

**Міністерство освіти й науки України**  
**Львівський національний університет імені Івана Франка**  
Факультет електроніки та комп'ютерних технологій  
*з предмета: Комп'ютерна лінгвістика*

**Звіт**  
про виконання лабораторної роботи № 4  
**«Закон Гіпса для слів в окремих текстах»**

Виконав:  
Студент групи  
Фес-32с  
Бойко Кирило

Львів 2024

## Завдання

Використовуючи програму +proj6stats&plots, на лінгвістичному рівневі слів дослідити виконання закону Гіпса для обраних Вами текстів у лабораторній роботі №1. Провести дослідження для однограм ( $n = 1$ ) за різних опцій (умов).

## Хід виконання роботи

Підготував тексти з першої лабораторної роботи

*Clemencia Novela de costumbres by Fernán Caballero*

*DONALD J. TRUMP January 20, 2017*

Написав код програми Python для визначення лінійної апроксимації зі збережених результатів

```
import numpy as np
import matplotlib.pyplot as plt

file_path = "1.1.txt"
fixed_file_path = "1.1fix.txt"

with open(file_path, 'r', encoding='utf-8') as input_file, open(fixed_file_path, 'w', encoding='utf-8') as output_file:
    for line in input_file:
        output_file.write(line.replace(',', ' '))

data = np.loadtxt(fixed_file_path, delimiter='\t', skiprows=1)
length = data[:, 0] # Довжина
volume = data[:, 1] # Об'єм
delta_volume = data[:, 2] # Різниця об'ємів

# Логарифмічні перетворення
log_length = np.log(length)
log_volume = np.log(volume)

# Лінійна апроксимація залежності  $V(L) \propto L^p$ 
coefficients = np.polyfit(log_length, log_volume, 1)
slope = coefficients[0]
intercept = coefficients[1]

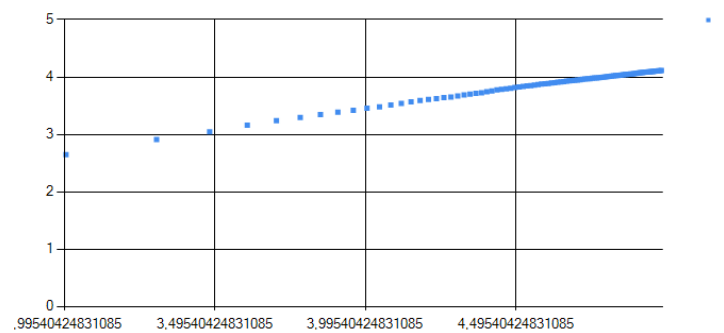
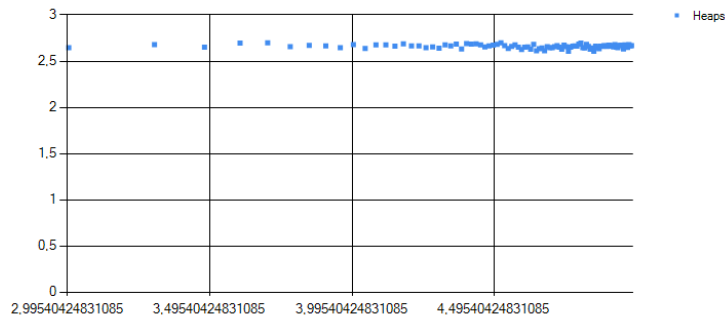
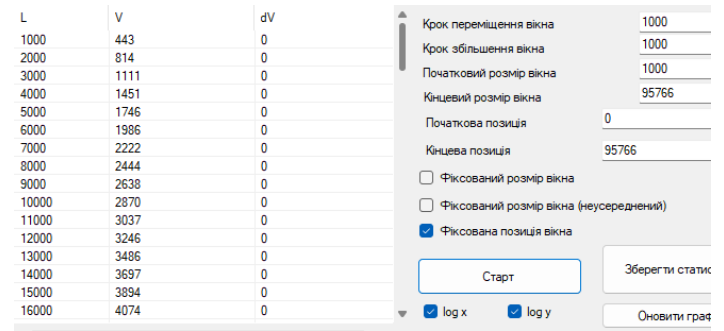
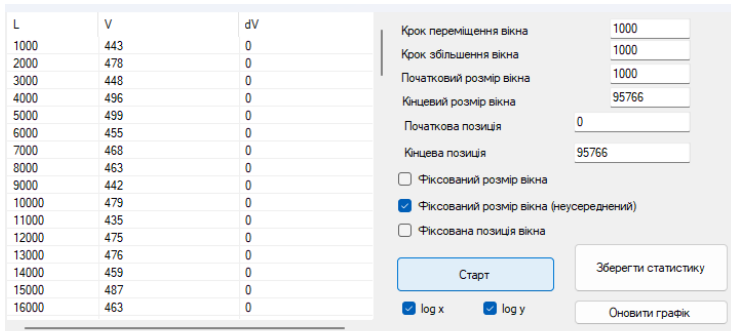
plt.figure(figsize=(10, 6))
plt.scatter(log_length, log_volume, color='green', label='Дані (логарифмічні координати)')
plt.plot(log_length, np.polyval(coefficients, log_length), color='orange', label=f"Апроксимація:  $p \approx \{slope:.2f\}")
plt.title("Лінійна апроксимація для закону Гіпса")
plt.xlabel("log(L)")
plt.ylabel("log(V)")
plt.legend()
plt.grid(True)
plt.show()

print(f"Коефіцієнт нахилу (p): {slope:.2f}")
print(f"Перетин з віссю (intercept): {intercept:.2f}")$ 
```

