(questionCleaned400, n=50)
twentyFreqWordsQuestion400
```

```
##         feature frequency rank docfreq group
## 1          one      2651     1    2534    all
## 2          href     2196     2    1798    all
## 3         first     1997     3    1950    all
## 4          name     1912     4    1872    all
## 5        target     1809     5    1447    all
## 6        _blank     1798     6    1437    all
## 7          city     1268     7    1236    all
## 8             2     1165     8    1122    all
## 9        called     1014     9    1003    all
## 10          can      994    10     970    all
## 11        named      992    11     977    all
## 12      country      969    12     964    all
## 13        state      958    13     911    all
## 14         seen      955    14     948    all
## 15          u.s      948    15     934    all
## 16         like      947    16     927    all
## 17          new      914    17     884    all
## 18         type      912    18     894    all
## 19          man      861    19     838    all
## 20         clue      857    20     786    all
## 21         film      822    21     805    all
## 22         made      804    22     794    all
## 23         used      757    23     747    all
## 24        title      712    24     706    all
## 25         also      670    25     669    all
## 26         said      668    26     665    all
## 27         crew      668    26     666    all
```

```
## 28       known       667   28     659    all
## 29         may       629   29     621    all
## 30           1       625   30     583    all
## 31       years       611   31     606    all
## 32           3       607   32     593    all
## 33      became       601   33     599    all
## 34        word       576   34     553    all
## 35      played       575   35     557    all
## 36       novel       548   36     545    all
## 37       wrote       544   37     532    all
## 38        term       535   38     534    all
## 39   president       534   39     522    all
## 40    american       513   40     509    all
## 41        show       506   41     491    all
## 42     capital       506   41     490    all
## 43       world       489   43     477    all
## 44        part       487   44     481    all
## 45         war       478   45     466    all
## 46        book       475   46     465    all
## 47        home       474   47     467    all
## 48        last       472   48     466    all
## 49        king       469   49     452    all
## 50        john       468   50     456    all
```

```
twentyFreqWordsQuestion600 = textstat_frequency(questionCleaned600, n=51)
twentyFreqWordsQuestion600
```

```
##        feature frequency rank docfreq group
## 1         href      1253    1    1042    all
## 2          one      1238    2    1189    all
## 3         name      1053    3    1031    all
## 4       target      1032    4     839    all
## 5       _blank      1027    5     834    all
## 6        first       977    6     948    all
## 7         seen       555    7     552    all
## 8         city       547    8     531    all
## 9            2       537    9     515    all
## 10      called       497   10     491    all
## 11         can       479   11     463    all
## 12        film       479   11     474    all
## 13        like       478   13     459    all
## 14         new       472   14     462    all
## 15         man       450   15     439    all
## 16         u.s       448   16     444    all
## 17     country       443   17     439    all
## 18        clue       443   17     404    all
## 19       named       441   19     433    all
## 20       state       423   20     400    all
## 21        type       414   21     406    all
## 22       title       381   22     379    all
```

```
## 23        made        360   23       356    all
## 24        used        351   24       345    all
## 25        crew        346   25       346    all
## 26       known        345   26       342    all
## 27        also        340   27       340    all
## 28           3        299   28       290    all
## 29        term        298   29       298    all
## 30        said        291   30       288    all
## 31         may        291   30       287    all
## 32        word        284   32       273    all
## 33       novel        284   32       283    all
## 34       years        276   34       270    all
## 35   president        275   35       263    all
## 36      became        271   36       267    all
## 37           1        268   37       255    all
## 38      played        263   38       254    all
## 39    american        261   39       258    all
## 40       world        260   40       255    all
## 41        john        258   41       252    all
## 42        king        256   42       250    all
## 43      famous        256   42       255    all
## 44       wrote        250   44       246    all
## 45        last        247   45       244    all
## 46        show        244   46       241    all
## 47        part        238   47       234    all
## 48      island        237   48       230    all
## 49        book        236   49       229    all
## 50         won        235   50       226    all
## 51        song        234   51       230    all
```

```
twentyFreqWordsQuestion800 = textstat_frequency(questionCleaned800, n=58)
twentyFreqWordsQuestion800
```

```
##        feature frequency rank docfreq group
## 1         href      2376    1    1909   all
## 2       target      1969    2    1550   all
## 3       _blank      1960    3    1541   all
## 4          one      1944    4    1865   all
## 5         name      1636    5    1603   all
## 6        first      1412    6    1384   all
## 7         city       914    7     895   all
## 8            2       912    8     880   all
## 9         clue       880    9     807   all
## 10       named       869   10     859   all
## 11        seen       823   11     822   all
## 12        like       814   12     785   all
## 13      called       781   13     769   all
## 14         can       744   14     724   all
## 15         man       732   15     714   all
## 16        type       732   15     716   all
```

```
## 17        new      719    17    695    all
## 18       crew      717    18    715    all
## 19    country      711    19    705    all
## 20       film      699    20    691    all
## 21        u.s      691    21    681    all
## 22      title      655    22    651    all
## 23      state      637    23    598    all
## 24       made      619    24    606    all
## 25       also      616    25    615    all
## 26       used      611    26    601    all
## 27      known      586    27    581    all
## 28       word      505    28    479    all
## 29      wrote      480    29    475    all
## 30          1      456    30    428    all
## 31     became      456    30    450    all
## 32          3      450    32    437    all
## 33      novel      446    33    445    all
## 34     played      445    34    436    all
## 35        may      441    35    433    all
## 36       part      437    36    423    all
## 37       said      435    37    428    all
## 38       king      435    37    416    all
## 39       term      431    39    429    all
## 40   american      429    40    419    all
## 41      years      417    41    414    all
## 42  president      414    42    400    all
## 43       book      413    43    401    all
## 44      world      408    44    404    all
## 45       last      402    45    396    all
## 46        now      389    46    381    all
## 47     island      387    47    370    all
## 48       show      383    48    376    all
## 49     french      381    49    377    all
## 50        won      378    50    370    all
## 51    capital      371    51    359    all
## 52        get      361    52    356    all
## 53        war      358    53    339    all
## 54       just      343    54    342    all
## 55       john      339    55    332    all
## 56       play      335    56    329    all
## 57       time      334    57    322    all
## 58     famous      333    58    332    all

twentyFreqWordsQuestion1000 = textstat_frequency(questionCleaned1000, n=60)
twentyFreqWordsQuestion1000

