



WHAT IS THE BEST WORDLE STARTING WORD?



Grand Challenges PRESENTATION - Georgia Mason

I N T R O D U C T I O N

In this presentation, I will be sharing the data found by analyzing Wordle, including:

- Most common letters
- Average amount of vowels
- Comparison with the English dictionary
- The best 2 starting words
- The worst starting word

Relevance:

- Using a better starting word, players have an increased chance to get the answer in a less amount of guesses.

METHODOLOGY

Coded a MATLAB program that read from the data source & found the number of occurrences of each letter in the alphabet:

```
clear;
clc;

alphabet = zeros(1, 26);

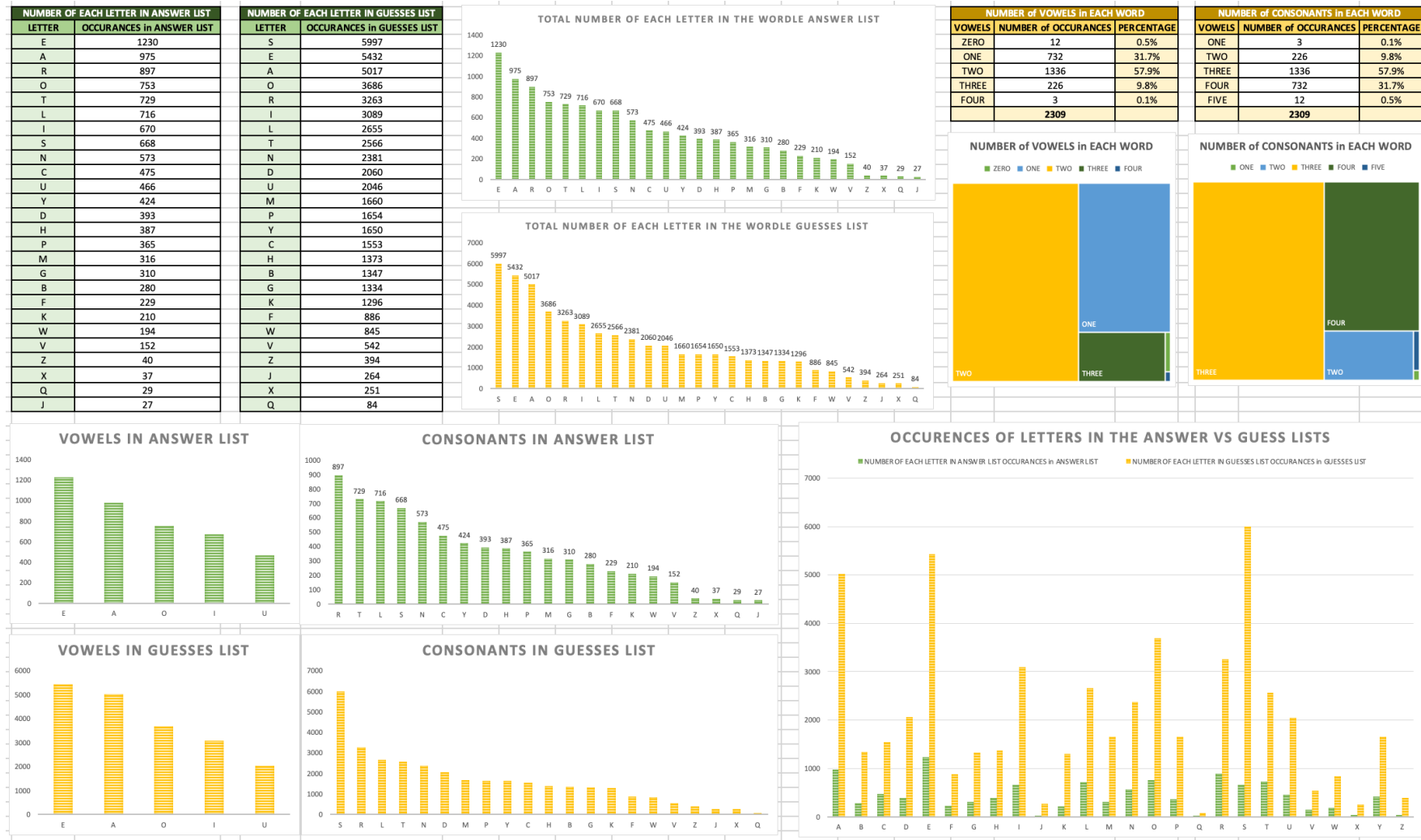
file = fopen('valid_guesses.txt', 'rb'); %reads from file
fseek(file, 0, 'eof');
fileSize = ftell(file);
frewind(file);
data = fread(file, fileSize, 'uint8');
numLines = sum(data == 10);
fclose(file);

file = fopen('valid_guesses.txt');
for i = 1:numLines
    line = fgetl(file); % read one line excluding newline character
    word = char(line);
    for i = 1:5
        if word(i) == 'a'
            alphabet(1) = alphabet(1) + 1;
        elseif word(i) == 'b'
            alphabet(2) = alphabet(2) + 1;
        elseif word(i) == 'c'
            alphabet(3) = alphabet(3) + 1;
        elseif word(i) == 'd'
            alphabet(4) = alphabet(4) + 1;
        elseif word(i) == 'e'
            alphabet(5) = alphabet(5) + 1;
        elseif word(i) == 'f'
            alphabet(6) = alphabet(6) + 1;
        elseif word(i) == 'g'
            alphabet(7) = alphabet(7) + 1;
```

```
elseif word(i) == 'h'
    alphabet(8) = alphabet(8) + 1;
elseif word(i) == 'i'
    alphabet(9) = alphabet(9) + 1;
elseif word(i) == 'j'
    alphabet(10) = alphabet(10) + 1;
elseif word(i) == 'k'
    alphabet(11) = alphabet(11) + 1;
elseif word(i) == 'l'
    alphabet(12) = alphabet(12) + 1;
elseif word(i) == 'm'
    alphabet(13) = alphabet(13) + 1;
elseif word(i) == 'n'
    alphabet(14) = alphabet(14) + 1;
elseif word(i) == 'o'
    alphabet(15) = alphabet(15) + 1;
elseif word(i) == 'p'
    alphabet(16) = alphabet(16) + 1;
elseif word(i) == 'q'
    alphabet(17) = alphabet(17) + 1;
elseif word(i) == 'r'
    alphabet(18) = alphabet(18) + 1;
elseif word(i) == 's'
    alphabet(19) = alphabet(19) + 1;
elseif word(i) == 't'
    alphabet(20) = alphabet(20) + 1;
elseif word(i) == 'u'
```