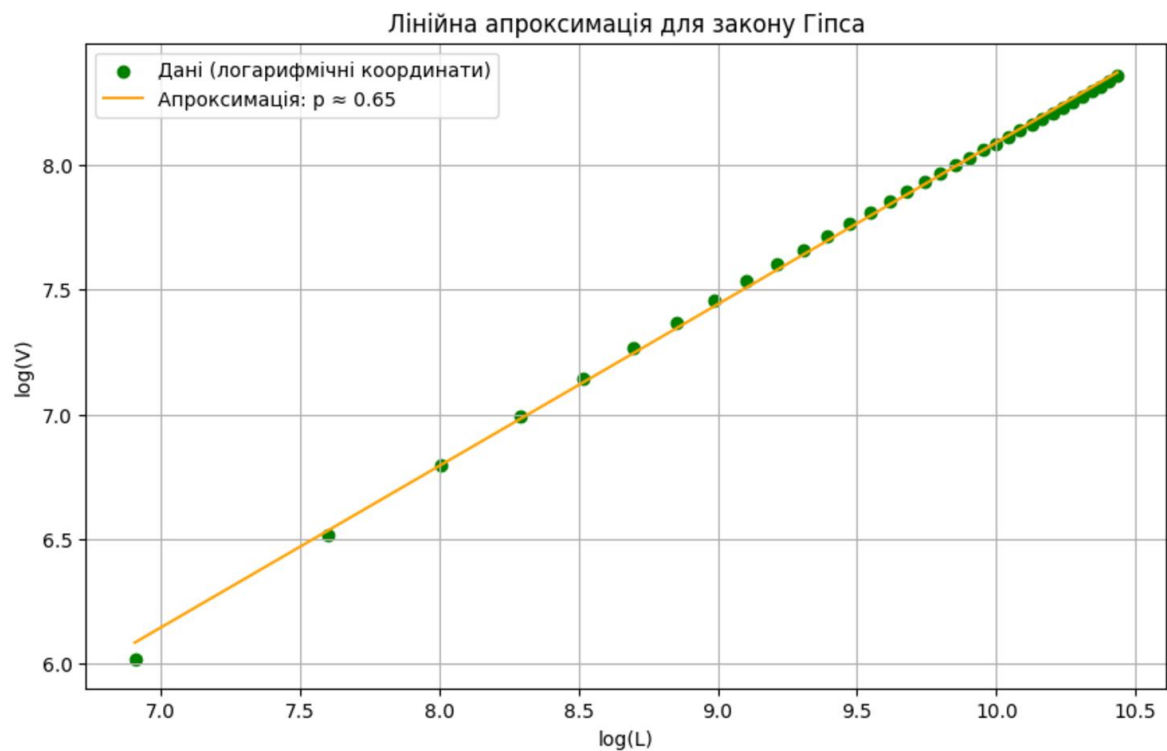
Запустив програму +proj6stats&plots і завантажив текстові файли для проведення дослідження закону Гіпса

