

Regression-Based Word Trials Distribution and Difficulty Classification Model

Summary

Question 1 is asking to develop a model to explain this variation and use the model to create a prediction interval for the number of reported results on March 1, 2023. Determine whether there are any attributes of the word affect the percentage of scores reported that were played in Hard Mode. For question 1, we fit the number of people in the report and the date with exponential regression, and add a confidence interval to obtain a prediction interval. At the same time, we calculated the correlation between the attributes of the word and the proportion of the number of people reporting in the hard mode, and thus came to the conclusion that none of the attributes we found has significant influence on the proportion of the number of people reporting in the hard mode. We think that this is reasonable, because we believe that the percentage of scores reported that were played in Hard Mode is proportional to the number of players who play the game in Hard Mode on that day. However, when choosing whether to play the game in Hard Mode or not, the players do not know anything about the word on that day, so the attributes of the word have no reason to affect this choice.

Question 2 is asking to develop a model to predict the distribution of the reported results. For question 2, we divide the attributes of words into three categories, and use time as one of the independent variables to establish a multiple linear regression model, and finally get a model that can predict the number of correct word attempts.

Question 3 is asking to develop and summarize a model to classify solution words by difficulty. Identify the attributes of a given word that are associated with each classification. For question 3, we use the number of correct word attempts as an indicator for dividing the difficulty of words, and use this indicator to divide the difficulty of words. For a new word, we only need to predict the distribution of the number of correct word attempts based on the model in question 2, and then obtain the difficulty based on the distribution, which means the solution is successful.

Finally, we synthesize the above three models and write a summary report.

Keywords: exponential regression fit; multiple linear regression

Contents

1	Introduction	2
1.1	Background	2
1.2	Questions	2
1.3	Solutions	3
2	Analysis of the Problem	3
2.1	Problem 1: The number of the reported results	3
2.2	Problem 2: The distribution of the reported results	4
2.3	Problem 3: The classification of solution words by difficulty . . .	4
3	Models and Results	4
3.1	Problem 1	4
3.1.1	Exponential regression fitting to predict the number of the reported results	4
3.1.2	Results of problem 1.1	6
3.1.3	Correlation test to explore attributes of words	7
3.1.4	Results of problem 1.2	9
3.2	Problem 2	9
3.2.1	Multiple linear regression forecasting model for distribution	9
3.2.2	Results of problem 2	11
3.2.3	Validation	12
3.3	Problem 3	12
3.3.1	Multiple Linear Regression Classification Model	12
3.3.2	Results of problem 3	13
4	Other Interesting Features of this Data Set	13
5	Summary Letter	14
6	References	17
	Appendices	18

1 Introduction

1.1 Background

Wordle is a popular puzzle currently offered daily by the New York Times. Players try to solve the puzzle by guessing a five-letter word in six tries or less, receiving feedback with every guess. For this version, each guess must be an actual word in English. Guesses that are not recognized as words by the contest are not allowed. Wordle continues to grow in popularity and versions of the game are now available in over 60 languages.

The New York Times website directions for Wordle state that the color of the tiles will change after you submit your word. A yellow tile indicates the letter in that tile is in the word, but it is in the wrong location. A green tile indicates that the letter in that tile is in the word and is in the correct location. A gray tile indicates that the letter in that tile is not included in the word at all.

Players can play in regular mode or “Hard Mode.” Wordle’s Hard Mode makes the game more difficult by requiring that once a player has found a correct letter in a word (the tile is yellow or green), those letters must be used in subsequent guesses. The example in Figure 1 was played in Hard Mode.

Many (but not all) users report their scores on Twitter. For this problem, MCM has generated a file of daily results for January 7, 2022 through December 31, 2022. This file includes the date, contest number, word of the day, the number of people reporting scores that day, the number of players on hard mode, and the percentage that guessed the word in one try, two tries, three tries, four tries, five tries, six tries, or could not solve the puzzle (indicated by X).

1.2 Questions

This contest requires solutions to several questions.

- Develop a model to explain this variation and use the model to create a prediction interval for the number of reported results on March 1, 2023. Determine whether there are any attributes of the word affect the percentage of scores reported that were played in Hard Mode.
- Develop a model to predict the distribution of the reported results. In other words, to predict the associated percentages of (1, 2, 3, 4, 5, 6, X) for a future date. Give a specific example of the prediction for the word EERIE on March 1, 2023.
- Develop and summarize a model to classify solution words by difficulty. Identify the attributes of a given word that are associated with each classification.

1.3 Solutions

For question 1, we fit the number of people in the report and the date with exponential regression, and add a confidence interval to obtain a prediction interval. At the same time, we calculated the correlation between the attributes of the word and the proportion of the number of people reporting in the hard mode, and thus came to the conclusion that none of the attributes we found has significant influence on the proportion of the number of people reporting in the hard mode. We think that this is reasonable, because we believe that the percentage of scores reported that were played in Hard Mode is proportional to the number of players who play the game in Hard Mode on that day. However, when choosing whether to play the game in Hard Mode or not, the players do not know anything about the word on that day, so the attributes of the word have no reason to affect this choice.

For question 2, we divide the attributes of words into three categories, and use time as one of the independent variables to establish a multiple linear regression model, and finally get a model that can predict the number of correct word attempts. After testing, the credibility of the model higher degree.

For question 3, we use the number of correct word attempts as an indicator for dividing the difficulty of words, and use this indicator to divide the difficulty of words. For a new word, we only need to predict the distribution of the number of correct word attempts based on the model in question 2, and then obtain the difficulty based on the distribution, which means the solution is successful.

Finally, we synthesize the above three models and write a summary report.

2 Analysis of the Problem

2.1 Problem 1: The number of the reported results

Problem 1 requires a model to explain the variation of the number of reported results. And this model will be used to create a prediction interval for the number of reported results on March 1, 2023.

In addition, it would like to determine whether there are any attributes of the word affect the percentage of scores reported that were played in Hard Mode.

The first question essentially asks to build a regression model based on the independent variable **date** and the dependent variable **the number of the reported results**, so that the future situation can be predicted according to the model, and a reasonable prediction interval is given.

The second question is essentially asking the correlation between each **attribute** of the word and the **hardmode report percentage**.

2.2 Problem 2: The distribution of the reported results

Problem 2 requires to predict the associated percentages of (1, 2, 3, 4, 5, 6, X) for a future date, hence the most important thing is to find all the variables contributing to the percentages or distribution. We can mainly consider these variables, which is similar to the problem 1 but in addition, the time also needed to be considered.

- 1. daily frequency of a word
- 2. repetition of letters in a word
- 3. similarity between a word and other words
- 4. date

According to the requirements of the topic, we can consider using multiple linear regression models. In regression analysis, if there are two or more independent variables, it is called multiple regression. In fact, a phenomenon is often associated with multiple factors, and the optimal combination of multiple independent variables to predict or estimate the dependent variable is more effective and more realistic than using only one independent variable to predict or estimate. Therefore, the practical significance of multiple linear regression is greater than that of single linear regression.

Once the regression model is derived, we use it to predict the distribution.

2.3 Problem 3: The classification of solution words by difficulty

Problem 3 requires us to develop and summarize a model to classify solution words by difficulty. According to the title, the difficulty of a word is mainly defined by the distribution of the number of answers, so it is considered to extract the elements in the distribution of the number of answers as an index of word difficulty. Then borrow the model for each attribute of the word in the second question to predict the difficulty of the new word.

3 Models and Results

3.1 Problem 1

3.1.1 Exponential regression fitting to predict the number of the reported results

Taking the date as x and the number of the reported results as y to make a $y - x$ relationship diagram, it can be concluded that the change trend of the number of reports is a sharp increase at

first, and then a significant decline after reaching the peak, and it shows a trend of slowing down. Therefore, we consider that the significant growth at the beginning is due to a large number of ne-tizens were attracted by curiosity when the game was first launched, and after the peak, the rate of decline slowed down. The primitive relation is shown in **Figure 1**.

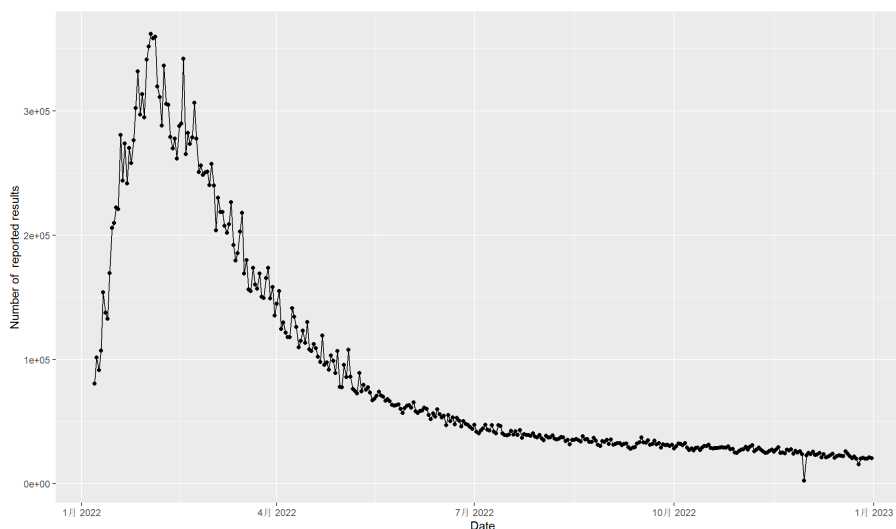


Figure 1: The mapping relation between the number of the reported results and the date

From this figure we see that there is an outlier when the contest number equals to 529, hence this data should be moved. Choose the data with the column number ≥ 115 and do exponential regression fitting.