##       feature frequency rank docfreq group
## 1        href      1363    1    1111    all
## 2        name      1230    2    1204    all
## 3      target      1131    3     900    all
```

```
## 4       _blank    1122    4    893   all
## 5         one     1002    5    955   all
## 6        first     826    6    810   all
## 7         city     618    7    597   all
## 8        named     579    8    567   all
## 9           2      541    9    511   all
## 10        seen     533   10    533   all
## 11      called     529   11    523   all
## 12     country     481   12    477   all
## 13        like     476   13    465   all
## 14        clue     472   14    433   all
## 15        type     459   15    447   all
## 16        film     450   16    443   all
## 17         u.s     445   17    438   all
## 18         new     437   18    426   all
## 19         man     426   19    420   all
## 20       state     393   20    371   all
## 21        used     390   21    383   all
## 22        crew     385   22    385   all
## 23         can     378   23    364   all
## 24       known     376   24    375   all
## 25       title     371   25    370   all
## 26        made     366   26    362   all
## 27        also     325   27    325   all
## 28    american     319   28    317   all
## 29        word     314   29    304   all
## 30       novel     306   30    304   all
## 31       wrote     303   31    295   all
## 32      french     297   32    294   all
## 33        term     293   33    292   all
## 34      became     287   34    283   all
## 35           3     279   35    268   all
## 36        last     278   36    274   all
## 37        king     276   37    271   all
## 38       years     275   38    270   all
## 39       world     275   38    273   all
## 40           1     270   40    253   all
## 41   president     265   41    253   all
## 42      played     256   42    248   all
## 43         may     256   42    253   all
## 44     capital     253   44    246   all
## 45       means     250   45    240   all
## 46      famous     248   46    248   all
## 47        said     245   47    243   all
## 48       whose     242   48    240   all
## 49         war     241   49    230   all
## 50         won     241   49    231   all
## 51        part     240   51    235   all
## 52         now     236   52    234   all
## 53       great     236   52    227   all
```

```
## 54       book        235    54      232    all
## 55       play        235    54      231    all
## 56     author        228    56      225    all
## 57       john        227    57      223    all
## 58      south        225    58      222    all
## 59     island        223    59      214    all
## 60       home        219    60      216    all
```

```
twentyFreqWordsQuestion1200 = textstat_frequency(questionCleaned1200, n=60)
twentyFreqWordsQuestion1200
```

```
##      feature frequency rank docfreq group
## 1        the     10158    1    6473   all
## 2          a      8029    2    4611   all
## 3       this      6968    3    6897   all
## 4         of      6567    4    4991   all
## 5         in      5677    5    4546   all
## 6         to      2830    6    2455   all
## 7        for      2099    7    1903   all
## 8         is      1868    8    1763   all
## 9        was      1491    9    1395   all
## 10      href      1377   10    1083   all
## 11        on      1375   11    1275   all
## 12      from      1260   12    1183   all
## 13    target      1177   13     904   all
## 14    _blank      1173   14     900   all
## 15      with       973   15     931   all
## 16        as       960   16     868   all
## 17       his       941   17     848   all
## 18        by       889   18     841   all
## 19      that       852   19     830   all
## 20        it       818   20     757   all
## 21        an       781   21     746   all
## 22        at       740   22     699   all
## 23     these       732   23     731   all
## 24       one       712   24     676   all
## 25        he       711   25     647   all
## 26       you       665   26     565   all
## 27      it's       635   27     610   all
## 28      clue       613   28     559   all
## 29      name       603   29     594   all
## 30        or       561   30     537   all
## 31      crew       504   31     503   all
## 32         i       492   32     370   all
## 33     first       486   33     474   all
## 34       are       451   34     426   all
## 35      here       429   35     423   all
## 36       its       420   36     403   all
## 37       who       395   37     389   all
## 38       and       369   38     338   all
```

```
## 39     when     366    39     358    all
## 40       be     361    40     350    all
## 41     like     351    41     335    all
## 42      has     333    42     328    all
## 43      her     330    43     293    all
## 44     city     328    44     307    all
## 45    about     327    45     318    all
## 46        2     323    46     311    all
## 47   called     304    47     303    all
## 48     type     302    48     300    all
## 49     seen     297    49     297    all
## 50      but     297    49     295    all
## 51    named     296    51     292    all
## 52      man     296    51     286    all
## 53      can     288    53     281    all
## 54     have     285    54     278    all
## 55      not     278    55     264    all
## 56     film     265    56     258    all
## 57      new     264    57     255    all
## 58 country     260    58     256    all
## 59    after     249    59     247    all
## 60      she     248    60     222    all
```

```
twentyFreqWordsQuestion1600 = textstat_frequency(questionCleaned1600, n=60)
twentyFreqWordsQuestion1600
```

```
##      feature frequency rank docfreq group
## 1       href      1077    1     859    all
## 2     target       904    2     706    all
## 3     _blank       902    3     704    all
## 4        one       615    4     590    all
## 5       name       588    5     578    all
## 6      first       466    6     458    all
## 7       clue       393    7     361    all
## 8       city       341    8     329    all
## 9       crew       321    9     319    all
## 10      like       300   10     291    all
## 11      type       297   11     294    all
## 12     named       294   12     292    all
## 13         2       291   13     278    all
## 14    called       287   14     280    all
## 15       man       268   15     260    all
## 16      seen       263   16     262    all
## 17       new       249   17     240    all
## 18      film       244   18     240    all
## 19      also       238   19     238    all
## 20       can       236   20     232    all
## 21     title       228   21     228    all
## 22      word       227   22     217    all
## 23      used       226   23     224    all
```