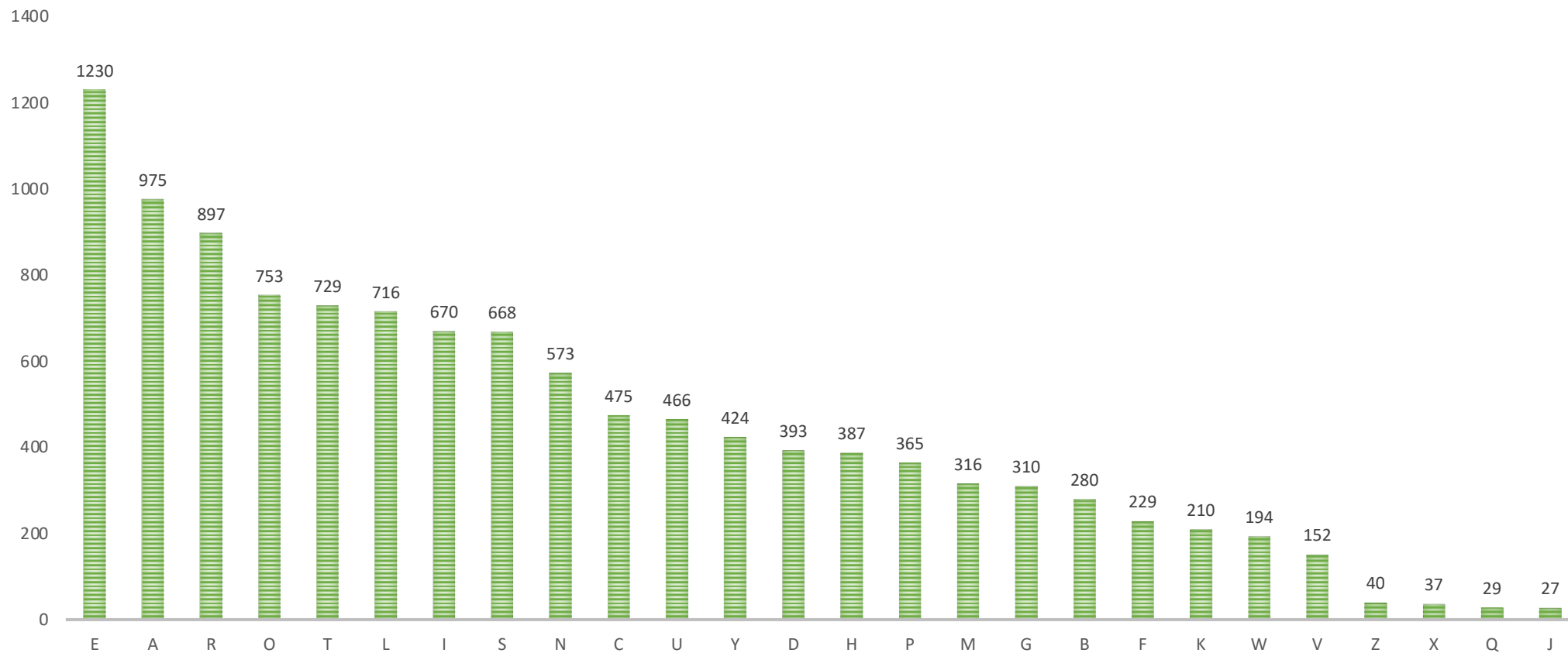
```
    alphabet(20) = alphabet(20) + 1;
elseif word(i) == 'u'
    alphabet(21) = alphabet(21) + 1;
elseif word(i) == 'v'
    alphabet(22) = alphabet(22) + 1;
elseif word(i) == 'w'
    alphabet(23) = alphabet(23) + 1;
elseif word(i) == 'x'
    alphabet(24) = alphabet(24) + 1;
elseif word(i) == 'y'
    alphabet(25) = alphabet(25) + 1;
elseif word(i) == 'z'
    alphabet(26) = alphabet(26) + 1;
end
end
end
for i = 1:26
    fprintf("%d\n", alphabet(i));
end
```

