

# Short Report on Caesar Cipher Frequency Analysis

Traian Sfarghiu

March 21, 2025

## 1 Introduction

I implemented a program to encrypt and decrypt text using the **Caesar cipher** and to **automatically break** the cipher using **frequency analysis**. The system tests all possible shifts (0–25), then computes how closely the resulting decrypted text’s letter distribution matches **standard English frequencies** using three distance metrics:

- **Chi-Squared ( $\chi^2$ ) Distance**
- **Cosine Distance**
- **Euclidean Distance**

Below, I present the **testing procedure**, **results**, and **conclusions**.

## 2 Testing Procedure

### 1. Text Length Variations

- *Short Texts (10–20 characters)*: Simple sentences or phrases, sometimes lacking a variety of letters.
- *Long Texts (200+ characters)*: Paragraphs from articles or classic literature, with richer letter distribution.

### 2. Random Shifts

- Applied shifts like 5, 13, 19, or 23 to create encrypted versions of the original English texts.
- Attempted to decrypt them using each of the three distance metrics, noting whether the correct shift was recovered.

### 3. Special Text Structures

- *Highly Repetitive Text*: e.g., repeated words/letters (“AAAAAA BBBBBB...”).
- *Unusual Letter Frequencies*: texts with very few occurrences of common letters like ‘E’ or ‘T’.

### 4. Measurement of Success

- For each test, we recorded which metric found the **correct** shift.
- We noted **incorrect** guesses and possible reasons (very short length, unusual distribution, etc.).

## 3 Results & Observations

### 3.1 Short Texts (10–20 characters)

- *Chi-Squared*: Often correct, but sometimes off by 1–2 shift values when text lacked enough variety.
- *Cosine*: Performed reasonably but occasionally tied with an incorrect shift due to insufficient data.
- *Euclidean*: Similar to Chi-Squared, but also struggled with extremely short texts.

For *very short* strings, all three methods can fail if the text is not representative of typical English frequencies.

### 3.2 Long Texts (200+ characters)

- *Chi-Squared*: Most consistently found the correct shift because longer texts closely match expected English letter frequencies.
- *Cosine*: Typically correct, sometimes second-best candidate.
- *Euclidean*: Often matched Chi-Squared’s success; minor deviations appeared in unusual distributions.

Overall, for *normal-length* English texts, all methods succeeded more often, with Chi-Squared slightly in the lead.

### 3.3 Edge Cases & Anomalies

- **High Repetition / Limited Variety**: All methods struggled if a text repeated only a few letters; the distribution did not match typical English.
- **Missing Common Letters**: If a text barely used ‘E’ or ‘T’, analysis was less reliable. Chi-Squared still often did best, but sometimes over-penalized missing letters.
- **Very Short Text** (<10 characters): All metrics can fail entirely because the sample is not representative enough.

## 4 Conclusions

**Text length** and **distribution** are key for successful Caesar cipher breaking via frequency analysis.

- *Chi-Squared* is typically most reliable for standard English.
- *Euclidean* also performs well for normal texts, often matching Chi-Squared.
- *Cosine* is reasonably effective, but slightly more prone to ties or confusion with peculiar distributions.

With longer texts, all three distance metrics almost always correctly identify the shift. For very short or unusual texts, all can fail.

## 5 Example Table of Results

Text & Shift	Length	Chi-Squared	Cosine	Euclidean
HELLO WORLD (Shift 5)	11 chars	Correct (5)	2nd best guess	Off by 1
Paragraph (Shift 13)	300+ chars	Correct (13)	Correct (13)	Correct (13)
AAAA BBBB CCCC (3)	14 chars	Correct (3)	Incorrect (8)	Correct (3)

## 6 Final Thoughts

This project demonstrates how **classical ciphers** like Caesar can be broken with simple frequency analysis. While modern ciphers are far more secure, this assignment highlights the *importance of distribution* in determining a shift. Future explorations might involve other substitution ciphers or analyzing how text domain (e.g., scientific vs. casual) affects frequency distributions.