```
## 24      known      219   24    218   all
## 25    country      208   25    206   all
## 26        u.s      193   26    192   all
## 27       made      190   27    188   all
## 28      state      185   28    181   all
## 29      novel      174   29    173   all
## 30     played      168   30    162   all
## 31     french      164   31    163   all
## 32       term      163   32    163   all
## 33       part      160   33    157   all
## 34       king      159   34    151   all
## 35      wrote      158   35    156   all
## 36      years      155   36    154   all
## 37     became      150   37    148   all
## 38       said      150   37    147   all
## 39      means      149   39    147   all
## 40          3      148   40    146   all
## 41     island      147   41    136   all
## 42        war      140   42    127   all
## 43       play      138   43    136   all
## 44        now      134   44    133   all
## 45        won      129   45    127   all
## 46       last      127   46    126   all
## 47      world      127   46    125   all
## 48   american      127   46    127   all
## 49          1      127   46    120   all
## 50          s      126   50    122   all
## 51      greek      122   51    121   all
## 52      latin      122   51    122   all
## 53       john      121   53    118   all
## 54      movie      121   53    120   all
## 55        get      121   53    118   all
## 56       show      120   56    119   all
## 57    meaning      119   57    117   all
## 58       born      119   57    118   all
## 59       work      119   57    117   all
## 60      south      117   60    114   all
```

```
twentyFreqWordsQuestion2000 = textstat_frequency(questionCleaned2000, n=60)
twentyFreqWordsQuestion2000
```

```
##        feature frequency rank docfreq group
## 1         href      1593    1    1228   all
## 2       target      1375    2    1040   all
## 3       _blank      1371    3    1036   all
## 4         name       735    4     717   all
## 5         clue       649    5     586   all
## 6          one       613    6     588   all
## 7         crew       542    7     540   all
## 8        first       515    8     500   all
```

```
## 9        seen      376   9   376   all
## 10     called      361  10   356   all
## 11       like      326  11   313   all
## 12      named      315  12   311   all
## 13       type      301  13   294   all
## 14       city      299  14   289   all
## 15       film      276  15   271   all
## 16      title      275  16   273   all
## 17          2      268  17   252   all
## 18        new      268  17   260   all
## 19       also      268  17   268   all
## 20      known      267  20   265   all
## 21       used      261  21   253   all
## 22       word      255  22   245   all
## 23        man      252  23   244   all
## 24        can      230  24   226   all
## 25        u.s      221  25   217   all
## 26       made      221  25   218   all
## 27      state      204  27   189   all
## 28      wrote      204  27   201   all
## 29     french      201  29   199   all
## 30      novel      200  30   200   all
## 31    country      197  31   195   all
## 32      jimmy      181  32   181   all
## 33       last      177  33   174   all
## 34       term      168  34   166   all
## 35   american      166  35   164   all
## 36    reports      164  36   161   all
## 37       king      160  37   154   all
## 38     became      160  37   160   all
## 39        won      156  39   152   all
## 40       said      155  40   155   all
## 41        war      153  41   146   all
## 42     island      152  42   145   all
## 43    capital      152  42   146   all
## 44      years      150  44   149   all
## 45       part      147  45   142   all
## 46      latin      145  46   145   all
## 47      means      145  46   140   all
## 48       play      145  46   141   all
## 49      sarah      145  46   145   all
## 50          3      144  50   139   all
## 51      greek      142  51   140   all
## 52      world      141  52   137   all
## 53       work      140  53   140   all
## 54    british      140  53   139   all
## 55        get      138  55   134   all
## 56  president      138  55   133   all
## 57      great      136  57   131   all
## 58    century      135  58   131   all
```

```
## 59      author      133   59      132   all
## 60         may      132   60      132   all
```

```
twentyFreqWordsQuestionNone = textstat_frequency(questionCleanedNone, n=60)
twentyFreqWordsQuestionNone
```

```
##         feature frequency rank docfreq group
## 1             2       360    1     344   all
## 2         first       351    2     338   all
## 3          name       311    3     295   all
## 4           one       308    4     285   all
## 5           u.s       228    5     226   all
## 6           man       207    6     200   all
## 7         named       169    7     164   all
## 8         state       155    8     141   all
## 9     president       153    9     139   all
## 10            1       145   10     134   all
## 11         last       137   11     136   all
## 12      country       132   12     123   all
## 13         city       131   13     125   all
## 14            3       124   14     120   all
## 15        years       119   15     114   all
## 16         said       117   16     117   all
## 17        title       112   17     111   all
## 18         film       111   18     109   all
## 19          new       104   19     101   all
## 20       called       101   20      99   all
## 21         word       100   21      95   all
## 22        wrote        96   22      96   all
## 23        novel        93   23      92   all
## 24       became        92   24      92   all
## 25        whose        92   24      91   all
## 26        world        87   26      85   all
## 27      century        86   27      85   all
## 28    character        83   28      82   all
## 29         made        82   29      81   all
## 30         book        76   30      74   all
## 31       famous        73   31      72   all
## 32         used        73   31      70   all
## 33        known        71   33      71   all
## 34         time        70   34      69   all
## 35         born        69   35      66   all
## 36            4        69   35      68   all
## 37         work        69   35      69   all
## 38     american        68   38      67   all
## 39       people        67   39      66   all
## 40        names        67   39      63   all
## 41         best        66   41      61   all
## 42          won        65   42      62   all
## 43          war        64   43      60   all
```

```
## 44      died        62   44     61    all
## 45       may        62   44     59    all
## 46    states        62   44     61    all
## 47       men        62   44     62    all
## 48      play        61   48     60    all
## 49     since        61   48     60    all
## 50   capital        61   48     56    all
## 51      year        61   48     57    all
## 52       now        60   52     60    all
## 53  national        58   53     57    all
## 54   largest        57   54     55    all
## 55      term        55   55     53    all
## 56   million        53   56     50    all
## 57      part        53   56     52    all
## 58   world's        52   58     52    all
## 59 countries        52   58     52    all
## 60    island        52   58     51    all
```

