UDC: 811.322:

## COMPUTATIONAL METHODS USING CHARACTER STATISTICS FOR THE WORD GAME

Ulugbek Salaev,

ulugbek0302@gmail.com

Urgench State University, Ph.D. student

**Abstract.** Wordle is a word game that intends six attempts to guess a five-letter word. This study uses character statistics of Uzbek five-letter words in Cyrillic script to determine the solving strategy of the word game. The proposed methodology covered computing the letter frequency (LF) and positional letter frequency (PLF) of Uzbek words in a Cyrillic script and calculated a word score based on the PLF to suggest the best probability words to obtain a solution in the game with minimal attempts.

**Keywords:** Wordle, Character statistics, word-game, Uzbek words

Introduction. The web-based word game Wordle [1] was implemented in November of 2021. This game is adopted into many languages of the world since it is open-source software. Among them, there is also the Uzbek language in both Latin and Cyrillic scripts. It gives players six attempts to try and guess a five-letter secret answer word. At each attempt, hints are offered to the player so that they may make subsequently "better" guesses by marking each letter with one of three possible colors, each has a specific meaning:

- 1) Green: The entered letter is in the expected solution and is in the expected position.
- 2) Yellow: The entered letter is in the expected solution but it is not in the expected position.
- 3) Gray: The entered letter is not in the expected solution.

Figure 3 shows examples of Wordle solutions with the tile colors providing hints to the player to progress. Also, there is no indication as to how many times the character would appear.

The work by Sidhu [3] attempted to find the best starting word from a linguistic perspective while the work by Anderson and Meyer [2] have used machine learning to find the optimal human strategy for solving the Wordle. The objective of this work is to derive the set of 3 optimum starting words for the English version of the game covering 15 different characters and ordered in the descending order of significance. Fully automatic solving software of Wordle has also been created for

the English language.

The proposed methodology in this article is targeted at the Cyrillic script version of the Uzbek language for the Wordle game. This version can be accessed via the webpage http://wordlar.uz/.

**Methodology.** The word database contains five-letter 5923 total words (W) in Cyrillic script that was manually prepared from the Uzbek dictionary [4] that has over 85k Uzbek root words. The majority portion of the dataset included root morphemes, because of the agglutinative character of the Uzbek language as many inflectional words can be formed by adding several suffixes to most root morphemes and it causes mismeasurements while calculating LF and PLF. From W, each of 3954 words has unique letters, which made up 67% of the total number of words. Table I shows the character occurrence count and its frequency, and the count of words that the character exists.

Table I: Calculated character frequencies

Letter	cnt.	freq.	*cnt.	Letter	cnt.	freq.	*cnt.	Letter	cnt.	freq.	*cnt.
a	4188	14.14	3439	M	1213	4.10	1167	Ш	714	2.41	676
б	1050	3.55	997	Н	1429	4.83	1365	Ъ	31	0.10	31
В	416	1.40	411	O	1826	6.17	1736	Ь	7	0.02	7
Γ	626	2.11	615	П	300	1.01	291	Э	114	0.38	114
Д	1070	3.61	1010	p	1400	4.73	1358	Ю	178	0.60	176
e	541	1.83	535	c	1142	3.86	1089	R	230	0.78	229
Ж	237	0.80	232	Т	1226	4.14	1166	ë	308	1.04	308
3	680	2.30	665	y	1022	3.45	926	ÿ	650	2.19	649
И	3639	12.29	3133	ф	213	0.72	210	F	318	1.07	316
й	647	2.18	632	X	289	0.98	287	Қ	929	3.14	853
К	829	2.80	763	Ц	7	0.02	7	X,	352	1.19	350
Л	1290	4.36	1229	Ч	504	1.70	488				

\*cnt – Word counts of the character exist

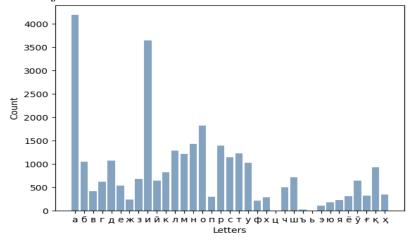


Figure I. Character frequencies

It can be seen from Table I that some letters are frequently used while some of them are very rare. Letters "a", " $\mu$ ", "o", " $\mu$ " and " $\mu$ ", " $\mu$ " occur 58%, 53%, 29%, 23%, 21% in words of W respectively, in which the 3 letters with the highest frequency correspond to the vowel letters. In the next step, we calculated the character distribution on each position (Table II) and kept it as a char frequency map.