OVERVIEW of the EXCEL SPREADSHEET:



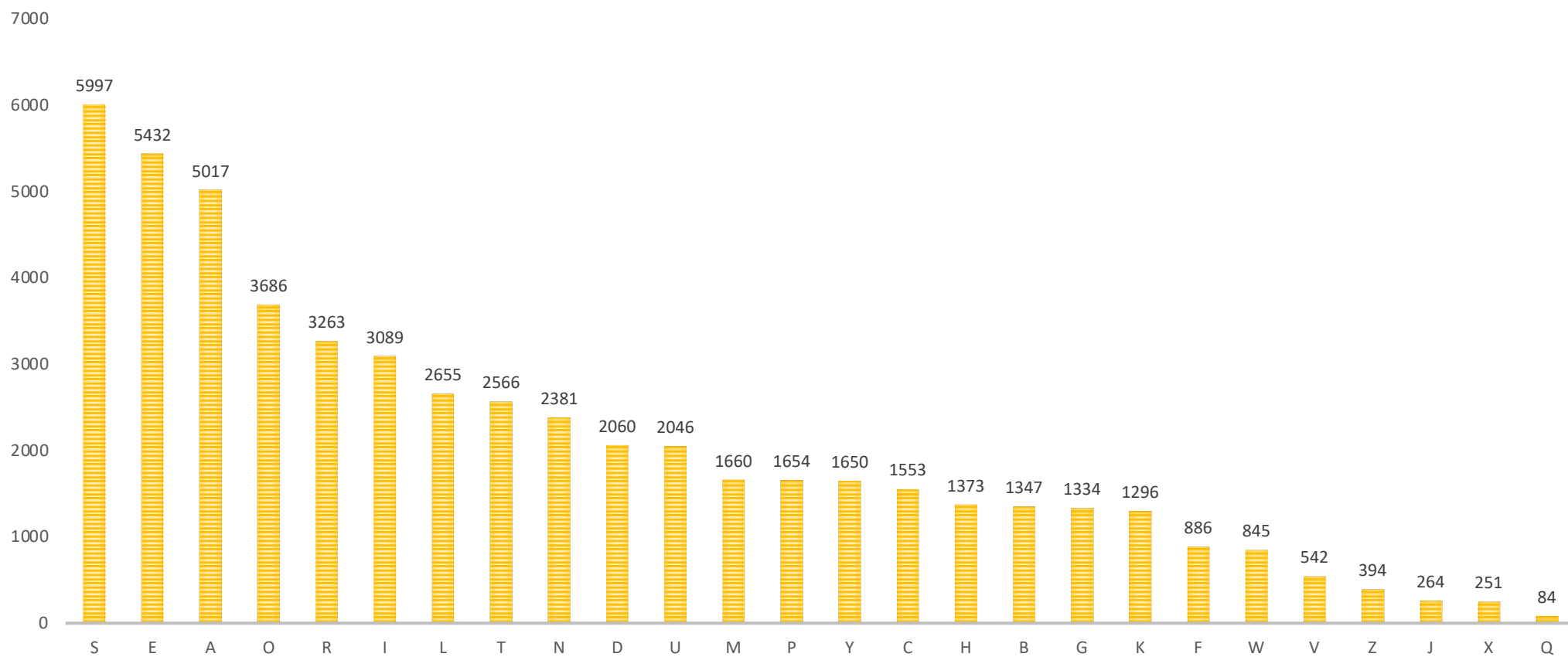
GRAPHICS

TOTAL NUMBER OF EACH LETTER IN THE WORDLE ANSWER LIST



GRAPHICS

TOTAL NUMBER OF EACH LETTER IN THE WORDLE GUESSES LIST