```r
#import cleaned csv's
allTopWords = read.csv('All_Top_Words_Jeopardy.csv')
valueTopWords = read.csv("Value_Question_Top_Words_Jeopardy.csv")
roundQuestionTopWords = read.csv("Round_Question_Top_Words_Jeopardy.csv")
roundCategoryTopWords = read.csv("Round_Category_Top_Words_Jeopardy.csv")

library(ggplot2)
library(dplyr)

#bar graph of all category top words
bar1 =ggplot(allTopWords, aes(ï..Top_All_Category, Top_All_Category_Freq)) +
  geom_bar(stat="identity", fill="#060CE9") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all question top words
bar2 =ggplot(allTopWords, aes(Top_All_Question, Top_All_Question_Count)) +
  geom_bar(stat="identity", fill="black") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

bar1
```

bar2

```r
library(ggplot2)
library(dplyr)

#bar graph of all value 400 Question top words
bar3 =ggplot(valueTopWords, aes(ï..400_Top_Q, X400_Top_Q_Freq)) +
  geom_bar(stat="identity", fill="green") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all value 600 Question top words
bar4 =ggplot(valueTopWords, aes(X600_Top_Q, X600_Top_Q_Freq)) +
  geom_bar(stat="identity", fill="green") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all value 800 Question top words
bar5 =ggplot(valueTopWords, aes(X800_Top_Q, X800_Top_Q_Freq)) +
  geom_bar(stat="identity", fill="green") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all value 1000 Question top words
bar6 =ggplot(valueTopWords, aes(X1000_Top_Q, X1000_Top_Q_Freq)) +
  geom_bar(stat="identity", fill="green") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all value 1200 Question top words
bar7 =ggplot(valueTopWords, aes(X1200_Top_Q, X1200_Top_Q_Freq)) +
  geom_bar(stat="identity", fill="green") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all value 1600 Question top words
bar8 =ggplot(valueTopWords, aes(X1600_Top_Q, X1600_Top_Q_Freq)) +
  geom_bar(stat="identity", fill="green") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all value 2000 Question top words
bar9 =ggplot(valueTopWords, aes(X2000_Top_Q, X2000_Top_Q_Freq)) +
  geom_bar(stat="identity", fill="green") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all value Final Jeopardy (None) Question top words
bar10 =ggplot(valueTopWords, aes(None_Top_Q, None_Top_Q_Freq)) +
  geom_bar(stat="identity", fill="green") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