#5 #1 #2 #3 #4 #5 #1 #2 #3 #4 #1 #2 #3 #4 #5 M Ш a б Ъ В Ь П Г Ю Д p e cЯ ë T Ж ÿ y ф F И й X Қ К Ц X Ч Л

Table II. Character count by positional

We calculated the score for words of W based on Table II, this score indicates possibilities of words that can be offered as a solution. In the first position, the letter 'a' has the highest probability in the letter frequency map, indicating that 'a' is the most common starting letter for 5-letter words. Similarly, the 5th position shows that 'a' is the most common ending letter. This distribution was calculated only on words found in the 5923 words database. Figure II illustrates the most common five letter frequencies for the position.

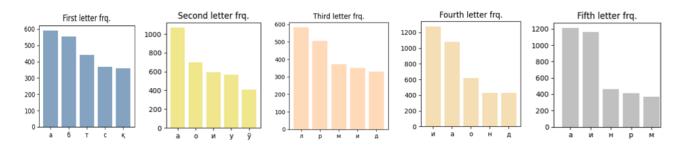


Figure II. Most common five letter frequencies for the position

To make a conclusion based on the frequency of letters in the Uzbek language, the highest frequency letters correspond to vowels. The best first guess has five different letters and it should include some of these three vowels. For this reason, it is suggested to start with vowels in the initial attempt when guessing a word. Using words with high frequency and not overlapping letters lead to a good process. We chose to focus only on the possible winning answers for this analysis because our goal is to find the best guesses for human players. Since machines can make perfect use of information from grays and yellows in choosing their next move, this method aims to optimize a Wordle strategy for a human player. Whereas with a green letter, a human can better construct their next move with the certainty that a green provides.

A detailed algorithmic description of the proposed methodology:

## Data collection:

Collected Uzbek (Cyrillic script) 5-letter words from the dictionary and manually added other words

Load data:

$$W \leftarrow \{word_1, ..., word_n\}, |W| = 5923, # List of five-letter words$$

$$C \leftarrow \{\text{"a", "6", "B", "r", ..., "x"}\}, |A| = 35 \# Alphabet, List of character$$

Calculate character frequencies:

$$\forall w \left[ \forall c \left[ F(c) = F(c) + \sum_{l \in w} \left\{ \begin{cases} 1 & \text{if } c = l \\ 0 & else \end{cases} \right] \right], \# 1 \text{ is a character in the word } w, \text{ F is }$$

character frequency map of the format <character, value>

$$\bigvee_{w \in W} \left[ \bigvee_{p \in \overline{1,5}} \left[ \bigvee_{l \in w} [P(p,l) = P(p,l) + 1] \right] \right], P \text{ is a count of the letter in the position}$$
of the format,  $l$  is a character in the word  $w$ ,  $p$  is a

position of l in the word w;

$$\bigvee_{w \in W} \left[ \bigvee_{p \in \overline{1,5}} \left[ \bigvee_{c \in C} \left[ W = \sum_{l \in w} \left\{ \begin{matrix} P(p,l) & \text{if } l = c \\ 0 & else \end{matrix} \right] \right] \right], \text{ update } W, \text{ word list map of the }$$

format <word, value> by adding the word w score, l is a character in the word w;

Sorted W by a score in descending order;

$$\hat{W} = \bigcup_{w \in W} \begin{cases} \{w\} & \text{if } len\left(\bigcup_{c \in w} \{c\}\right) = 5, \text{ the function makes } \hat{W} \text{ only contains words} \end{cases}$$

that have five unique characters by removing the words which have repeated characters;

Output data:

Output top **k** highest scored words in a format <word, value>;

*Input parameters: non\_exist\_letters, exist\_letters, positional\_letters;* 

Filter  $\hat{\mathbf{W}}$  by conditions given by input parameters;