## Результати



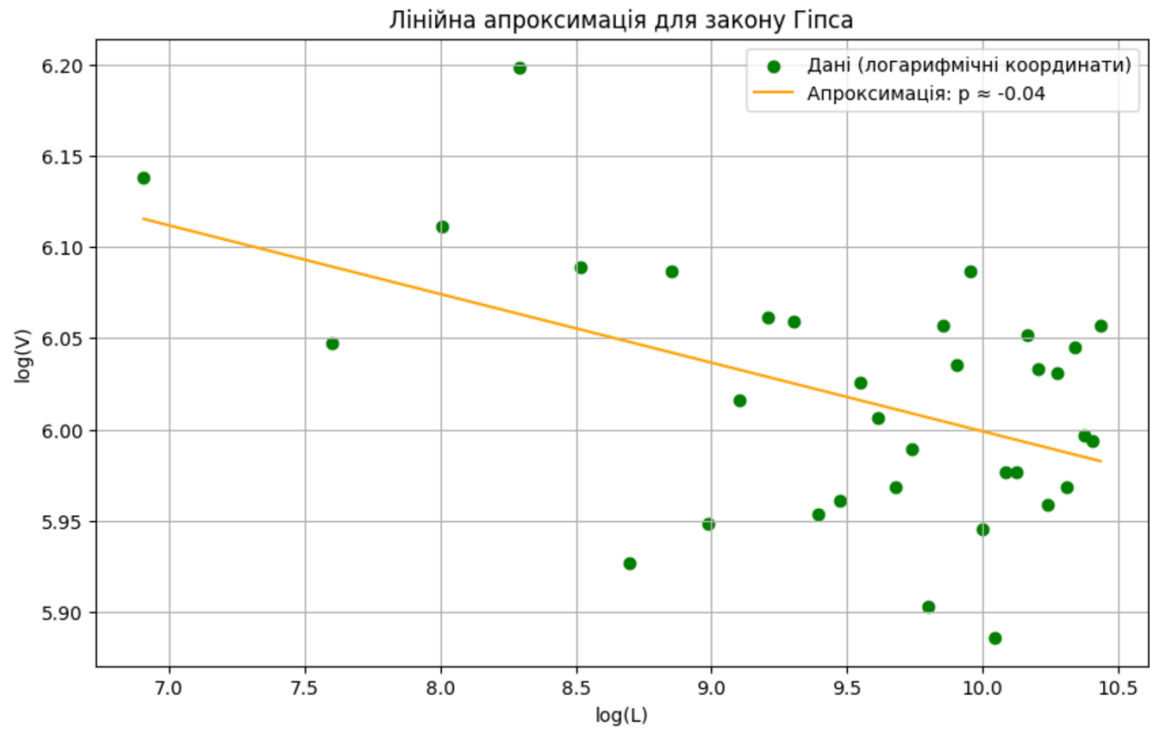


Результати виконання закону Гіпса для файлу ***Clemencia Novela de costumbres by Fernán Caballero***