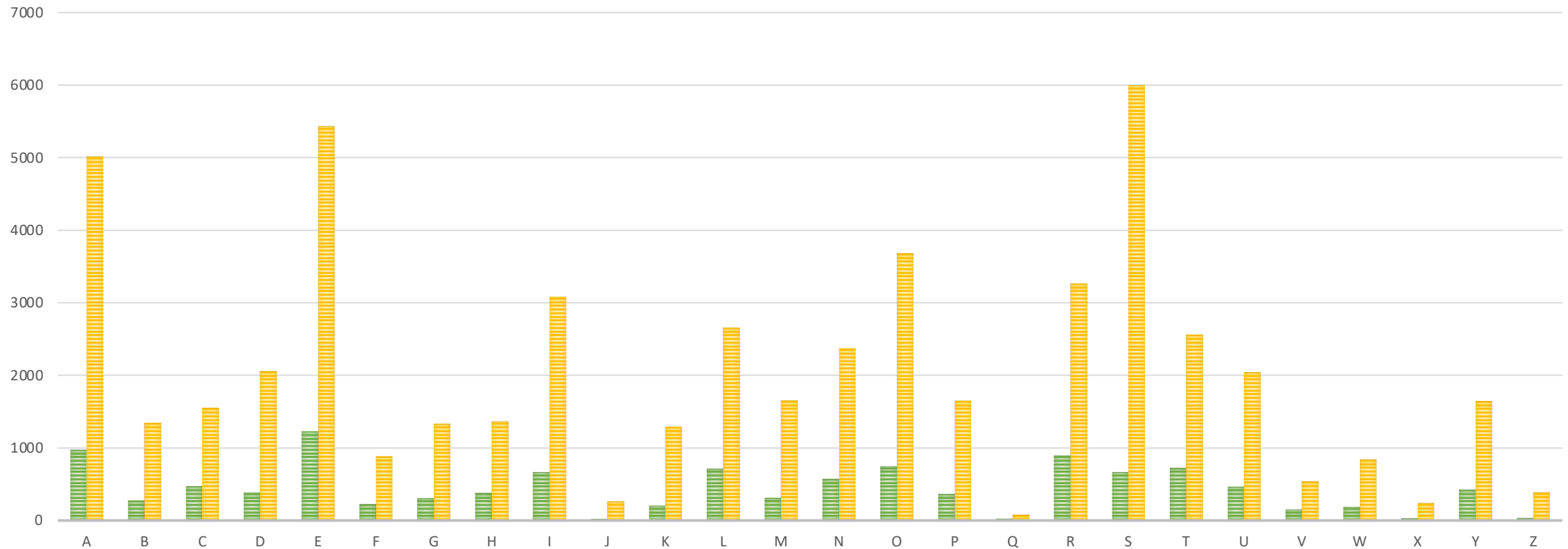


GRAPHICS

OCCURENCES OF LETTERS IN THE ANSWER VS GUESS LISTS

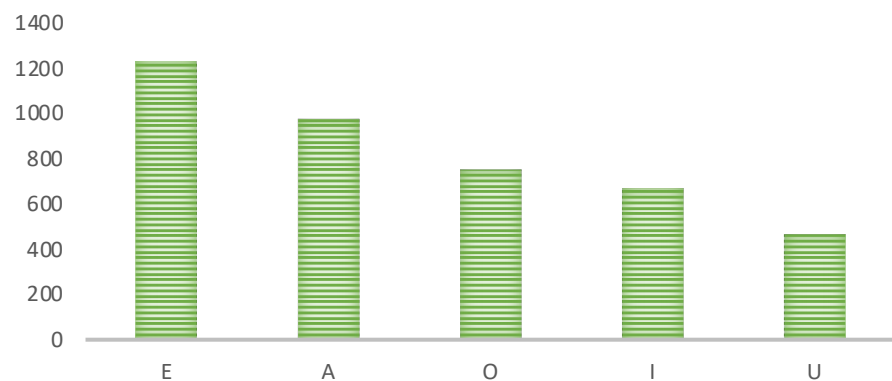
■ NUMBER OF EACH LETTER IN ANSWER LIST OCCURANCES in ANSWER LIST