Output top k highest scored words that matched the conditions;

Go to input parameters and continue;

Exit when a user terminates;

Experiment and Result. When we proceeded with the algorithm, best probabilities words by computed positional letter occurrences were extracted for suggestions to use in an initial attempt. At this point, we came across a problem, in which some of the words that were suggested as options were invalid words by Wordle as they are not included in its dataset. To input, we will choose one of the highest-ranking suggestions until it is accepted in Wordle. However, it is avoided to output letter overlapping words to have to cover more characters. The highest-ran 2 words set accepted by wordle was {'қалин', 'балиқ'}. Logically, these are the best words to use in the first attempt. Figure III shows an example of using one of the suggested words at the beginning of the game and the process of continuation based on the methodology. Additionally, any words from W may not be acceptable in the Wordle game.

Suggested words for initial attempt:

{'word': 'қалин', 'score': 3747} {'word': 'балиқ', 'score': 3688} {'word': 'салим', 'score': 3666} {'word': 'бадир', 'score': 3641} {'word': 'бақир', 'score': 3515}



Figure III. Example of Wordle Solution

Left 220/5923 words
{'word': 'бодом', 'score': 2568}
{'word': 'бозор', 'score': 2565}
{'word': 'тумор', 'score': 2410}

Left 8/5923 words
{'word': 'дутор', 'score': 2123}
{'word': 'дучор', 'score': 1918}
{'word': 'девор', 'score': 1759}

{'word': 'бемор', 'score': 2304} {'word': 'дугох', 'score': 1617} {'word': 'бедор', 'score': 2261} {'word': 'ёдгор', 'score': 1481}

After providing the word on the third attempt, we got the following information: the solution does not have 9 characters  $[\kappa, a, \pi, \mu, \mu, \delta, M, y, \tau]$  and characters  $[\pi, a, \pi, \mu, \mu, \delta, M, y, \tau]$  and characters  $[\pi, a, \mu, \mu, \mu, \delta, M, y, \tau]$  and characters  $[\pi, a, \mu, \mu, \mu, \delta, M, \mu, \mu, \tau]$  and characters  $[\pi, a, \mu, \mu, \mu, \mu, \tau]$  and characters  $[\pi, a, \mu, \mu, \mu, \tau]$  and characters  $[\pi, a, \tau]$  and characters  $[\pi,$ 

Conclusion. This article explored the letter frequency of Uzbek words in a Cyrillic script and calculate word scores based on the positional letter frequency among words. This methodology is aimed at offering the most matching words to win in the game with minimal attempts. According to the proposed algorithm the software was created. The complete list of words and the software is available at the address: https://github.com/UlugbekSalaev/wordle. It is a standalone program that is created in Python language and converted to an executable file.

## References

- [1] J. Wardle, "Wordle a daily word game," https://www.powerlanguage.co.uk/wordle/, accessed: 2022-02-04.
- [2] B. J. Anderson and J. G. Meyer, "Finding the optimal human strategy for wordle using maximum correct letter probabilities and reinforcement learning," arXiv preprint arXiv:2202.00557, 2022.
- [3] D. Sidhu, "Wordle the best word to start the game, according to a language researcher," <a href="https://bit.ly/3qItHsI">https://bit.ly/3qItHsI</a>, 2022, [Online; accessed 07-February-2022].
- [4] А.Мадвалиев, Э.Бегматов, "Ўзбек тилининг имло луғати", Тошкент: Akademnashr, 2013. 520 б.
- [5] M. Butterfield, "Science has determined the worst Wordle starting word," <a href="https://bit.ly/3J9FOWi">https://bit.ly/3J9FOWi</a>, 2022, [Online; accessed 07-February-2022].
- [6] The Best Wordle Starting Word Has Been Figured Out With (Computer) Science. In: GameSpot [Internet]. [cited 28 Jan 2022]. https://bit.ly/35k7EB9
- [7] Groux C. The 20 best Wordle starting words, according to science. In: Inverse [Internet]. [cited 28 Jan 2022]. https://bit.ly/3ISDcvn