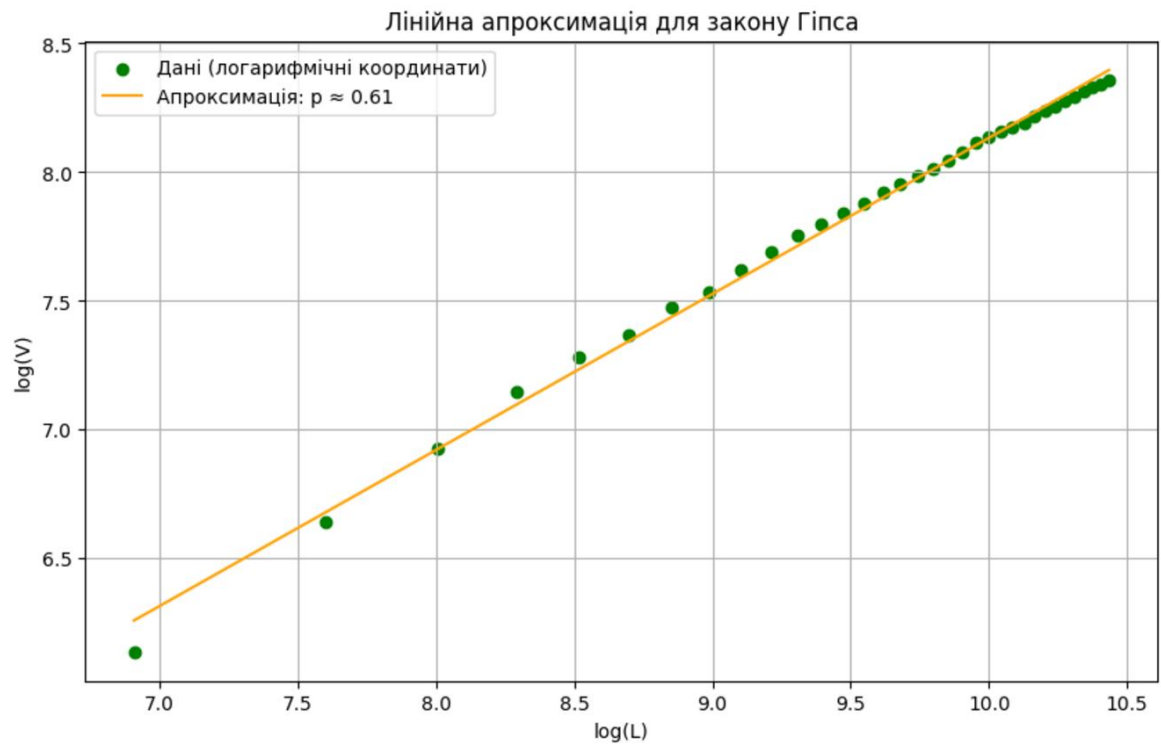
Коефіцієнт нахилу ( $p$ ): 0.65  
Перетин з віссю (intercept): 1.61