When fitting a model, estimated parameter values and covariance matrices are usually available. The covariance matrix reflects the uncertainty in the estimated parameter values, and the standard error is a common measure of this uncertainty.

For a linear regression model, the estimated parameter values can be solved by the least squares method, and the covariance matrix can be calculated from the estimated residual sum of squares, explanatory variables, and sample size. In this case, for each parameter estimate, we can quantify its uncertainty by obtaining its standard error via the covariance matrix.

In this exponential fitting model, the `nls()` function is used for fitting, which returns an object containing the model parameter estimates and the covariance matrix.

The function needed to be fitted is considered as:

$$y = e^{a \cdot x + b} + c$$

After fitting we get the function `fit`, use the function to predic y for each x , then calculate the standard error at each x value. Standard error is an indicator used to describe the reliability of an estimator, and it is usually used to measure the difference between a sample statistic and a population parameter. In regression analysis, standard errors are often used to calculate confidence intervals

and significance tests. Details are as follows:

$$\left\{ \begin{array}{l} y_{\text{new}} = \text{predict}(fit) \\ \text{resid} = y - \text{predict}(fit) \\ s = \text{standard}(\text{resid}) \\ \text{mean_resid} = \text{mean}(\text{resid}) \\ \text{standard_error} = s * \sqrt{1 + 1/\text{length}(x) + \frac{(x - \text{mean}(x))^2}{\text{sum}((x - \text{mean}(x))^2)}}} \\ \text{free} = \text{length}(x) - 2 \\ t_{\text{val}} = \text{qt}(0.025, \text{free}) \\ \text{lower} = y_{\text{new}} - t_{\text{val}} * \text{se} \\ \text{upper} = y_{\text{new}} + t_{\text{val}} * \text{se} \end{array} \right.$$

and add the confidence interval to the new data frame, finally get the value and confidence interval at the point March 1, 2023, to get the prediction interval for the number of reported results.

3.1.2 Results of problem 1.1

The exponential regression fitting results to the data are as follows. Among them, with regard to the uncertainty measurement of parameters a, b, and c, for each parameter estimate, we obtain its standard error through the covariance matrix, thereby quantifying its uncertainty.

	Estimate	Std. Error	t value	Pr(> t)
a	-1.570e-02	5.356e-04	-29.32	<2e-16 ***
b	1.601e+01	1.722e-01	92.95	<2e-16 ***
c	2.311e+04	5.438e+02	42.49	<2e-16 ***

After this, normalize the raw data, calculate the standard error at each value of x, and use the qt() function to look up the value of t in the t-distribution table. We calculate the sample mean and degrees of freedom, and finally draw the confidence interval curve, which is shown in Figure 2.

The predicted interval for the number of the reported scores on March 1, 2023 is **[17302, 29975]**, which is a 95% confidence interval. The middle number of it is **23639**, which is the expected value.

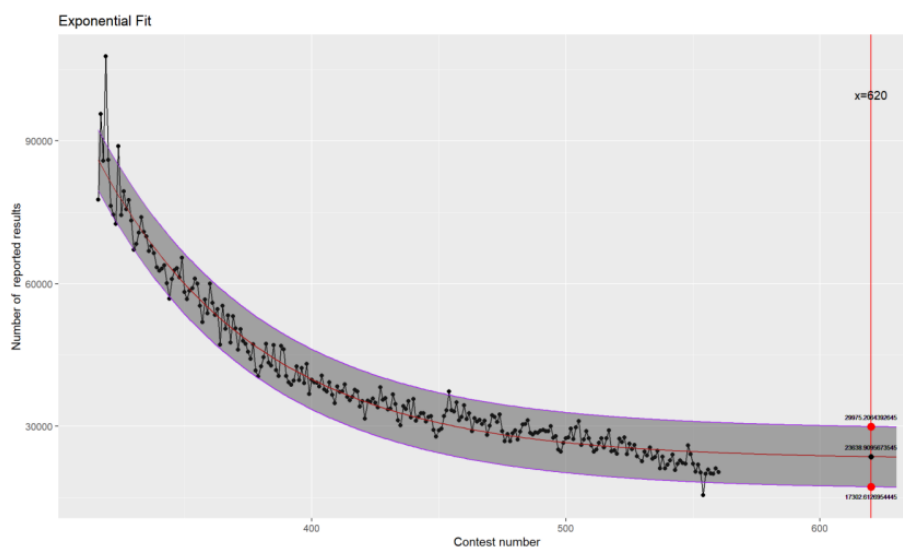


Figure 2: 95% confidence interval

3.1.3 Correlation test to explore attributes of words

Considering the specificity of words, we set the attributes of the words as:

1. word frequency
2. number of unique letters in the word
3. similarity between the word and other words

For the first attribute, the word frequency is an important variable in cognitive processing. Words of high-frequency are perceived and produced faster and more efficiently than words of low-frequency. At the same time, they are easier to recall but more difficult to recognize in episodic memory tasks. For the daily frequency of a word, we get the frequency data from the website [SUBTLEXus](#). In Van Heuven, Mandera, Keuleers, & Brysbaert (QJEP, 2014) we proposed a new frequency measure, the Zipf scale, which is much easier to understand than the usual frequency measures. Zipf values range from 1 to 7, with the values 1-3 indicating low-frequency words (with frequencies of 1 per million words and lower) and the values 4-7 indicating high-frequency words (with frequencies of 10 per million words and higher). For each valid word, we explored the correlation between word frequency and the proportion of hard mode reports and found a small correlation **0.08641802**, which shows that the word frequency does not significantly affect the percentage of scores reported that were played in Hard Mode.

For the second attribute, the number of unique letters in a word may affect the result of the game since repetition may add to the difficulty by leading to more letters to choose from next time. For this attribute, we define a function called *uniqueness*, which shows how many unique letters in a word, e.g. $\text{uniqueness}(\text{happy}) = 4$, while $\text{uniqueness}(\text{tepid}) = 5$. Then we explored the corre-

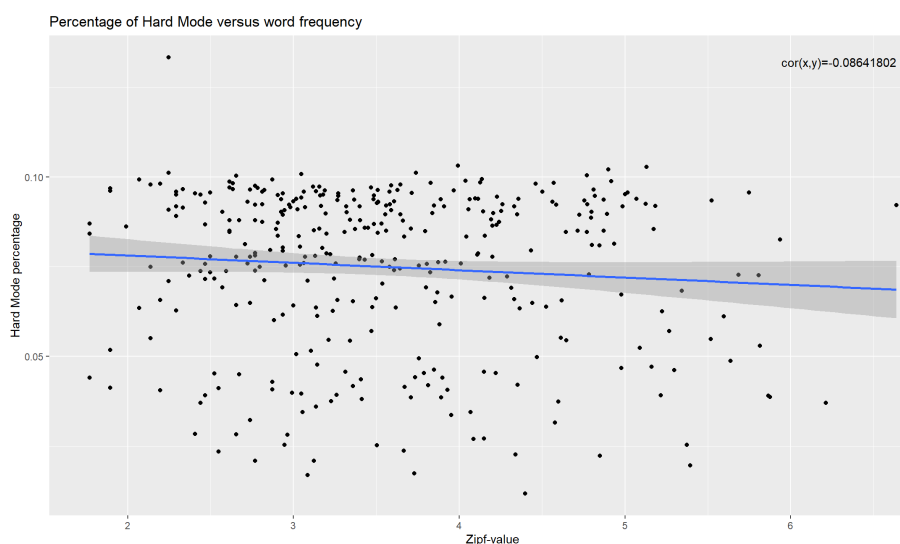


Figure 3: percentage of Hard Mode versus word frequency

lation between word uniqueness and the proportion of hard mode reports and found a correlation **0.07373057**, which shows that the number of unique letters in a word may not have significant effect on the percentage of scores reported that were played in Hard Mode.

