

Implementing a Random Passphrase Generator using Markov Chains

Alan Matthew

Abstract

This document presents the design and implementation of a Markov Phrase Generator, a simple yet effective tool for generating passphrases using Markov chains based on a given training text. The project demonstrates the application of Markov chains in text generation by creating a transition matrix from the training data, which is then used to generate new phrases.

Contents

1	Introduction	3
2	Data Collection and Preprocessing	3
3	Problem Statement	3
4	Assumptions and Limitations	4
5	Core Concepts	4
5.1	Transition Matrix Creation	4
5.2	Markov Chain for Passphrase Generation	5
5.3	Generating the Passphrase	5
6	Conclusion	6

1 Introduction

The Markov Phrase Generator is a passphrase generation tool that leverages the statistical model of Markov chains. By analyzing a corpus of text, it creates a transition matrix that models the likelihood of sequences of words, allowing for the generation of new, semantically plausible phrases.

One application of this tool is in the context of cryptocurrency wallets, where secure and memorable passphrases are essential for authorizing transactions and recovering accounts. The Markov Phrase Generator aims to address the challenge of creating such passphrases by offering a balance between security and memorability, leveraging the linguistic patterns found in a chosen corpus of text.

Note that when using this tool, we need to specify the length of the passphrase to be generated. Furthermore, if you want to use a different training corpus, you can add your text file into the `traindata` directory and modify the `main.rs` file to use the new file. A `README.md` file is also included in the accompanying code repository, providing detailed instructions on how to run the project.

2 Data Collection and Preprocessing

The training data for the Markov Phrase Generator project was sourced from a collection of verses from Kanye West's songs, available on the Kaggle dataset titled Kanye West Verses. This dataset comprises various lines from Kanye West's discography, providing a rich corpus of contemporary language use and stylistic phrasing. The choice of this dataset was motivated by its diversity in vocabulary and syntax, enabling the generation of creative and memorable passphrases. Furthermore, with regards to security, the dataset's vast collection of verses increases the likelihood that the tool will not generate the same passphrase twice, enhancing the security of the generated passphrases.

Prior to usage, we need to preprocess the data to ensure that the Markov chain model can effectively capture the linguistic patterns in the training text. This involves cleaning and normalizing the text by removing non-alphabetic characters and converting all words to lowercase. The following Rust function, `clean_word`, accomplishes this task efficiently.

```
1  fn clean_word(word: &str) -> String {
2      word.chars()
3          .filter(|c| c.is_ascii_alphabetic())
4          .collect::<String>()
5          .to_lowercase()
6  }
```

Listing 1: Preprocessing Function

3 Problem Statement

The problem of generating secure and memorable passphrases for cryptocurrency wallets can be formally stated as a network problem in the context of Markov chains.

This problem is formalized as follows: given a network which consists of states represented by pairs of words from the training corpus, and the transitions between these states are determined by the subsequent word that follows a given pair, how can we navigate this network (or graph) to construct a sequence of words (passphrase) that is both unpredictable and semantically coherent?

4 Assumptions and Limitations

The Markov Phrase Generator operates under several assumptions:

- **Linguistic Coherence:** It assumes that a sequence of words generated based on the statistical patterns in the training corpus will be somewhat coherent and resemble natural language most of the time.
- **Security through Randomness:** The model presumes that by generating passphrases based on a probabilistic method, the outputs will be sufficiently random to ensure security against brute-force or dictionary attacks.

However, there are limitations to this approach:

- **Predictability:** The reliance on a fixed corpus, especially one with a distinctive style like Kanye West's lyrics, may introduce patterns that could potentially be exploited by sophisticated attackers.
- **Edge Cases of Passphrase Quality:** Since the passphrases are generated probabilistically, there's a risk of edge cases where the generated phrases are too predictable, potentially compromising security, or too nonsensical to be memorable.
- **Cultural Sensitivity:** The content of the training data may reflect the biases or explicit themes of the source material, which might not be appropriate for all users.