Лінійна апроксимація до результатів першого файлу першого тексту



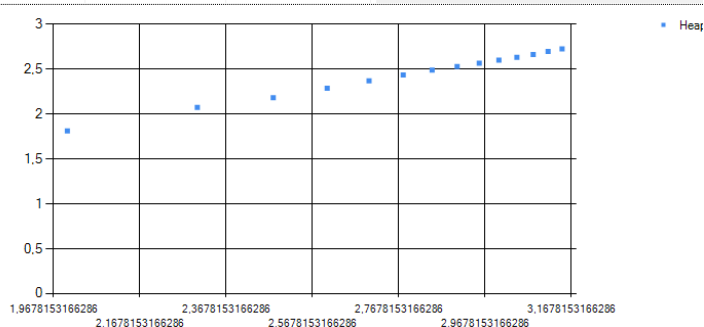
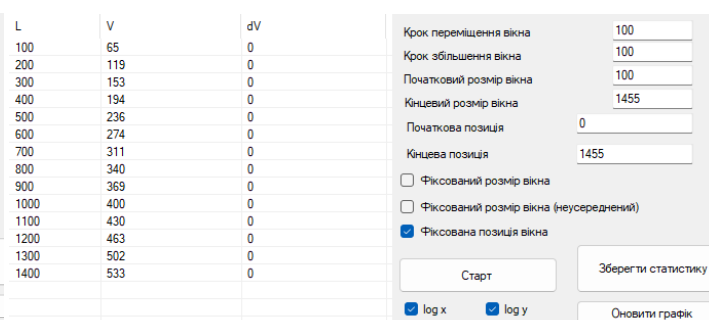
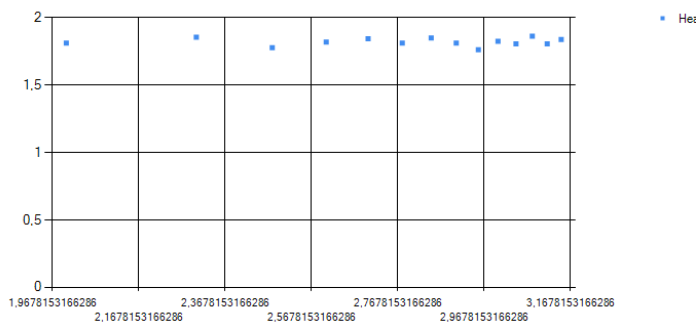
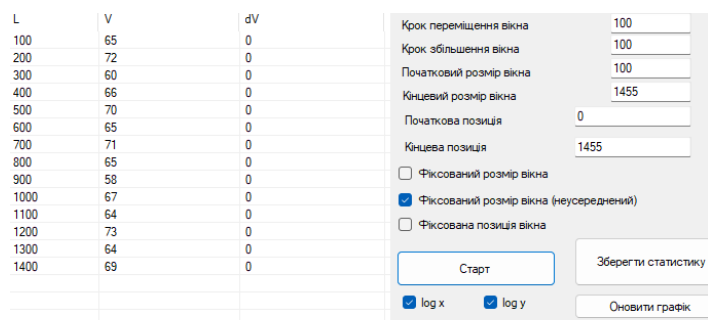
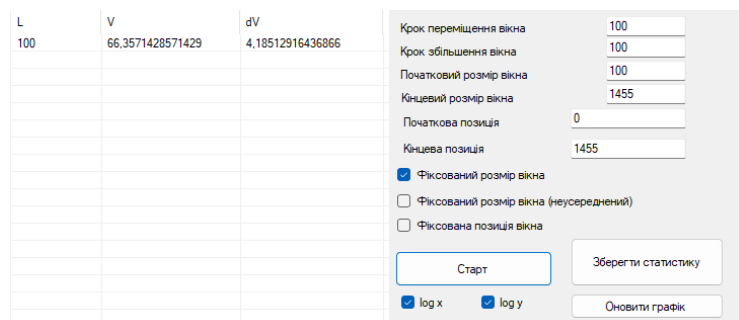
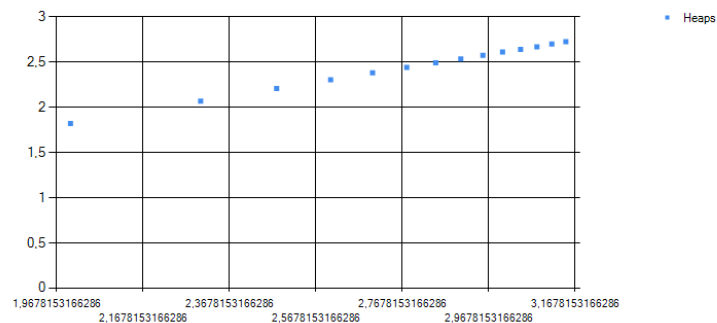
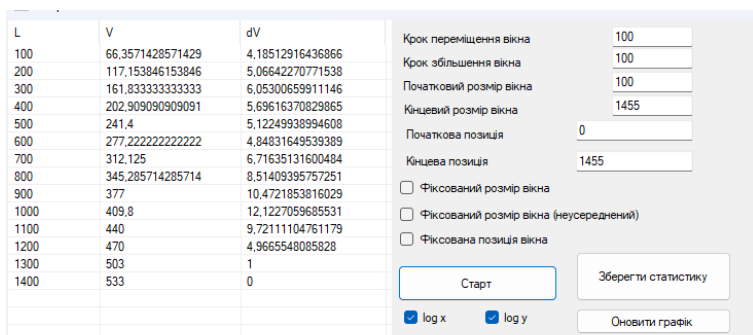
Коефіцієнт нахилу ( $p$ ):  $-0.04$   
Перетин з віссю (intercept): 6.38