■ NUMBER OF EACH LETTER IN GUESSES LIST OCCURANCES in GUESSES LIST

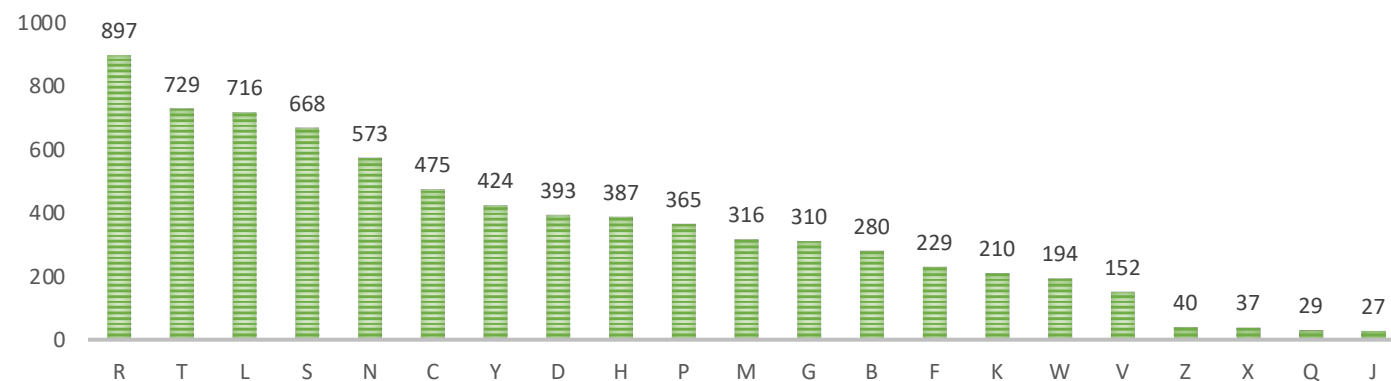


GRAPHICS

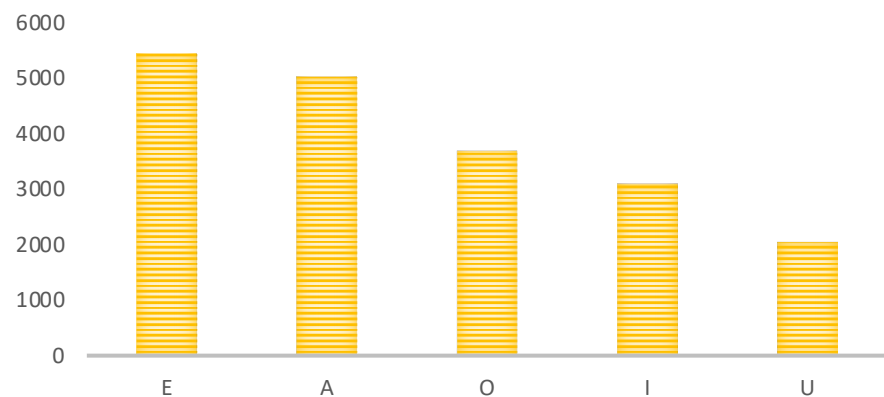
VOWELS IN ANSWER LIST



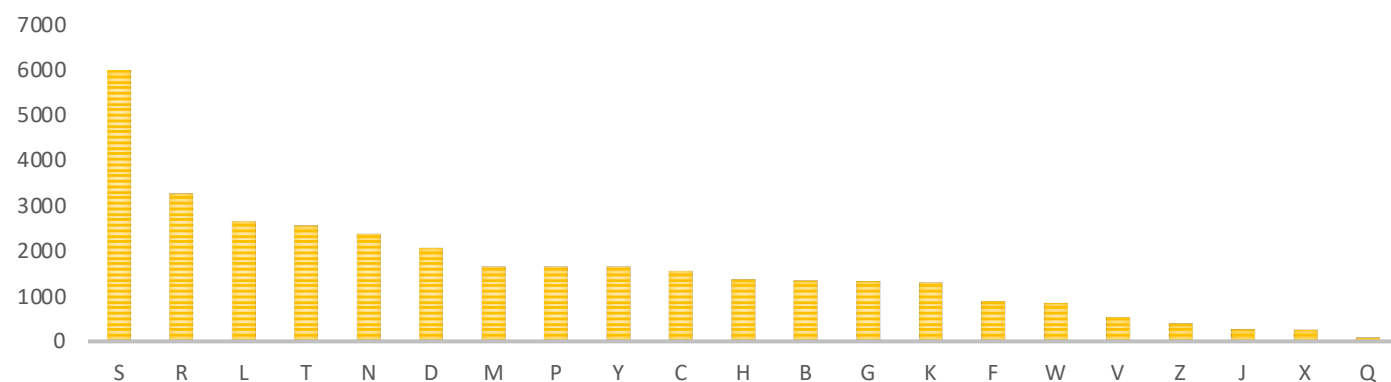