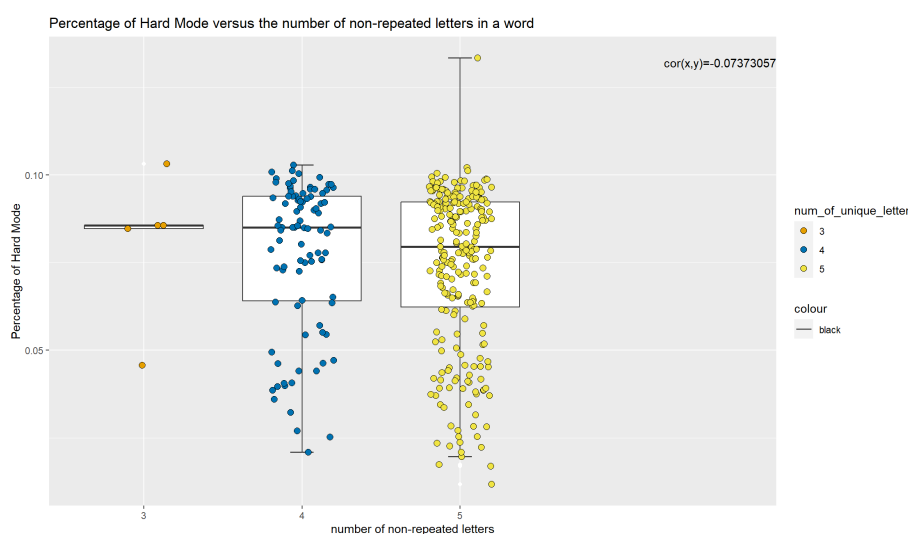


Figure 4: percentage of Hard Mode versus the number of non-repeated letters in a word

For the third attribute, we calculate one word's similarity to another one by function `num_same()`, the result is the number of the letter at same position by other words by the sum of the number of the same letter at same position. For example, in a word set { happy, lucky, puppy}, the similarity between "happy" and "lucky" is $2^1 = 1$, and the similarity between "happy" and "puppy" is $2^3 = 8$, then in this set, "happy"'s similarity to other words is $1 + 8 = 9$. Build a letter set with all the five-letter words and calculate all words' similarity, then explored the correlation between word uniqueness and the proportion of hard mode reports and found a correlation **0.04626648**, which shows that the similarity between a word and other words may not have significant effect on the

percentage of scores reported that were played in Hard Mode.

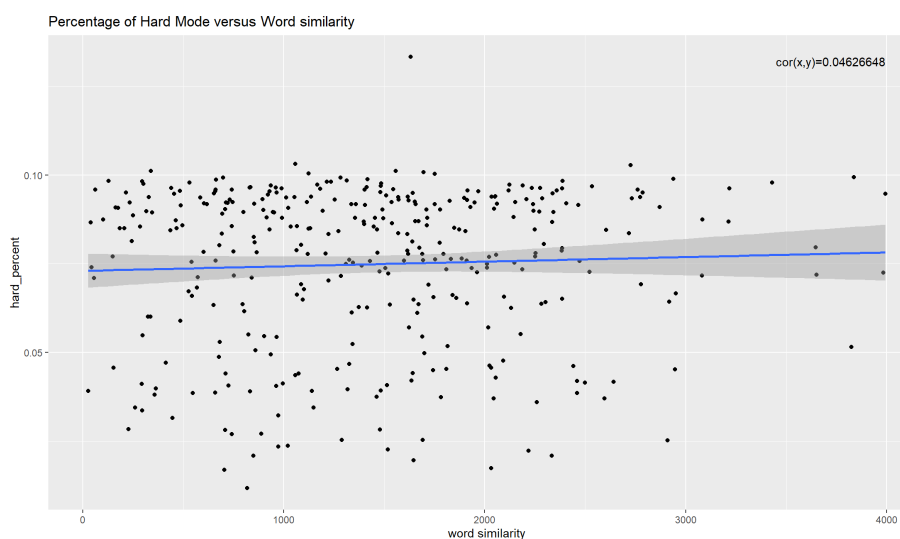


Figure 5: percentage of Hard Mode versus word similarity

3.1.4 Results of problem 1.2

In conclusion, the correlation measuring results are as follows. It shows that there is **no attribute of the word** affect the percentage of scores reported that were played in Hard Mode.

Content	Correlation coefficient
percentage of hard mode - word frequency	0.08641802
percentage of hard mode - number of unique letters	0.07373057
percentage of hard mode - similarity	0.04626648

In summary, we did not find any attributes of the word which significantly affect this percentage. We think that this is reasonable, because we believe that the percentage of scores reported that were played in Hard Mode is proportional to the number of players who play the game in Hard Mode on that day. However, when choosing whether to play the game in Hard Mode or not, the players do not know anything about the word on that day, so the attributes of the word have no reason to affect this choice.

3.2 Problem 2

3.2.1 Multiple linear regression forecasting model for distribution

First, we calculated the relationship between the distribution and time from one to six times, and performed a linear regression analysis of a single variable (contest number) on a multivariate

(a matrix composed of a "distribution" distribution). In the table below, V2 refers to similarity. z refers to word frequenc. num_unique refer to number of unique letters in the word. time refers to contest number.

one	Estimate	Std. Error	t value	Pr(> t)
V2	-2.735e-05	4.924e-05	-0.555	0.57903
z	2.074e-01	4.165e-02	4.981	1.00e-06 ***
num_unique	3.240e-01	8.018e-02	4.041	6.56e-05 ***
time	-1.634e-03	3.766e-04	-4.339	1.89e-05 ***

two	Estimate	Std. Error	t value	Pr(> t)
V2	1.471e-04	2.379e-04	0.618	0.537
z	1.341e+00	2.012e-01	6.663	1.07e-10 ***
num_unique	2.794e+00	3.874e-01	7.213	3.53e-12 ***
time	4.376e-04	1.820e-03	0.240	0.810

three	Estimate	Std. Error	t value	Pr(> t)
V2	-4.685e-04	4.563e-04	-1.027	0.3053
z	1.916e+00	3.859e-01	4.966	1.08e-06 ***
num_unique	6.519e+00	7.430e-01	8.774	< 2e-16 ***
time	7.685e-03	3.490e-03	2.202	0.0283 *

four	Estimate	Std. Error	t value	Pr(> t)
V2	-0.0008076	0.0003388	-2.384	0.0177 *
z	-0.5329498	0.2865252	-1.860	0.0637
num_unique	0.9517038	0.5516177	1.725	0.0854
time	0.0114989	0.0025913	4.438	1.23e-05 ***

five	Estimate	Std. Error	t value	Pr(> t)
V2	2.455e-06	3.488e-04	0.007	0.9944
z	-1.716e+00	2.950e-01	-5.816	1.38e-08 ***
num_unique	-4.933e+00	5.679e-01	-8.687	< 2e-16 ***
time	-4.470e-03	2.668e-03	-1.676	0.0947

six	Estimate	Std. Error	t value	Pr(> t)
V2	0.0005504	0.0003745	1.470	0.142555
z	-1.1158347	0.3167114	-3.523	0.000484 ***
num_unique	-4.2614406	0.6097321	-6.989	1.45e-11 ***
time	-0.0118930	0.0028643	-4.152	4.16e-05 ***

more	Estimate	Std. Error	t value	Pr(> t)
V2	0.0004515	0.0002213	2.040	0.042133 *
z	-0.1589039	0.1871877	-0.849	0.396528
num_unique	-1.3096702	0.3603733	-3.634	0.000322 ***
time	-0.0031946	0.0016929	-1.887	0.059994

It was found that time and Pr(>|t|) of 1, 4, and 6 have a larger Signif., indicating that time is indeed one of the factors affecting the distribution. Time has also become a factor like, similarity, the number of unique letters in the word, and word frequency, so add time to the variables of the multiple linear regression model.

First, perform basic processing on the data, such as removing the data with a large offset and removing words whose length is not 5. After word processing, 348 valid words remain from 359 words.

The similarity, the number of unique letters in the word, word frequency, and time are used as independent variables, and the data frame composed of `cbind(one,two,three,four,five,six,more)` is used as dependent variables, and `lm()` linear regression and prediction are performed.

Then we calculating 95% confidence intervals for multiple linear regression models.

3.2.2 Results of problem 2

The linear regression results are as follows:

Attempts	1	2	3	4	5	6	7
Coefficient	-0.267	2.877	19.173	36.334	27.266	11.854	2.483

The result of solving for the word "eerie" is as follows:

Attempts	1	2	3	4	5	6	7
Percentage	0	3	19	36	27	12	2

3.2.3 Validation

In order to verify the reliability of Model 2, we plotted the residuals of multiple linear regression.