bar3
```
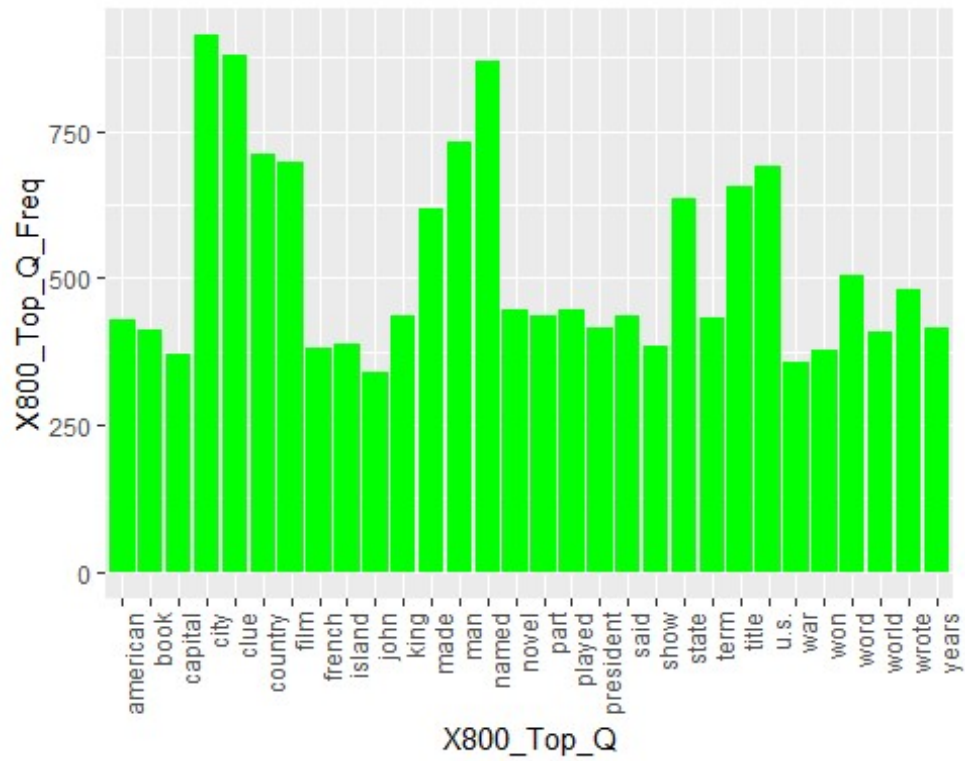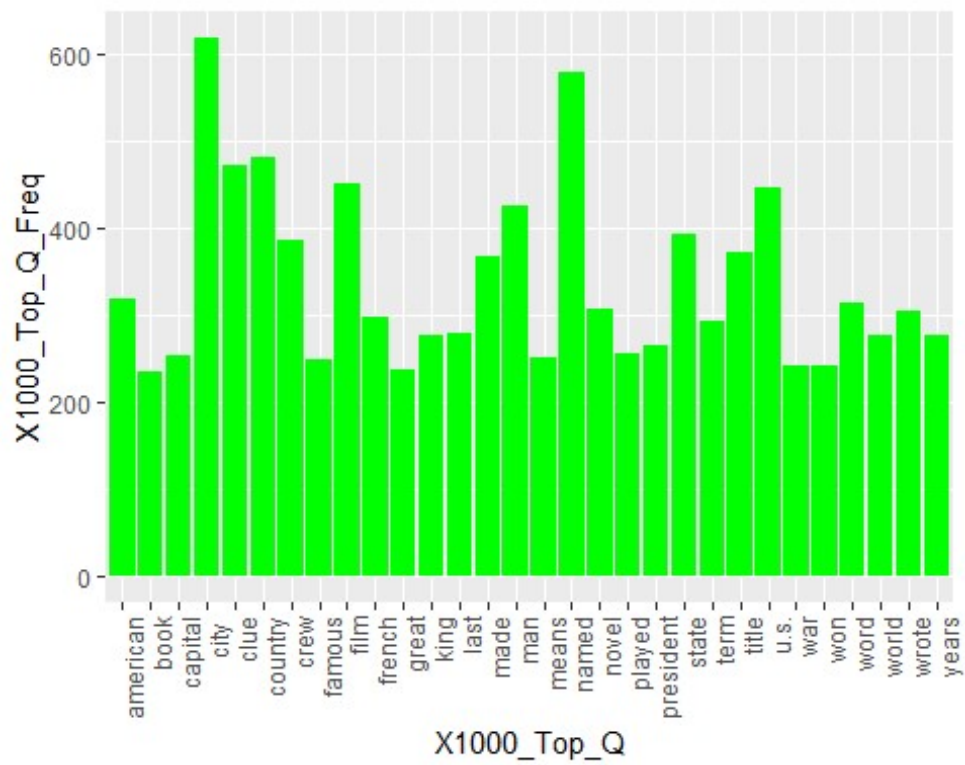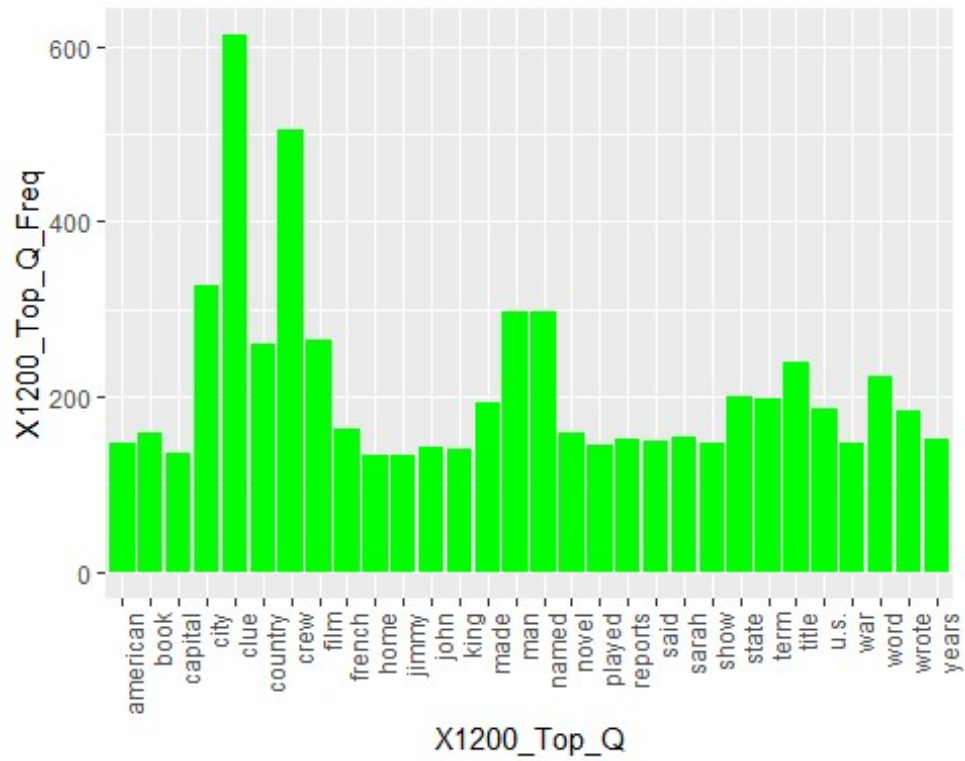