CONSONANTS IN ANSWER LIST



VOWELS IN GUESSES LIST

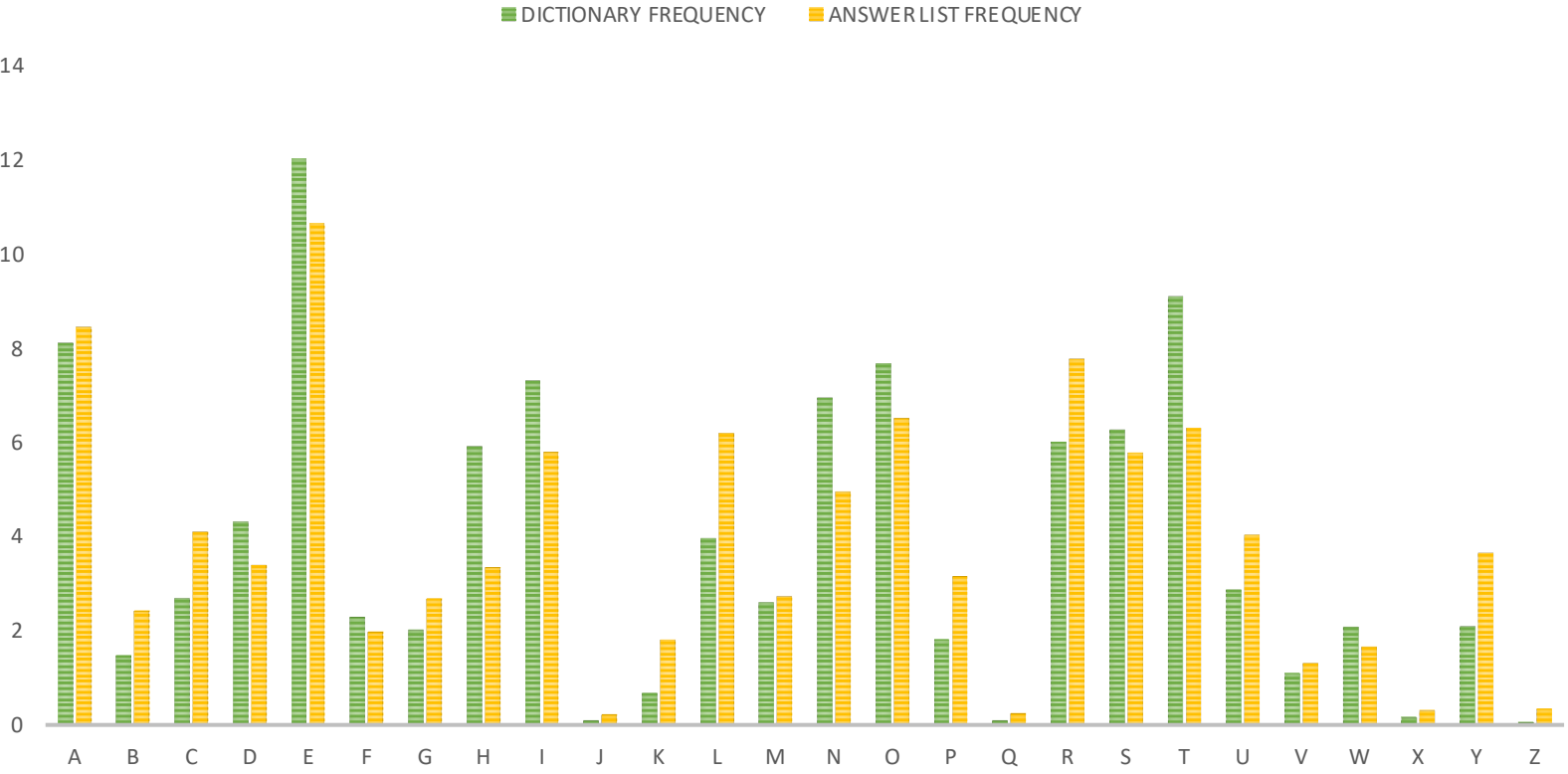


CONSONANTS IN GUESSES LIST



DATA SET USED: Cornell Department of Mathematics – 2004 – English Letter Frequency