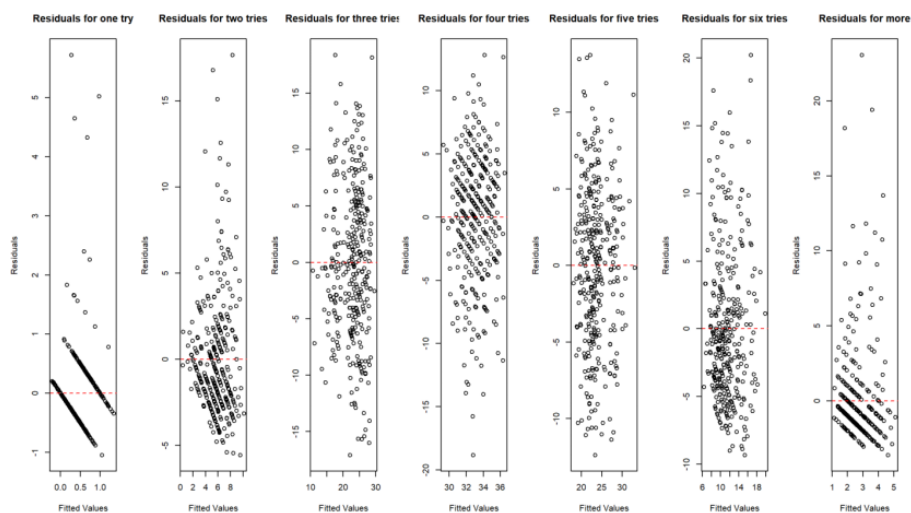


Figure 6: the residuals of multiple linear regression

According to the residual plot, our model has some systematic error and is not accurate enough.

3.3 Problem 3

3.3.1 Multiple Linear Regression Classification Model

Similarly, we preprocess the word data similar to the second question, and thus build a difficulty classification model. Arrange the processed data according to the average number of rounds, and divide the data into five equal parts. Specifically, for 348 rounds, the classification method of 69, 70, 70, 70, 69 can be used to obtain the numerical classification of difficulty.

After eliminating the time variable with the second question multiple linear regression model, insert the word to be predicted, calculate the distribution of the number of attempts, and test its validity. Then use this distribution data to obtain the average number of rounds, put it into the numerical classification model of difficulty, and obtain the difficulty of the word.

To identify the attributes of a given word that are associated with each classification. We solve the correlation between the hard classification and the similarity between a word and other words, word frequency, and the number of unique letters in the word, and finally get the result.

3.3.2 Results of problem 3

Establish the average number of rounds $E(x)$ to measure the difficulty, replace “X” by 7, arrange the processed data according to the average number of rounds, and divide the data into five equal parts, the results are as follows.

Difficulty	1	2	3	4	5
Interval	3.10-3.89	3.90-4.06	4.07-4.25	4.26-4.50	4.51-5.84

Use the multiple linear regression model of the second question, put the attributes of “EERIE” into the number of predicted rounds, and get the number of rounds as follows.

Attempts	0	1	2	3	4	5	6	more
Percentage	1	0	1	10	31	33	20	6

Calculate the average number of rounds on the distribution data of the number of attempts to obtain 4.83. According to the difficulty level, the final result is that the word “EERIE” has a difficulty of 5, which is the hardest.

To identify the attributes of a given word that are associated with each classification. We obtained the linear correlation coefficient between difficulty classification and similarity as 0.02437253, obtained the linear correlation coefficient between difficulty classification and word frequency as -0.2270678, and obtained the linear correlation coefficient between difficulty classification and the number of unique letters as -0.3817751.

Content	Correlation coefficient
hard classification - similarity	0.02437253
hard classification - word frequency	-0.2270678
hard classification - number of unique letters in the word	-0.3817751

This data shows that difficulty has little to do with similarity, but difficulty has a negative correlation with word frequency, and the number of unique letters in the word. The higher the word frequency, the simpler the word is. The more unique letters in the word, the simpler it is. The detailed figures for analyzing are shown in Figure 6, 7, 8.

4 Other Interesting Features of this Data Set

Success requires the average skewness and average kurtosis of the bureau number distribution data to be 0.3026415, -1.748139, respectively. For skewness, $0.3026415 > 0$, indicating that the right

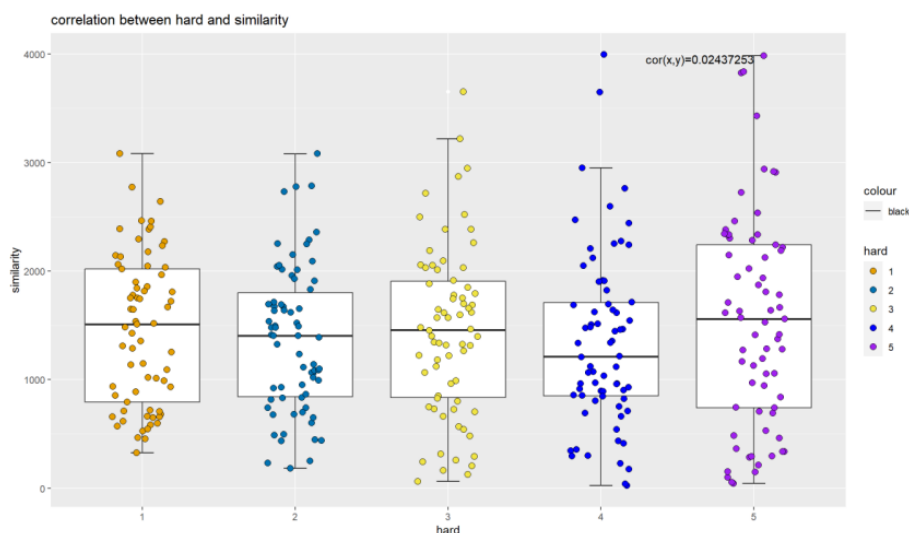


Figure 7: correlation between hard and similarity

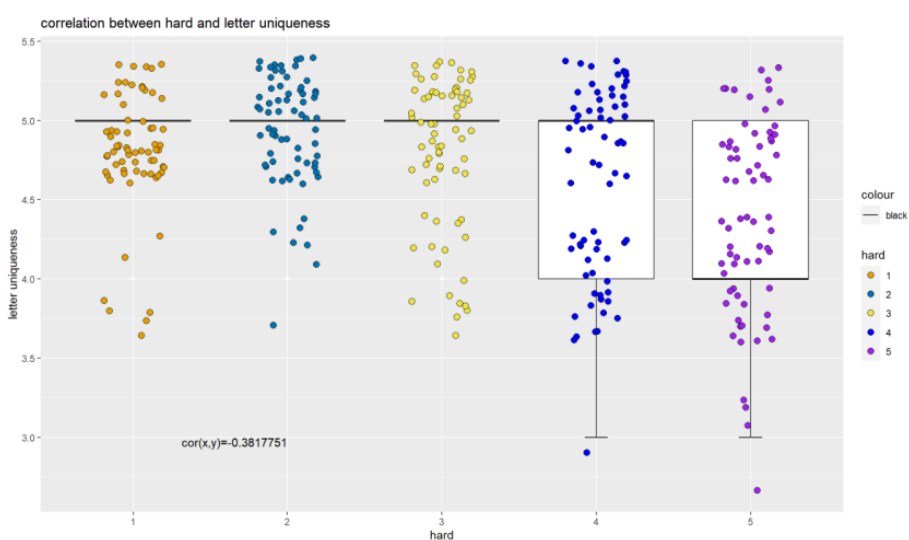


Figure 8: correlation between hard and letters uniqueness

tail is longer than the left, the data distribution is skewed to the right, and the mean is biased to the right of the data. Success requires that the game number distribution data be a positively skewed data distribution. For kurtosis, $-1.748139 < 0$, indicating that the kurtosis of the data is lower than that of the normal distribution, that is, the distribution of the data is relatively flat. This information can help us understand the distribution characteristics of the game winning rate data, and then adjust and optimize the winning rate prediction model. These are the interesting aspects of the data.