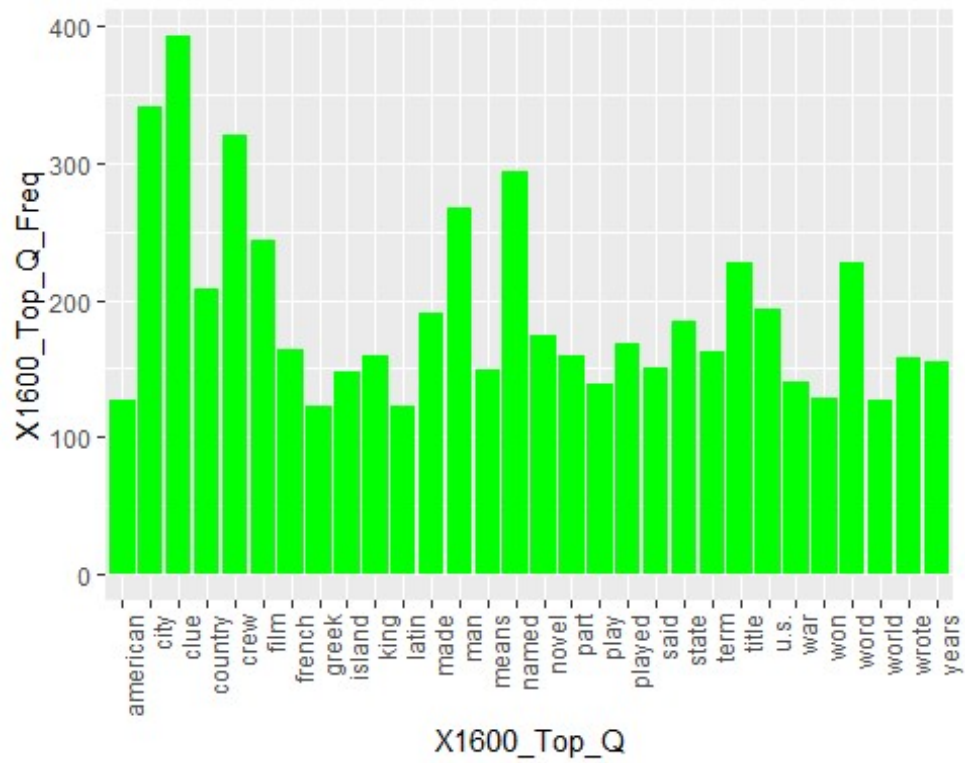
bar4



bar5

bar6



bar7

bar8



bar9

bar10