Лінійна апроксимація до результатів другого файлу першого тексту

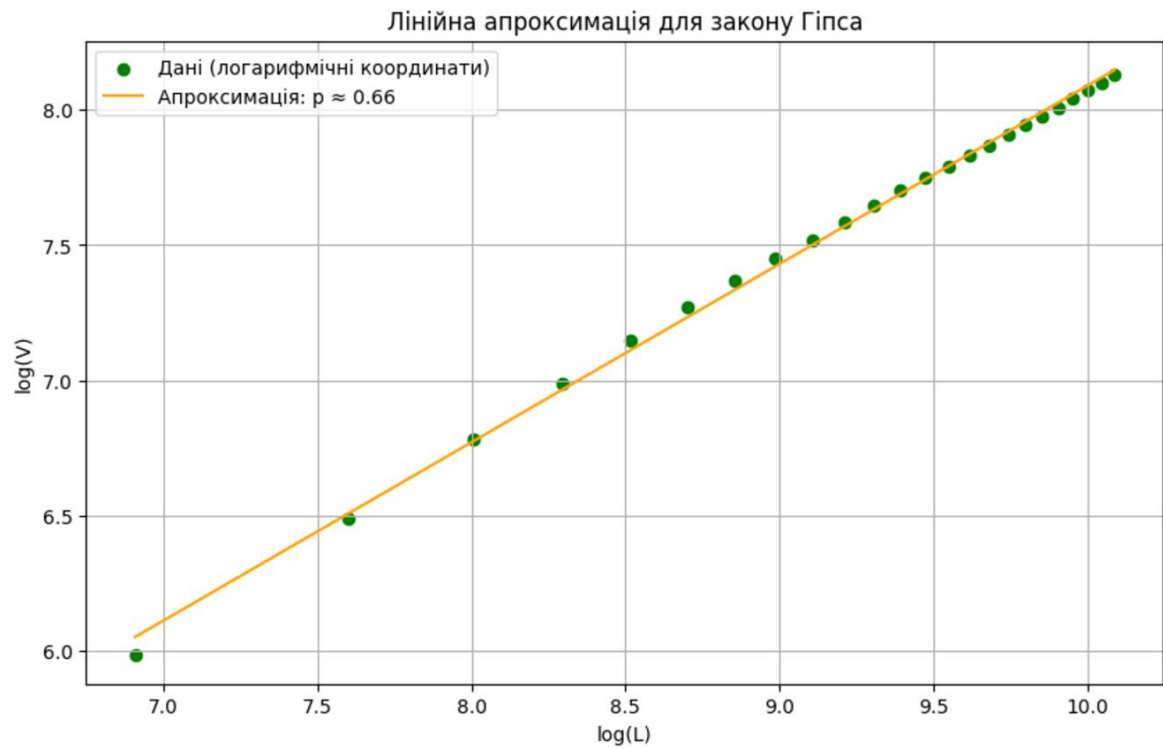


Коефіцієнт нахилу ( $p$ ):  $0.61$   
Перетин з віссю (intercept): 2.07

Лінійна апроксимація до результатів третього файлу першого тексту

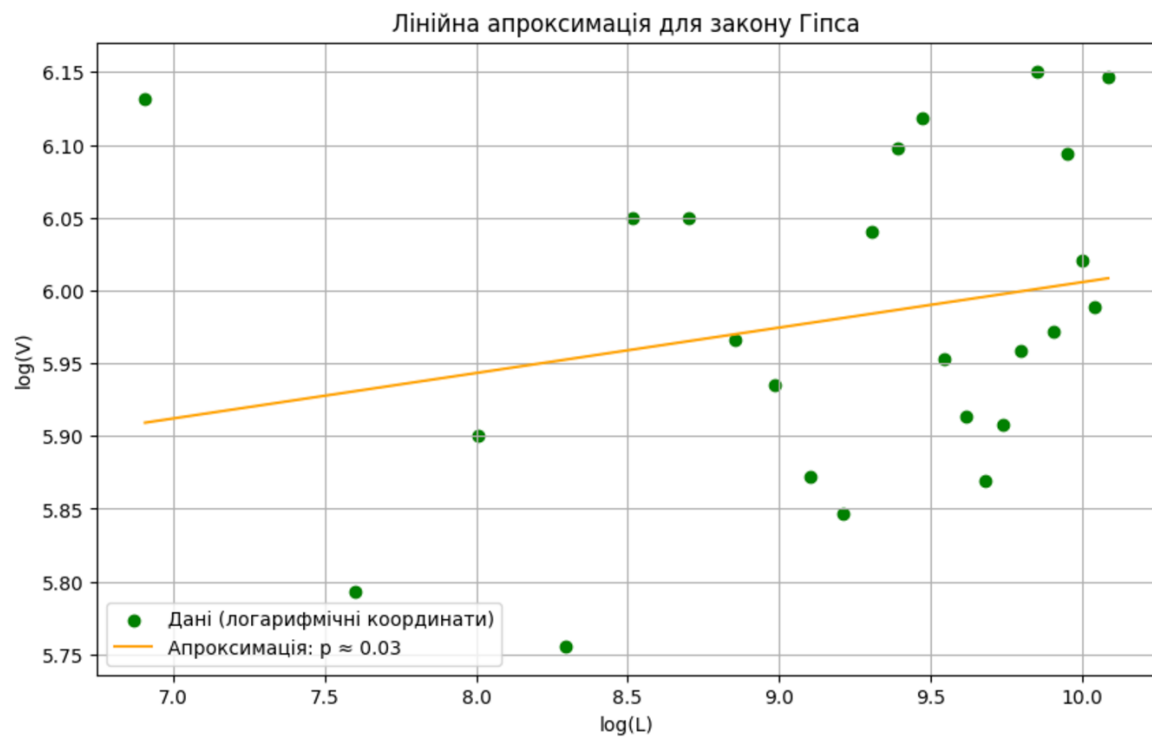


Результати виконання закону Гіпса для **DONALD J. TRUMP January 20, 2017**



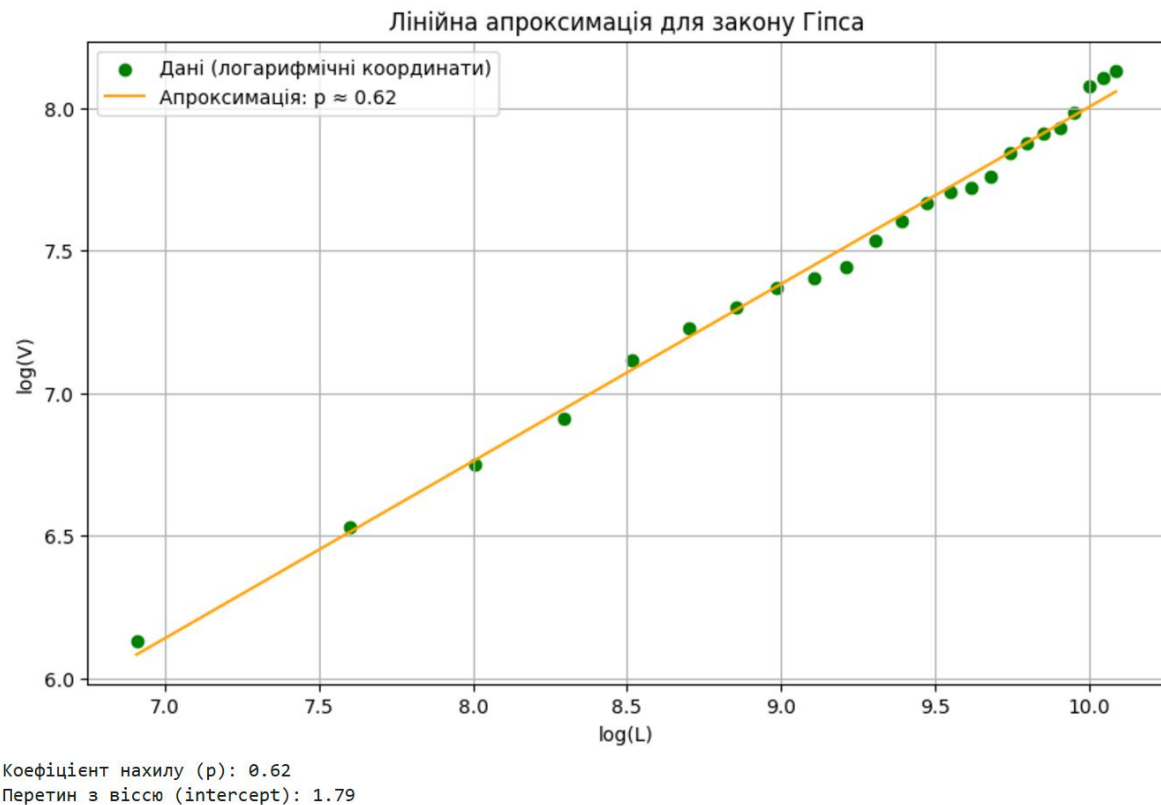
Коефіцієнт нахилу ( $r$ ): 0.66  
Перетин з віссю (intercept): 1.50

Лінійна апроксимація до результатів першого файлу другого тексту



Коефіцієнт нахилу ( $r$ ): 0.03  
Перетин з віссю (intercept): 5.69

Лінійна апроксимація до результатів другого файлу другого тексту



#### Лінійна апроксимація до результатів третього файлу другого тексту

Кожен з текстів має характерну експоненційну форму кривої, що підтверджує відповідність розподілу частот закону Гіпса. Це означає, що найбільш вживані слова мають значно вищу частоту, ніж менш вживані, що свідчить про ієрархічну структуру лексики.

Малі відхилення від ідеального експоненційного розподілу можуть бути зумовлені стилістичними особливостями текстів та авторським слововживанням. Однак загальна тенденція розподілу частот у обох текстах підтверджує відповідність закону Гіпса, що свідчить про те, що лексика кожного з текстів підпорядковується цьому закону.

**Висновок:** У цій лабораторній роботі я дослідив виконання закону Гіпса для текстів на рівні однограм. Результати показали, що частоти слів у текстах загалом слідують експоненційному розподілу, що підтверджує виконання закону Гіпса. Хоча на графіках можна помітити деякі відхилення, вони, ймовірно, зумовлені стилістичними особливостями авторів та структурою текстів. Отже, проведене дослідження підтверджує, що частоти слів у літературних текстах відповідають закономірностям, описаним законом Гіпса.