5 Summary Letter

Dear Puzzle Editor,

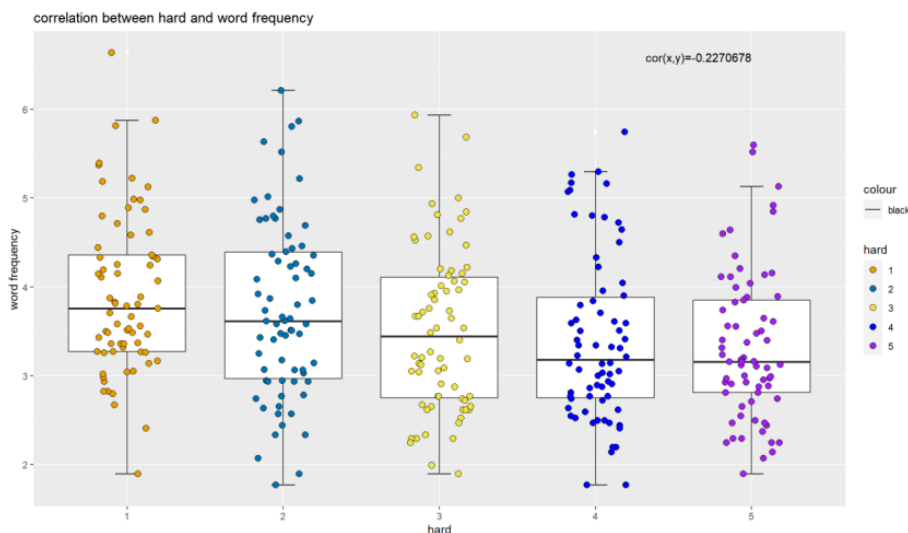


Figure 9: correlation between hard and word frequency

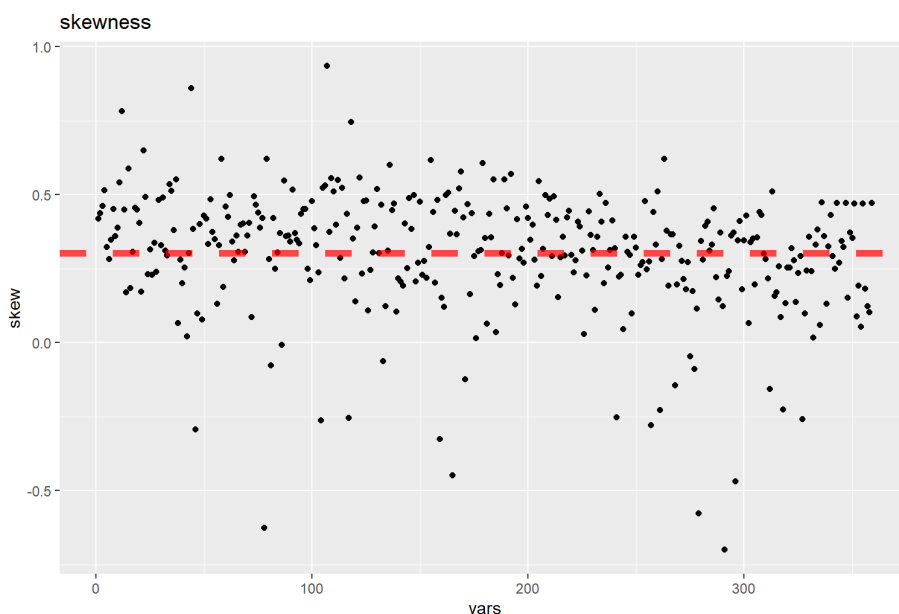


Figure 10: skewness

According to your requirements, here are the results given by us.

Firstly, you asked us to create a prediction interval for the number of reported results on March 1, 2023. We developed an exponential regression model of the form

$$y = e^{a \cdot x + b} + c$$

to explain the daily variation of the decreasing period of the number of reported results. We used the data you gaved us to train the parameters. The estimated value of each parameters are as follows:

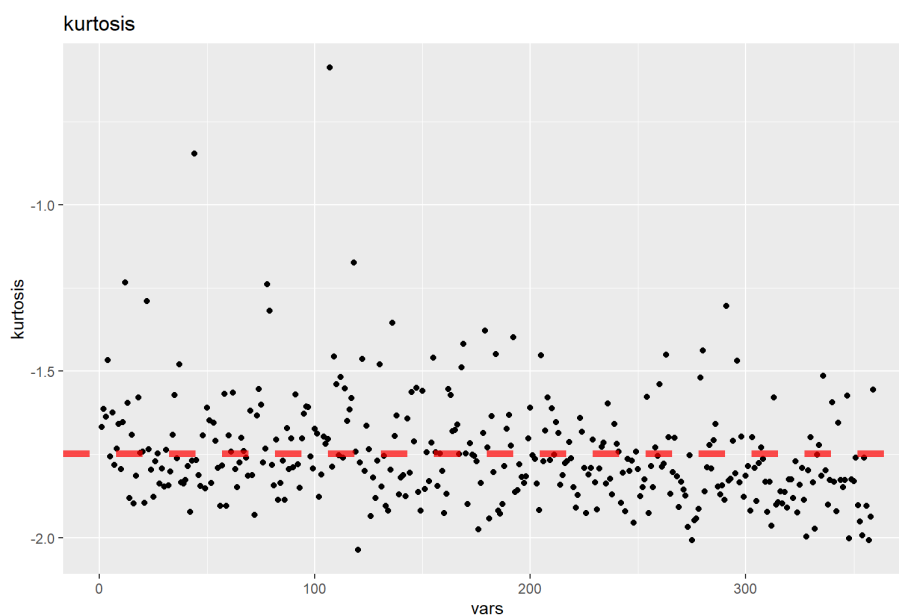


Figure 11: kurtosis

	Estimate	Std. Error	t value	Pr(> t)
a	-1.570e-02	5.356e-04	-29.32	<2e-16 ***
b	1.601e+01	1.722e-01	92.95	<2e-16 ***
c	2.311e+04	5.438e+02	42.49	<2e-16 ***

After this, we calculated the sample mean and degrees of freedom, and finally got the confidence interval. The predicted interval for the number of the reported scores on March 1, 2023 is **[17302, 29975]**, which is a 95% confidence interval. The middle number of it is **23639**, which is the expected value.

Then you asked us to find some attributes of the word affect the percentage of scores reported that were played in Hard Mode. We examined three attributes, which are the word frequency, the number of unique letters in the word, and the similarity of the word and other words. However, we find that all these three attributes have little correlation with the percentage of scores reported that were played in Hard Mode. Therefore, basically, we did not find any attributes of the word which significantly affect this percentage. We think that this is reasonable, because we believe that the percentage of scores reported that were played in Hard Mode is proportional to the number of players who play the game in Hard Mode on that day. However, when choosing whether to play the game in Hard Mode or not, the players do not know anything about the word on that day, so the attributes of the word have no reason to affect this choice.

After that, you asked us to develop a model to predict the distribution of the reported results. We developed a multiple linear regression model to predict the distribution. The variable we took into account are the word frequency, the number of unique letters in the word, the similarity of the

word and other words, and the date of the word being used for puzzle. The result of for the word “EERIE” on March 1, 2023 is as follows:

Attempts	1	2	3	4	5	6	X
Percentage	0	3	19	36	27	12	2

Then as you asked, we developed a model to classify solution words by difficulty. We arrange the processed data according to the average number of rounds, and divide the data into five equal parts. Specifically, for 348 rounds, the classification method of 69, 70, 70, 70, 69 can be used to obtain the numerical classification of difficulty. Establish the average number of rounds $E(x)$ to measure the difficulty, replace “X” by 7, arrange the processed data according to the average number of rounds, and divide the data into five equal parts, the results are as follows.

Difficulty	1	2	3	4	5
Interval	3.10-3.89	3.90-4.06	4.07-4.25	4.26-4.50	4.51-5.84

Calculate the average number of rounds on the distribution data of the number of attempts to obtain 4.83. According to the difficulty level, the final result is that the word “EERIE” has a difficulty of 5, which is the hardest.

For other interesting features of this data set, we found that success requires the average skewness and average kurtosis of the bureau number distribution data to be 0.3026415, -1.748139, respectively. For skewness, $0.3026415 > 0$, indicating that the right tail is longer than the left, the data distribution is skewed to the right, and the mean is biased to the right of the data. Success requires that the game number distribution data be a positively skewed data distribution. For kurtosis, $-1.748139 < 0$, indicating that the kurtosis of the data is lower than that of the normal distribution, that is, the distribution of the data is relatively flat. This information can help us understand the distribution characteristics of the game winning rate data, and then adjust and optimize the winning rate prediction model.

Above are all of our responses to your questions. Hope that they will be helpful for you.

Best,

Team 2311332

6 References

- [1] Van Heuven, W. J., Mandera, P., Keuleers, E., & Brysbaert, M. (2014). "Subtlex-UK: A new and improved word frequency database for British English". *Quarterly Journal of Experimental Psychology*, 67(6), 1176-1190.