```
#bar graph of all round top category words (Jeopardy)
bar11 =ggplot(roundCategoryTopWords, aes(ï..Jeopardy_Top_Cat,
Jeopardy_Top_Cat_Freq)) +
  geom_bar(stat="identity", fill="blue") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all round top category words (Double Jeopardy)
bar12 =ggplot(roundCategoryTopWords, aes(Double_Jeopardy_Top_Cat,
Double_Jeopardy_Top_Cat_Freq)) +
  geom_bar(stat="identity", fill="blue") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all round top category words (Final Jeopardy)
bar13 =ggplot(roundCategoryTopWords, aes(Final_Jeopardy_Top_Cat,
Final_Jeopardy_Top_Cat_Freq)) +
  geom_bar(stat="identity", fill="blue") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

bar11
```
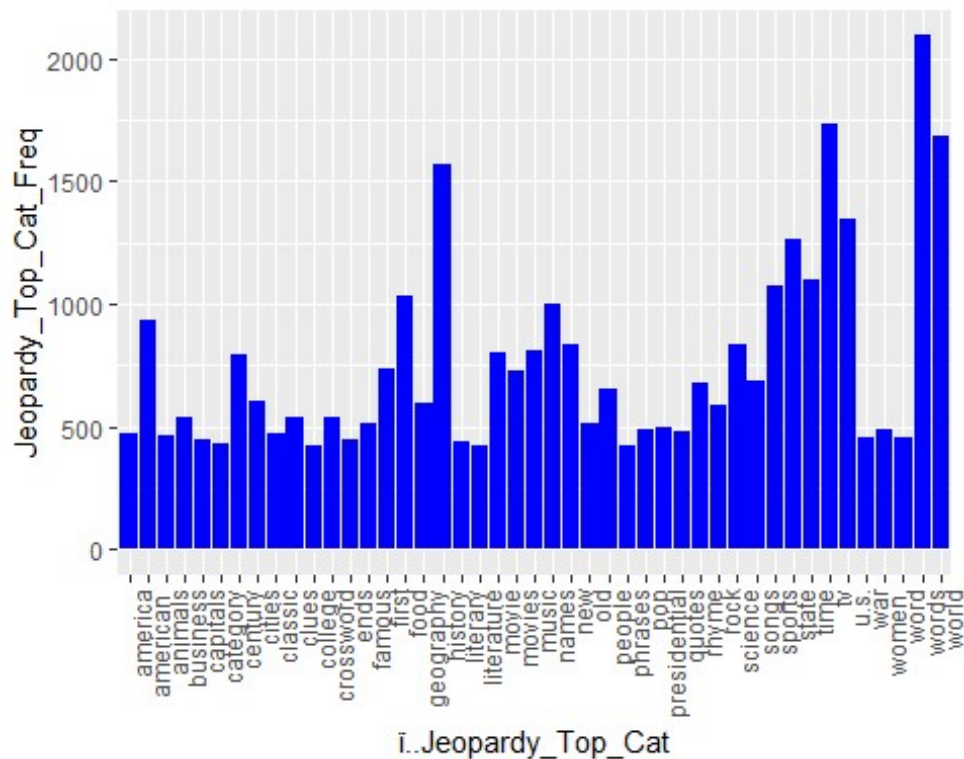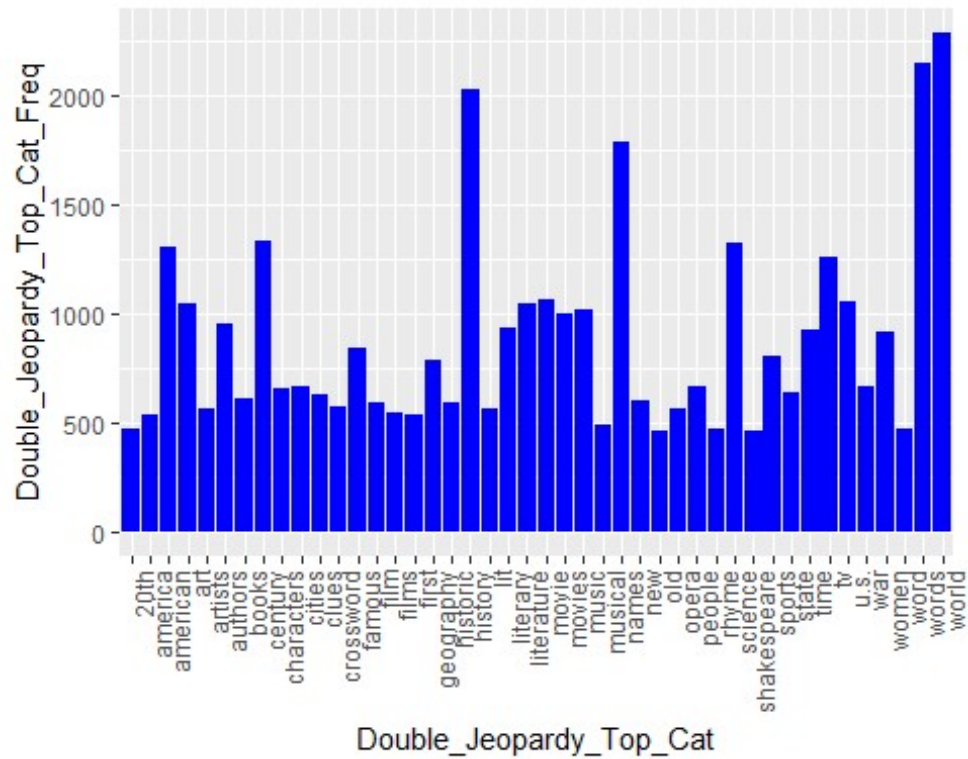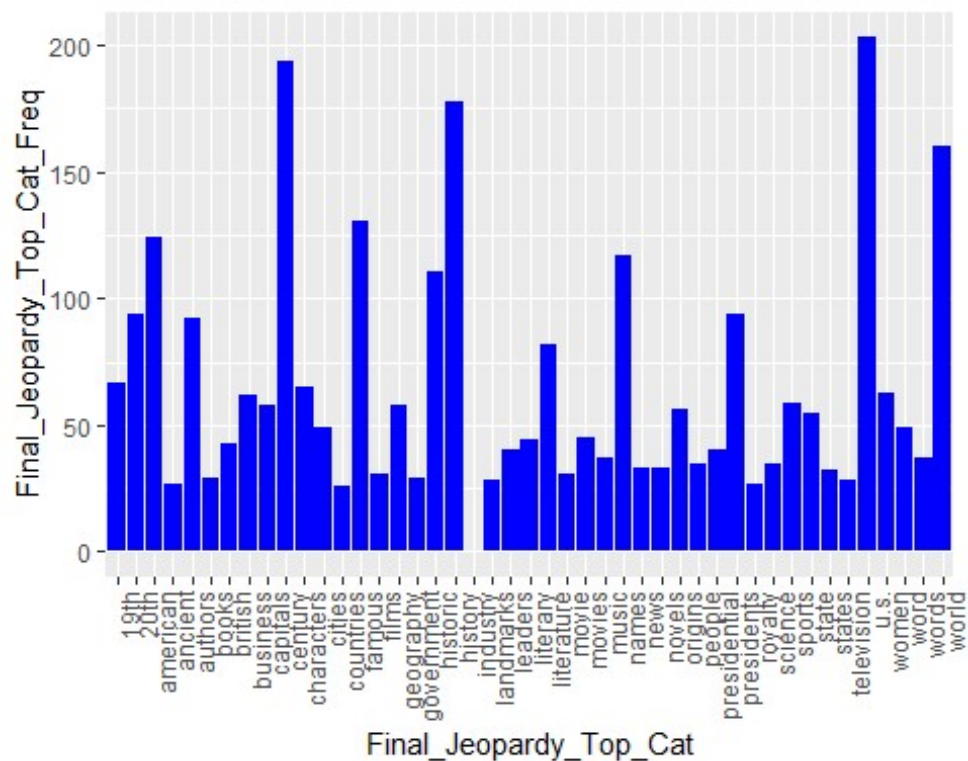


bar12

bar13

## Warning: Removed 1 rows containing missing values (position_stack).