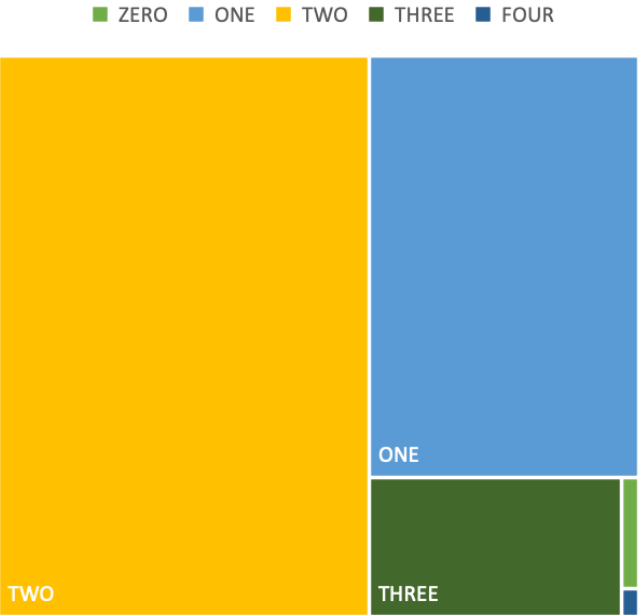
FREQUENCY OF EACH LETTER IN THE WORDLE ANSWER LIST
VS ENGLISH DICTIONARY



DICTIONARY		ANSWER LIST	
LETTER	FREQUENCY	LETTER	FREQUENCY
A	8.12	A	8.45
B	1.49	B	2.43
C	2.71	C	4.11
D	4.32	D	3.4
E	12.02	E	10.65
F	2.3	F	1.98
G	2.03	G	2.69
H	5.92	H	3.35
I	7.31	I	5.8
J	0.1	J	0.23
K	0.69	K	1.82
L	3.98	L	6.2
M	2.61	M	2.74
N	6.95	N	4.96
O	7.68	O	6.52
P	1.82	P	3.16
Q	0.11	Q	0.25
R	6.02	R	7.77
S	6.28	S	5.79
T	9.1	T	6.31
U	2.88	U	4.04
V	1.11	V	1.32
W	2.09	W	1.68
X	0.17	X	0.32
Y	2.11	Y	3.67
Z	0.07	Z	0.35

AVERAGE number of VOWELS + CONSONANTS per word in the answer list:

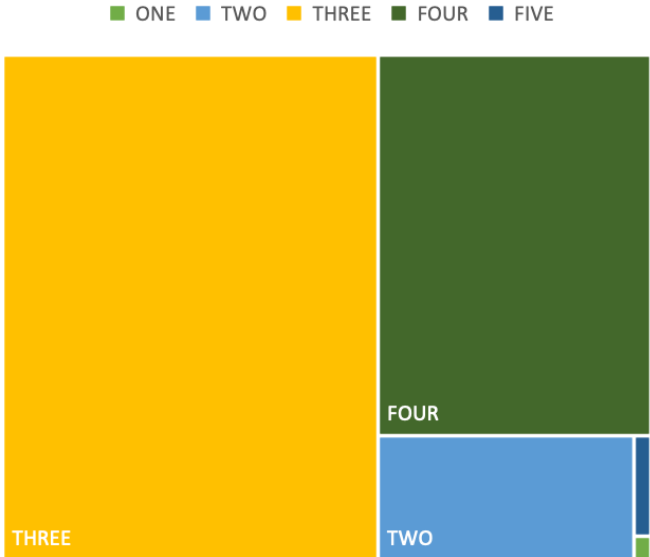
NUMBER of VOWELS in EACH WORD



NUMBER of CONSONANTS in EACH WORD		
VOWELS	NUMBER of OCCURANCES	PERCENTAGE
ONE	3	0.1%
TWO	226	9.8%
THREE	1336	57.9%
FOUR	732	31.7%
FIVE	12	0.5%
	2309	

NUMBER of VOWELS in EACH WORD		
VOWELS	NUMBER of OCCURANCES	PERCENTAGE
ZERO	12	0.5%
ONE	732	31.7%
TWO	1336	57.9%
THREE	226	9.8%
FOUR	3	0.1%
	2309	

NUMBER of CONSONANTS in EACH WORD



FINDING THE BEST STARTING WORD

1. Coded a program to find most common positions of common letters:

A COMMON POSITIONS		E COMMON POSITIONS		R COMMON POSITIONS		O COMMON POSITIONS		T COMMON POSITIONS	
POSITION	OCCURANCES	POSITION	OCCURANCES	POSITION	OCCURANCES	POSITION	OCCURANCES	POSITION	OCCURANCES
3	306	5	422	2	267	2	279	5	253
2	304	4	318	5	212	3	243	1	149
4	162	2	241	3	163	4	132	4	139
1	140	3	177	4	150	5	58	3	111
5	63	1	72	1	105	1	41	2	77