<https://www.ugent.be/pp/experimentele-psychologie/en/research/documents/subtlexus/overview.htm>

[2] lee. (2021). 103976 [Code repository]. GitHub. <https://github.com/leez/103976>

Appendices

First appendix

Here are simulation programmes we used in our model as follow.

initial_number_data_plot.R:

Plot the exponential fit and confidence intervals from the first question, including predicted values

```

1 library(readxl)
2 library(ggplot2)
3 library(stats)
4 Problem_C_Data_Wordle <- read_excel("C:/Users/11872/Desktop/mcm_icm/2023_MCM-
   ICM_Problems/Problem_C_Data_Wordle.xlsx", skip = 1)
5 pp<-Problem_C_Data_Wordle[order(Problem_C_Data_Wordle$Date),]
6 ggplot(pp, aes(x=Date, y='Number of reported results')) +
7   geom_point() +geom_line()

```

plot_exp.R:

The initial scatter distribution chart drawing code in the first question

```

1 library(readxl)
2 library(ggplot2)
3 library(stats)
4 Problem_C_Data_Wordle <- read_excel("C:/Users/11872/Desktop/mcm_icm/2023_MCM-
   ICM_Problems/Problem_C_Data_Wordle.xlsx", skip = 1)
5 pp<-Problem_C_Data_Wordle[order(Problem_C_Data_Wordle$Date),]
6 pp2<-pp[-328,]
7 bb<-pp2[c(115:dim(pp2)[1]),]
8 x <- bb$'Contest number'
9 y <- bb$'Number of reported results'
10 df <- data.frame(x = x, y = y)
11 nlc <- nls.control(maxiter = 1000, tol = 1e-2)
12 fit <- nls(y ~ (exp(a * x + b) + c), start = list(a = -0.001, b = 13, c=2000),
   control = nlc, algorithm = "port")
13 summary(fit)
14 se <- sqrt(diag(vcov(fit)))
15 se
16 newdata <- data.frame(x = seq(min(df[,1]), 630), length.out = dim(df)[1])
17 y_hat <- predict(fit, newdata = newdata)
18 resid <- y - predict(fit)
19 s <- sd(resid)
20 mean_resid <- mean(resid)
21 se <- s * sqrt(1 + 1/length(x) + (x - mean(x))^2/sum((x - mean(x))^2))
22 free <- length(x) - 2
23 t_val <- qt(0.025, free)
24 lower <- y_hat - t_val * se
25 upper <- y_hat + t_val * se
26 newdata$y<-y_hat
27 newdata$lower <- lower
28 newdata$upper <- upper

```

```

29 p <- ggplot(df, aes(x, y)) +
30   geom_point() + geom_line() + geom_line(data = newdata, aes(x = x, y = y),
      color = "red") +
31   ggtitle("Exponential Fit") +
32   xlab("Contest number") + ylab("Number of reported results")
33 p <- p + geom_ribbon(data = newdata, aes(x = x, ymin = lower, ymax = upper),
      alpha = 0.4, size=0.5, color="purple")
34 p<-p+geom_vline(
35   aes(xintercept = 620),
36   color = "red")
37 p
38 new_observation <- data.frame(x = 620)
39 y_predict <- predict(fit, newdata = new_observation)
40 y_predict
41 lower_predict <- y_predict - t_val * se
42 upper_predict <- y_predict + t_val * se
43 lower_predict <- mean(lower_predict)
44 upper_predict <- mean(upper_predict)
45 p <- p + geom_point(aes(x = 620, y = upper_predict), color = "red", size = 3)
46 p <- p + geom_point(aes(x = 620, y = lower_predict), color = "red", size = 3)
47 p <- p + geom_point(aes(x = 620, y = y_predict), color = "black", size = 2)
48 p<-p+geom_text(aes(x = 620, y = upper_predict-2000, label = upper_predict),
      size = 2, color = "black")+
49   geom_text(aes(x = 620, y = lower_predict+2000, label = lower_predict), size
      = 2, color = "black")+
50   geom_text(aes(x = 620, y = y_predict+2000, label = y_predict), size = 2,
      color = "black")
51 p<-p+annotate("text", x = 620, y = max(df$y)-10000, label = "x=620", vjust =
      -0.5)
52 p

```

correlation.R:

Draw a correlation diagram for the second question

```

1 library(readxl)
2 library(ggplot2)
3 SUBTLEX_US_frequency_list_with_PoS_and_Zipf_information <- read_excel("C:/
  Users/11872/Desktop/mcm_icm/SUBTLEX-US frequency list with PoS and Zipf
  information.xlsx")
4 frequency_word=SUBTLEX_US_frequency_list_with_PoS_and_Zipf_information
5 Problem_C_Data_Wordle <- read_excel("C:/Users/11872/Desktop/mcm_icm/2023_MCM-
  ICM_Problems/Problem_C_Data_Wordle.xlsx", skip = 1)
6 table(Problem_C_Data_Wordle$Word %in% frequency_word$Word)
7 Problem_C_Data_Wordle$Word[!Problem_C_Data_Wordle$Word %in% frequency_word$
  Word]
8 valid_data<-Problem_C_Data_Wordle[Problem_C_Data_Wordle$Word %in% frequency_
  word$Word,]
9 table(frequency_word$Word %in% Problem_C_Data_Wordle$Word)
10 valid_frequency<-frequency_word[frequency_word$Word %in% Problem_C_Data_Wordle
  $Word,]
11 x
12 df<-merge(valid_data, valid_frequency, by.x = "Word", by.y = "Word", all.x =
  TRUE)
13 df<-subset(df, 'Number in hard mode'/'Number of reported results'<0.7)
14 x=df$'Zipf-value'
15 y=df$'Number in hard mode'/'Number of reported results'
16 ggplot(df, aes(x='Zipf-value', y='Number in hard mode'/'Number of reported
  results'))+geom_point()+labs(title="Percentage of Hard Mode versus word
  frequency", y="Hard Mode percentage")+geom_smooth(method = "lm")+annotate("
  text", x = max(df$'Zipf-value'), y = max(y), label = "cor(x,y)=-0.08641802"
  , hjust = 1, vjust = 1)
17 cor(x,y)
18 cor.test(x,y)

```

```

19 length=nchar(Problem_C_Data_Wordle$Word)
20 table(length)
21 ggplot(df)+geom_point(aes(x='Zipf-value',y='7 or more tries (X)'),color="red")
22 cor(df$'Zipf-value',df$'7 or more tries (X)')
23 cor(df$'Zipf-value',df$'1 try')
24 cor(df$'Zipf-value',df$'2 tries')
25 cor(df$'Zipf-value',df$'3 tries')
26 cor(df$'Zipf-value',df$'4 tries')
27 cor(df$'Zipf-value',df$'5 tries')
28 cor(df$'Zipf-value',df$'6 tries')
29 df$new1=df$'1 try'/(100-df$'7 or more tries (X)')
30 df$new2=df$'2 tries'/(100-df$'7 or more tries (X)')
31 df$new3=df$'3 tries'/(100-df$'7 or more tries (X)')
32 df$new4=df$'4 tries'/(100-df$'7 or more tries (X)')
33 df$new5=df$'5 tries'/(100-df$'7 or more tries (X)')
34 df$new6=df$'6 tries'/(100-df$'7 or more tries (X)')
35 df$average<-(df$new1*1+df$new2*2+df$new3*3+df$new4*4+df$new5*5+df$new6*6)
36 ggplot(df,aes(x='Zipf-value',y='average'))+geom_point()+labs(title="Average
  tries of winners versus word frequency")+geom_smooth(method="lm")+
  annotate("text",x=max(df$'Zipf-value'),y=max(df$average),label="
  cor(x,y)=-0.08473747",hjust=1,vjust=1)
37 cov(df$'Zipf-value',df$average)
38 ggplot(df,aes(x='Zipf-value',y='7 or more tries (X)'))+geom_point()+labs(title
  ="7 or more tries (X)~percentage versus word frequency")+geom_smooth(
  method="lm")+annotate("text",x=max(df$'Zipf-value'),y=max(df$'7 or
  more tries (X)'),label="cor(x,y)=-0.09161143",hjust=1,vjust=1)
39 cov(df$'Zipf-value',df$'7 or more tries (X)')
40 df<-df[nchar(df$Word)==5,]
41 a<-lapply(strsplit(df$Word,""),unique)
42 b<-lapply(a,length)
43 df$uniqueness<-as.vector(unlist(b))
44 ggplot(df,aes(x=uniqueness,y=average))+geom_point()+labs(title="Average
  tries of winners versus the number of non-repeated letters in a word",x="
  number of non-repeated letters")
45 cor(df$uniqueness,df$average)
46 cor(df$uniqueness,df$'7 or more tries (X)')
47 df$num_of_unique_letter<-as.factor(df$uniqueness)
48 stat_boxplot(geom="errorbar",width=0.15,aes(color="black"))+geom_boxplot(
  size=0.5,fill="white",outlier.fill="white",outlier.color="white")+geom_
  jitter(aes(fill=num_of_unique_letter),width=0.2,shape=21,size=2.5)+
  scale_fill_manual(values=c("#E69F00","#0072B2","#F0E442"))+scale_
  color_manual(values=c("black","black","black"))+labs(title="Average tries
  of winners versus the number of non-repeated letters in a word",x="
  number of non-repeated letters")+
49 annotate("text",x=max(df$uniqueness),y=max(df$average),label="cor
  (x,y)=-0.3896368",hjust=1,vjust=1)
50 df$hard_percent<-df$'Number in hard mode'/df$'Number of reported results'
51 ggplot(df,aes(x=num_of_unique_letter,y=hard_percent,fill=num_of_unique_letter)
  )+stat_boxplot(geom="errorbar",width=0.15,aes(color="black"))+geom_
  boxplot(size=0.5,fill="white",outlier.fill="white",outlier.color="white")+
  geom_jitter(aes(fill=num_of_unique_letter),width=0.2,shape=21,size=2.5)+
  scale_fill_manual(values=c("#E69F00","#0072B2","#F0E442"))+scale_color
  _manual(values=c("black","black","black"))+labs(title="Percentage of Hard
  Mode versus the number of non-repeated letters in a word",x="number of non
  -repeated letters",y="Percentage of Hard Mode")+annotate("text",x=max(df
  $uniqueness),y=max(df$hard_percent),label="cor(x,y)=-0.07373057",
  hjust=1,vjust=1)
52 cor(df$uniqueness,df$hard_percent)
53 ggplot(df,aes(x=num_of_unique_letter,y='7 or more tries (X)',fill=num_of_
  unique_letter))+stat_boxplot(geom="errorbar",width=0.15,aes(color="black
  "))+geom_boxplot(size=0.5,fill="white",outlier.fill="white",outlier.color=
  "white")+geom_jitter(aes(fill=num_of_unique_letter),width=0.2,shape=21,
  size=2.5)+scale_fill_manual(values=c("#E69F00","#0072B2","#F0E442"))+
  scale_color_manual(values=c("black","black","black"))+labs(title="7 or