```
#bar graph of all round top question words (Jeopardy)
bar14 =ggplot(roundQuestionTopWords, aes(ï..Jeopardy_Top_Q,
Jeopardy_Top_Q_Freq)) +
  geom_bar(stat="identity", fill="black") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all round top question words (Double Jeopardy)
bar15 =ggplot(roundQuestionTopWords, aes(Double_Jeopardy_Top_Q,
Double_Jeopardy_Top_Q_Freq)) +
  geom_bar(stat="identity", fill="black") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

#bar graph of all round top question words (Final Jeopardy)
bar16 =ggplot(roundQuestionTopWords, aes(Final_Jeopardy_Top_Q,
Final_Jeopardy_Top_Q_Freq)) +
  geom_bar(stat="identity", fill="black") +
  theme(axis.text.x=element_text(angle = 90, hjust = 1.0))

bar14
```
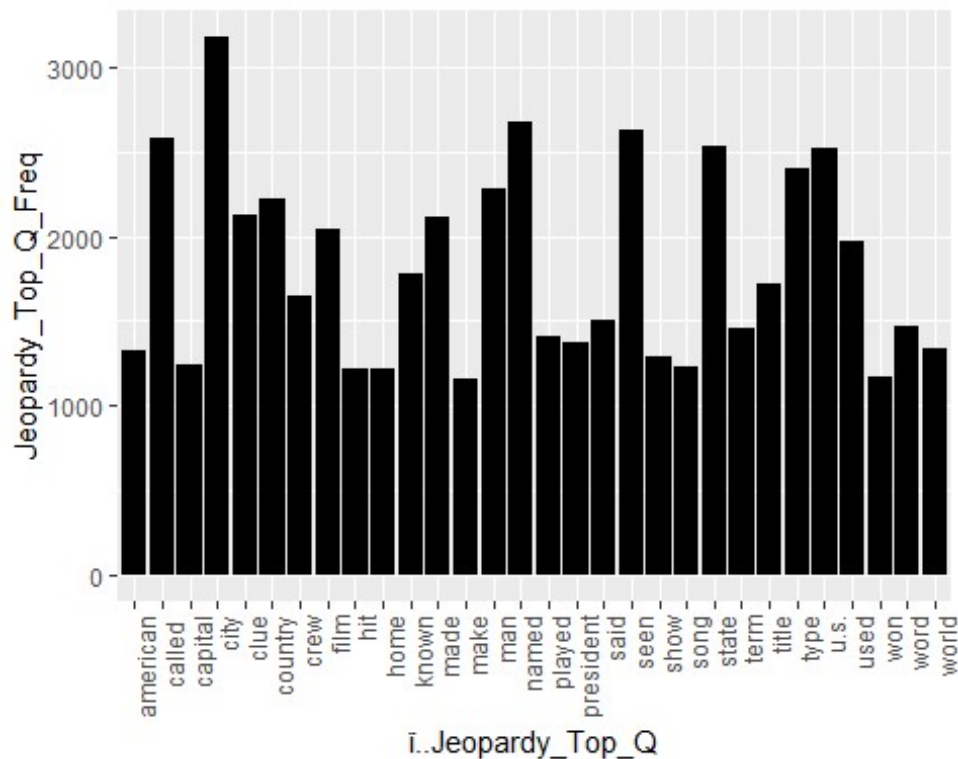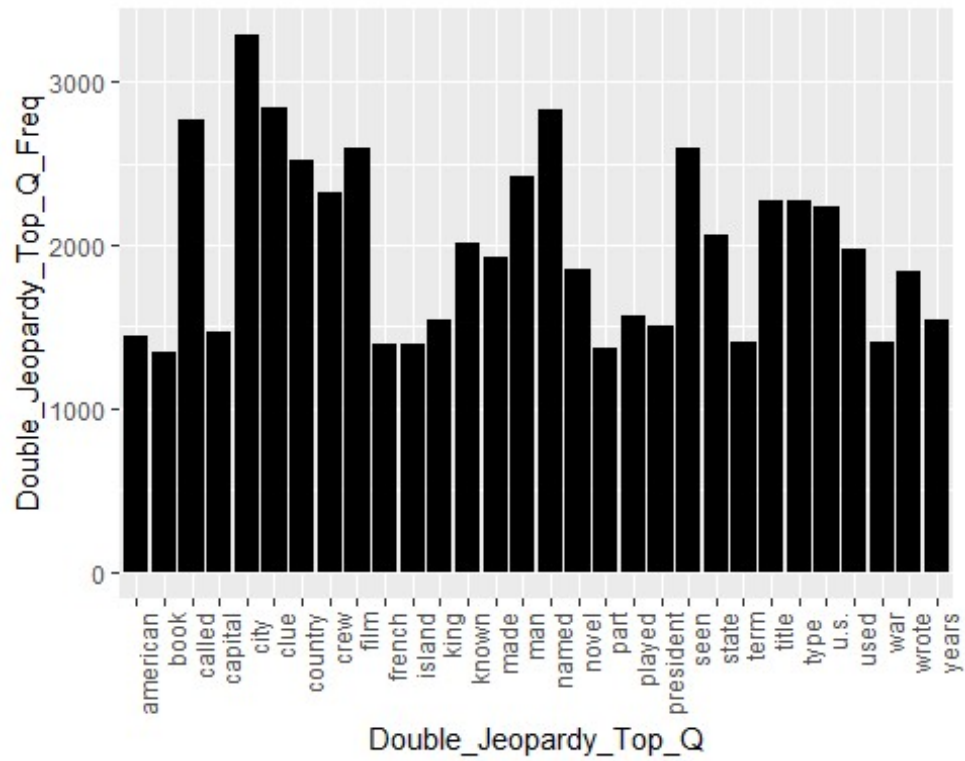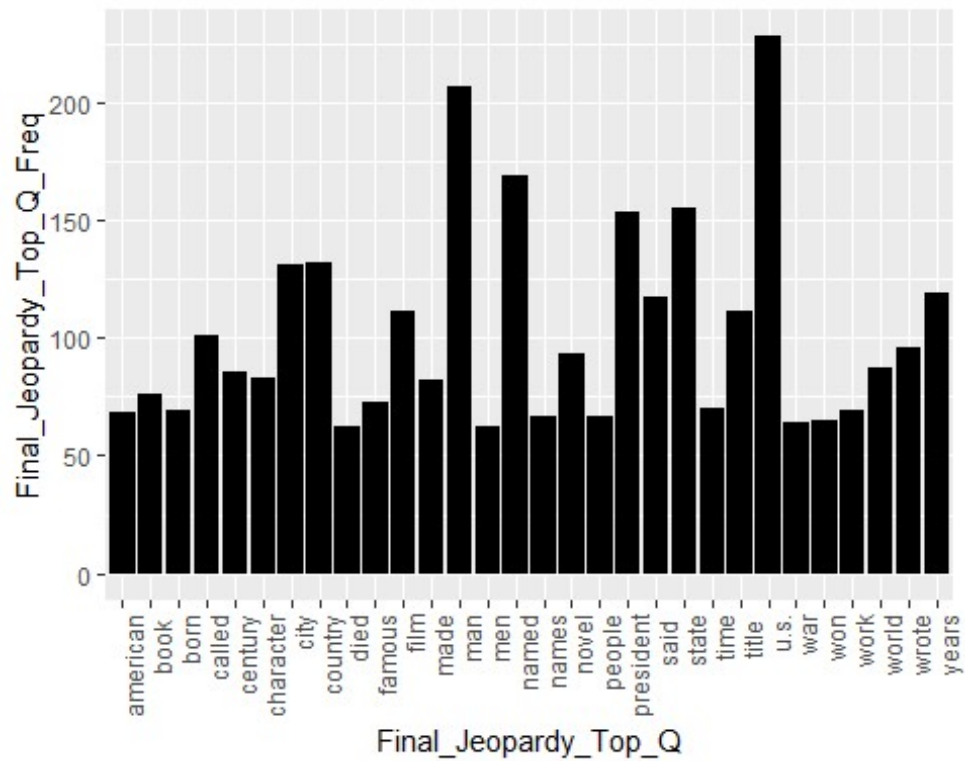


bar15

bar16

```r
#combining word columns from cleaned dataframes
  allTopWordsList = c(allTopWords$ï..Top_All_Category)
  roundCategoryTopWordsList = c(roundCategoryTopWords$ï..Jeopardy_Top_Cat)
  roundDJCategoryTopWordsList =
c(roundCategoryTopWords$Double_Jeopardy_Top_Cat)
  roundFJCategoryTopWordsList =
c(roundCategoryTopWords$Final_Jeopardy_Top_Cat)

  everythingList = c(allTopWordsList, roundCategoryTopWordsList,
roundDJCategoryTopWordsList, roundFJCategoryTopWordsList)
  allCatDf = data.frame(words = everythingList)
  write.csv(allCatDf, "test.csv")
```