2. Found the anagrams of common letters and found which word has the letters in the most common positions:

OATER:		ORATE:		ROATE:	
O	5	O	5	R	5
A	2	R	1	O	1
T	4	A	1	A	1
E	2	T	3	T	3
R	2	E	1	E	1
TOTAL	15	TOTAL	11	TOTAL	11

```
clear;
clc;

position = zeros(1, 5);

file = fopen('valid_answers.txt', 'rb'); %reads from file
fseek(file, 0, 'eof');
fileSize = ftell(file);
frewind(file);
data = fread(file, fileSize, 'uint8');
numLines = sum(data == 10);
fclose(file);

file = fopen('valid_answers.txt');

for i = 1:numLines
    line = fgetl(file); % read one line excluding newline character
    word = char(line);
    num_of_vowels = 0;
    for j = 1:5
        if word(j) == 'c'
            if j == 1;
                position(1) = position(1) + 1;
            elseif j == 2
                position(2) = position(2) + 1;
            elseif j == 3
                position(3) = position(3) + 1;
            elseif j == 4
                position(4) = position(4) + 1;
            elseif j == 5
                position(5) = position(5) + 1;
            end
        end
    end

    for i = 1:5
        fprintf("%d\n", position(i));
    end
end
```

I added in the sixth most common letter to the anagram solver, and condensed this list to words that ARE in the answer list.

I repeated the same process, and found that ALERT is the best starting word. This word has also not been the answer yet, so it is possible to get this word in one some day.

L COMMON POSITIONS	
POSITION	OCCURANCES
2	200
4	162
5	155
3	112
1	87

ALERT:	
A	4
L	1
E	4
R	4
T	1
TOTAL	14

ALTER:	
A	4
L	1
T	4
E	2
R	2
TOTAL	15

LATER:	
L	5
A	2
T	4
E	2
R	2
TOTAL	15

I then found another word to use after ALERT.

I repeated the same process as finding the best word, but using the letters, S, I, N, C, and O, as this letter was not used in the best guess.

The two anagrams in the answer list found both achieved the same score, meaning either could be used in a second guess.

S COMMON POSITIONS		I COMMON POSITIONS	
POSITION	OCCURANCES	POSITION	OCCURANCES
1	365	3	266
4	171	2	201
3	80	4	158
5	36	1	34
2	16	5	11

N COMMON POSITIONS		C COMMON POSITIONS	
POSITION	OCCURANCES	POSITION	OCCURANCES
4	182	1	198
3	137	4	150
5	130	3	56
2	87	2	40
1	37	5	31

SCION		SONIC	
S	1	S	1
C	4	O	1
I	1	N	2
O	3	I	3
N	3	C	5
TOTAL	12	TOTAL	12

FINDING THE WORST WORD

I compiled a list of the words with very uncommon letters and multiple double letters, then added up the score of the position in the most common letters list.

The highest score, with the letters that are the least common was FUZZY.

This method does not take into account double letters. In the future, I would be to code a program that calculates the average percentage of receiving a green or yellow letter using the guess.

JAZZY		XYLYL		FUFFY	
J	26	X	24	F	22
A	2	Y	12	U	11
Z	23	L	6	F	22
Z	23	Y	12	F	22
Y	12	L	6	Y	12
TOTAL	86	TOTAL	60	TOTAL	89
FUZZY		JEEZE		PHPHT	
F	22	J	26	P	15
U	11	E	1	H	14
Z	23	E	0	P	25
Z	23	Z	23	H	14
Y	12	E	0	T	5
TOTAL	91	TOTAL	50	TOTAL	73
QAJAQ		ZIZIT		MAMMA	
Q	25	Z	23	M	16
A	2	I	7	A	2
J	26	Z	23	M	16
A	2	I	7	M	16
Q	25	T	8	A	2
TOTAL	80	TOTAL	68	TOTAL	52

C O N C L U S I O N S

- ALERT is the best first guess
- SONIC or SCION is the best second guess
- There is a slight difference in common letters from the GUESS and ANSWER lists
- There is a correlation between common letters in Wordle and the English dictionary
- FUZZY is one of the worst starting words
- There is an average of TWO vowels and THREE consonants per word



T H A N K Y O U

FOR LISTENING 😊

B I B L I O G R A P H Y

CORNEL DEPARTMENT of MATHEMATICS – English Letter Frequency - 2004

- <https://www.statista.com/statistics/1328012/popularity-of-wordle-by-age-group-usa/>

KAGGLE - Wordle Valid Guesses & Answers List - Lucas Hohmann – 2022

- <https://www.kaggle.com/datasets/lucashohmann/wordle-valid-guesses-and-answers?resource=download>

ANAGRAM SOLVER – Word Tips

- <https://word.tips/anagram-solver/>