```

```

more tries (X) ~ percentage versus the number of non-repeated letters in a
word ",x="number of non-repeated letters",y="7 or more tries (X)")+annotate
("text", x = max(df$ununiqueness), y = max(df$`7 or more tries (X)`), label
= "cor(x,y)=-0.1800359", hjust = 1, vjust = 1)
54 cor(df$ununiqueness,df$`7 or more tries (X)`)

```

similarity_correlation.R:

Draw similarity and hard mode percentages, the average number of rounds of winners, and the percentage of losers

```

1 output_1 <- read.csv("C:/Users/11872/Desktop/mcm_icm/output_1.txt", header=
FALSE)
2 Problem_C_Data_Wordle <- read_excel("C:/Users/11872/Desktop/mcm_icm/2023_MCM-
ICM_Problems/Problem_C_Data_Wordle.xlsx", skip = 1)
3 Problem_C_Data_Wordle$Word[!Problem_C_Data_Wordle$Word %in% output_1$V1]
4 valid_data<-Problem_C_Data_Wordle[Problem_C_Data_Wordle$Word %in% output_1$V1
,]
5 table(output_1$V1 %in% Problem_C_Data_Wordle$Word)
6 valid_repetition<-output_1[output_1$V1 %in% valid_data$Word,]
7 df<-merge(valid_data,valid_repetition, by.x = "Word", by.y = "V1", all.x =
TRUE)
8 dim(df)
9 df<-subset(df,`Number in hard mode`/`Number of reported results`<0.7)
10 x=df$V2
11 df$hard_percent<-df$`Number in hard mode`/df$`Number of reported results`
12 y=df$`Number in hard mode`/df$`Number of reported results`
13 ggplot(df, aes(x=V2,y=hard_percent))+geom_point()+labs(title="Percentage of
Hard Mode versus Word similarity",x="word similarity")+geom_smooth(method =
"lm")+annotate("text", x = max(df$V2), y = max(df$hard_percent), label = "
cor(x,y)=0.04626648", hjust = 1, vjust = 1)
14 cor(x,y)
15 df$new1=df$`1 try`/(100-df$`7 or more tries (X)`)
16 df$new2=df$`2 tries`/(100-df$`7 or more tries (X)`)
17 df$new3=df$`3 tries`/(100-df$`7 or more tries (X)`)
18 df$new4=df$`4 tries`/(100-df$`7 or more tries (X)`)
19 df$new5=df$`5 tries`/(100-df$`7 or more tries (X)`)
20 df$new6=df$`6 tries`/(100-df$`7 or more tries (X)`)
21 df$average<-(df$new1*1+df$new2*2+df$new3*3+df$new4*4+df$new5*5+df$new6*6)
22 ggplot(df, aes(x=`V2`,y=`average`))+geom_point()+labs(title="Average tries of
winners versus Word similarity",x="word similarity")+geom_smooth(method = "
lm")+annotate("text", x = max(df$V2), y = max(df$average), label = "cor(x,y)
)=0.03651125", hjust = 1, vjust = 1)
23 cor(df$`V2`,df$average)
24 ggplot(df, aes(x=`V2`,y=`7 or more tries (X)`) )+geom_point()+labs(title="`7 or
more tries (X) ~ versus Word similarity",x="word similarity")+geom_smooth(
method = "lm")+annotate("text", x = max(df$V2), y = max(df$`7 or more tries
(X)`), label = "cor(x,y)=0.1114209", hjust = 1, vjust = 1)
25 cor(df$`V2`,df$`7 or more tries (X)`)

```

distribution_factor_correlation.R:

The correlation between the distribution and the three factors (common frequency, uniqueness of word letters, similarity)

```

1 library(readxl)
2 library(ggplot2)
3 SUBTLEX_US_frequency_list_with_PoS_and_Zipf_information <- read_excel("C:/
Users/11872/Desktop/mcm_icm/SUBTLEX-US frequency list with PoS and Zipf
information.xlsx")
4 frequency_word=SUBTLEX_US_frequency_list_with_PoS_and_Zipf_information
5 output_1 <- read.csv("C:/Users/11872/Desktop/mcm_icm/output_1.txt", header=
FALSE)

```

```

6 | Problem_C_Data_Wordle <- read_excel("C:/Users/11872/Desktop/mcm_icm/2023_MCM-
  |   ICM_Problems/Problem_C_Data_Wordle.xlsx", skip = 1)
7 | 'trite' %in% frequency_word$Word
8 | 'trite' %in% output_1$V1
9 | frequency_word$`Zipf-value`[frequency_word$Word=="trite"]
10 | output_1$V2[output_1$V1=="trite"]
11 | valid_data<-Problem_C_Data_Wordle[Problem_C_Data_Wordle$Word %in% output_1$V1
  |   ,]
12 | valid_similarity<-output_1[output_1$V1 %in% valid_data$Word,]
13 | df<-merge(valid_data,valid_similarity , by.x = "Word", by.y = "V1", all.x =
  |   TRUE)
14 | valid_data_2<-df[df$Word %in% frequency_word$Word,]
15 | valid_frequency<-frequency_word[frequency_word$Word %in% df$Word,]
16 | df<-merge(valid_data_2,valid_frequency , by.x = "Word", by.y = "Word", all.x =
  |   TRUE)
17 | df<-df[nchar(df$Word)==5,]
18 | a<-lapply(str_split(df$Word,""),unique)
19 | b<-lapply(a,length)
20 | df$num_unique<-as.vector(unlist(b))
21 | df$z=df$`Zipf-value`
22 | df$one<-df$`1 try`
23 | df$two<-df$`2 tries`
24 | df$three<-df$`3 tries`
25 | df$four<-df$`4 tries`
26 | df$five<-df$`5 tries`
27 | df$six<-df$`6 tries`
28 | df$more<-df$`7 or more tries (X)`
29 | modell<-lm(cbind(one,two,three,four,five,six,more)~z,data=df)
30 | summary(modell)
31 | model2<-lm(cbind(one,two,three,four,five,six,more)~num_unique,data=df)
32 | summary(model2)
33 | model3<-lm(cbind(one,two,three,four,five,six,more)~V2,data=df)
34 | summary(model3)

```

distribution_time_correlation.R:

The correlation of distribution in time

```

1 | library(readxl)
2 | library(ggplot2)
3 | Problem_C_Data_Wordle <- read_excel("C:/Users/11872/Desktop/mcm_icm/2023_MCM-
  |   ICM_Problems/Problem_C_Data_Wordle.xlsx", skip = 1)
4 | df<-Problem_C_Data_Wordle
5 | df$one<-df$`1 try`
6 | df$two<-df$`2 tries`
7 | df$three<-df$`3 tries`
8 | df$four<-df$`4 tries`
9 | df$five<-df$`5 tries`
10 | df$six<-df$`6 tries`
11 | df$more<-df$`7 or more tries (X)`
12 | df$time<-df$`Contest number`
13 | modell<-lm(cbind(one,two,three,four,five,six,more)~time,data=df)
14 | summary(modell)

```

linear_predict_model_withtime.R

Prediction results when multiple linear regression counts time as the dependent variable

```

1 | library(readxl)
2 | library(ggplot2)
3 | library(stats)
4 | SUBTLEX_US_frequency_list_with_PoS_and_Zipf_information <- read_excel("C:/
  |   Users/11872/Desktop/mcm_icm/SUBTLEX-US frequency list with PoS and Zipf