5 Core Concepts

5.1 Transition Matrix Creation

The transition matrix is a pivotal component in the Markov Phrase Generator. It records the probability of transitioning from one pair of words to another within the corpus. This matrix is constructed by parsing through each triplet of consecutive words in the training text, treating the first two words as the key and the third word as a possible continuation. The Rust function `create_transition_matrix` accomplishes this task efficiently.

```
1 fn create_transition_matrix(words: Vec<String>) -> HashMap<(String, String),  
    Vec<String>> {  
2     let mut transition_matrix = HashMap::new();  
3     for window in words.windows(3) {  
4         if let [w0, w1, w2] = &window {  
5             let entry = transition_matrix  
6                 .entry((w0.clone(), w1.clone()))  
7                 .or_insert_with(Vec::new);  
8             entry.push(w2.clone());  
9         }  
10    }  
11    transition_matrix  
12 }
```

Listing 2: Creating the Transition Matrix

This function iterates over each triplet (or window of 3 words) in the training data. For each triplet, it maps the first two words to a list of possible third words, effectively capturing the essence of a Markov chain's transition probabilities in a discrete and manageable form.

5.2 Markov Chain for Passphrase Generation

The core of passphrase generation lies in simulating a Markov chain process, where the next word is chosen based on the transition matrix's probabilities. The `markov_chain` function implements this idea by selecting the next word in the sequence with the help of the transition matrix created earlier.

```

1 fn markov_chain(
2     transition_matrix: &HashMap<(String, String), Vec<String>>,
3     length: usize,
4     w0: String,
5     w1: String,
6     w2: String,
7 ) -> Vec<String> {
8     let mut rng = rand::thread_rng();
9     let mut result = vec![w2.clone()];
10    let (mut _w0, mut w1, mut w2) = (w0, w1, w2);
11    for _ in 0..length - 1 {
12        if let Some(next_words) = transition_matrix.get(&(w1.clone(), w2.
13            clone())) {
14            let next_word = next_words[rng.gen_range(0..next_words.len())].
15                clone();
16            result.push(next_word.clone());
17            _w0 = w1;
18            w1 = w2;
19            w2 = next_word;
20        }
21    }
22    result
23 }

```

Listing 3: Markov Chain Function

This function starts with an initial seed (the last word of the starting triplet) and generates a sequence of words by randomly selecting from the list of possible continuations for the current state (the last two words). This randomness introduces variety in the output, maximizing the likelihood that each generated passphrase is unique.

5.3 Generating the Passphrase

The `generate_passphrase` function orchestrates the passphrase generation process, beginning with reading and cleaning the training data, followed by constructing the transition matrix, and finally generating the passphrase using the Markov chain.

```

1 fn generate_passphrase(words: Vec<String>, length: usize) -> String {
2     let clean_words = words
3         .into_iter()
4         .map(|w| clean_word(&w))
5         .collect::<Vec<String>>();
6     let transition_matrix = create_transition_matrix(clean_words.clone());
7     let start_index = rand::thread_rng().gen_range(0..clean_words.len() - 3);
8     let chain = markov_chain(
9         &transition_matrix,
10        length,
11        clean_words[start_index].clone(),
12        clean_words[start_index + 1].clone(),
13        clean_words[start_index + 2].clone(),
14    );
15    chain.join(" ")
16 }

```

Listing 4: Generating the Passphrase

This function showcases the end-to-end process of generating a passphrase. It starts by cleaning the input words to ensure consistency, creates a transition matrix to capture the linguistic structure of the input text, and then employs the Markov chain to generate a sequence of words that form the passphrase. The choice of starting words (seed) for the Markov chain is randomized to ensure diversity in the generated passphrases.

6 Conclusion

In conclusion, we were able to implement a Markov Phrase Generator that effectively generates passphrases based on a given training corpus. The project demonstrates the application of Markov chains in text generation, offering a practical tool for enhancing security for cryptocurrency wallets. While mindful of its assumptions and limitations, this project underscores the potential of applying natural language processing techniques in the realm of cybersecurity.

"the city what they do it my way to proud to admit" – A sample generated passphrase.