```

```

    information.xlsx")
5 frequency_word=SUBTLEX_US_frequency_list_with_PoS_and_Zipf_information
6 output_1 <- read.csv("C:/Users/11872/Desktop/mcm_icm/output_1.txt", header=
  FALSE)
7 Problem_C_Data_Wordle <- read_excel("C:/Users/11872/Desktop/mcm_icm/2023_MCM-
  ICM_Problems/Problem_C_Data_Wordle.xlsx", skip = 1)
8 'trite' %in% frequency_word$Word
9 'trite' %in% output_1$V1
10 frequency_word$`Zipf-value`[frequency_word$Word=="trite"]
11 output_1$V2[output_1$V1=="trite"]
12 valid_data<-Problem_C_Data_Wordle[Problem_C_Data_Wordle$Word %in% output_1$V1
  ,]
13 valid_similarity<-output_1[output_1$V1 %in% valid_data$Word,]
14 df<-merge(valid_data,valid_similarity, by.x = "Word", by.y = "V1", all.x =
  TRUE)
15 valid_data_2<-df[df$Word %in% frequency_word$Word,]
16 valid_frequency<-frequency_word[frequency_word$Word %in% df$Word,]
17 df<-merge(valid_data_2,valid_frequency, by.x = "Word", by.y = "Word", all.x =
  TRUE)
18 df<-df[nchar(df$Word)==5,]
19 a<-lapply(str_split(df$Word,""),unique)
20 b<-lapply(a,length)
21 df$num_unique<-as.vector(unlist(b))
22 df$z=df$`Zipf-value`
23 df$one<-df$`1 try`
24 df$two<-df$`2 tries`
25 df$three<-df$`3 tries`
26 df$four<-df$`4 tries`
27 df$five<-df$`5 tries`
28 df$six<-df$`6 tries`
29 df$more<-df$`7 or more tries (X)`
30 df$time<-df$`Contest number`
31 model1<-lm(`1 try` ~ V2+z+num_unique+time, data=df)
32 new_data<-data.frame(V2=462,z=2.904618,num_unique=4,time=620)
33 predicted_1 <- predict(model1, newdata = new_data)
34 predicted_1
35 model2<-lm(`2 tries` ~ V2+z+num_unique+time, data=df)
36 predicted_2 <- predict(model2, newdata = new_data)
37 predicted_2
38 model3<-lm(`3 tries` ~ V2+z+num_unique+time, data=df)
39 predicted_3 <- predict(model3, newdata = new_data)
40 predicted_3
41 model4<-lm(`4 tries` ~ V2+z+num_unique+time, data=df)
42 predicted_4 <- predict(model4, newdata = new_data)
43 predicted_4
44 model5<-lm(`5 tries` ~ V2+z+num_unique+time, data=df)
45 predicted_5 <- predict(model5, newdata = new_data)
46 predicted_5
47 model6<-lm(`6 tries` ~ V2+z+num_unique+time, data=df)
48 predicted_6 <- predict(model6, newdata = new_data)
49 predicted_6
50 model7<-lm(`7 or more tries (X)` ~ V2+z+num_unique+time, data=df)
51 predicted_7 <- predict(model7, newdata = new_data)
52 predicted_7
53 linear_result<-c(predicted_1,predicted_2,predicted_3,predicted_4,predicted_5,
  predicted_6,predicted_7)
54 round_linear_result<-round(linear_result)
55 linear_result
56 round_linear_result
57 sum(round_linear_result)
58 model<-lm(cbind(one,two,three,four,five,six,more) ~ V2+z+num_unique+time, data=
  df)
59 new_data<-data.frame(V2=462,z=2.904618,num_unique=4,time=620)
60 predicted <- predict(model, newdata = new_data)

```



```

61 predicted
62 result<-round(predicted)
63 result
64 summary(model)
65 df$pred <- predict(model)
66 df$resid <- df[,c("one","two","three","four","five","six","more")] - df$pred
67 a<-c("one try","two tries","three tries","four tries","five tries","six tries"
68     ,"more")
69 par(mfrow = c(1, 7))
70 for (i in 1:7) {
71   plot(df$pred[,i], residuals(model)[,i], main = paste("Residuals for", a[i]),
72         xlab = "Fitted Values", ylab = "Residuals")
73   abline(h = 0, col = "red", lty = 2)
74 }
75 pred1 <- predict.lm(model1, new_data, interval = "confidence", level = 0.95)
76 pred2 <- predict.lm(model2, new_data, interval = "confidence", level = 0.95)
77 pred3 <- predict.lm(model3, new_data, interval = "confidence", level = 0.95)
78 pred4 <- predict.lm(model4, new_data, interval = "confidence", level = 0.95)
79 pred5 <- predict.lm(model5, new_data, interval = "confidence", level = 0.95)
80 pred6 <- predict.lm(model6, new_data, interval = "confidence", level = 0.95)
81 pred7 <- predict.lm(model7, new_data, interval = "confidence", level = 0.95)
82 exact_interval<-rbind(pred1, pred2, pred3, pred4, pred5, pred6, pred7)
83 exact_interval
84 round(exact_interval)

```

mcm3.R

The third question is all the codes, including drawing

```

1 library(readxl)
2 library(ggplot2)
3 library(stats)
4 SUBTLEX_US_frequency_list_with_PoS_and_Zipf_information <- read_excel("C:/
5 Users/11872/Desktop/mcm_icm/SUBTLEX-US frequency list with PoS and Zipf
6 information.xlsx")
7 frequency_word=SUBTLEX_US_frequency_list_with_PoS_and_Zipf_information
8 output_1 <- read.csv("C:/Users/11872/Desktop/mcm_icm/output_1.txt", header=
9 FALSE)
10 Problem_C_Data_Wordle <- read_excel("C:/Users/11872/Desktop/mcm_icm/2023_MCM-
11 ICM_Problems/Problem_C_Data_Wordle.xlsx", skip = 1)
12 'trite' %in% frequency_word$Word
13 'trite' %in% output_1$V1
14 frequency_word$`Zipf-value`[frequency_word$Word=="trite"]
15 output_1$V2[output_1$V1=="trite"]
16 frequency_word[frequency_word$Word=="eerie",]$`Zipf-value`
17 valid_data<-Problem_C_Data_Wordle[Problem_C_Data_Wordle$Word %in% output_1$V1
18 ,]
19 valid_similarity<-output_1[output_1$V1 %in% valid_data$Word,]
20 df<-merge(valid_data, valid_similarity, by.x = "Word", by.y = "V1", all.x =
21 TRUE)
22 valid_data_2<-df[df$Word %in% frequency_word$Word,]
23 valid_frequency<-frequency_word[frequency_word$Word %in% df$Word,]
24 df<-merge(valid_data_2, valid_frequency, by.x = "Word", by.y = "Word", all.x =
25 TRUE)
26 df<-df[nchar(df$Word)==5,]
27 a<-lapply(strsplit(df$Word, ""), unique)
28 b<-lapply(a, length)
29 df$num_unique<-as.vector(unlist(b))
30 df$z=df$`Zipf-value`
31 df$one<-df$`1 try`
32 df$two<-df$`2 tries`
33 df$three<-df$`3 tries`
34 df$four<-df$`4 tries`
35 df$five<-df$`5 tries`

```

```

29 df$six<-df$'6 tries '
30 df$more<-df$'7 or more tries (X) '
31 df$time<-df$'Contest number '
32 df$ave<-(df$one*1+df$two*2+df$three*3+df$four*4+df$five*5+df$six*6+df$more*7)/
100
33 df<-df[order(df$ave),]
34 df$hard<-c(seq(1,1,length.out=69),seq(2,2,length.out=70),seq(3,3,length.out
=70),seq(4,4,length.out=70),seq(5,5,length.out=69))
35 df$hard<-as.factor(df$hard)
36 model<-lm(cbind(one,two,three,four,five,six,more) ~ V2+z+num_unique, data=df)
37 new_data<-data.frame(V2=2937,z=3.246077,num_unique=3)
38 predicted <- predict(model, newdata = new_data)
39 result<-round(predicted)
40 sum(result)
41 result
42 avg_result<-(result[1]*1+result[2]*2+result[3]*3+result[4]*4+result[5]*5+
result[6]*6+result[7]*7)/100
43 avg_result
44 df$ave[1]
45 df$ave[69]
46 df$ave[139]
47 df$ave[209]
48 df$ave[279]
49 df$ave[348]

```

generator.py

Preprocessing of mono-lexical data

```

1 import csv
2 def num_same(item1, item2):
3     res = 0
4     add = 1
5     for i in range(5):
6         if (item1[i] == item2[i]):
7             res += add
8             add += 1
9     return res
10 csvFile = open("EnWords.csv", "r")
11 reader = csv.reader(csvFile)
12 result = []
13 for item in reader:
14     if (len(item[0]) == 5):
15         result.append(item[0])
16 csvFile.close()
17 new_res = [0 for _ in range(len(result))]
18 for i in range(len(result)):
19     for j in range(i+1, len(result)):
20         new_res[i] += num_same(result[i], result[j])
21 print(len(result))
22 for i in range(len(result)):
23     print(result[i], new_res[i])

```

Second appendix

Here are some files we used in our model as follow.

SUBTLEX-US frequency list with PoS and Zipf information.xlsx

common level source data

Enwords.xlsx

Glossary data