# Text Processing and Visualization with Python: Exercises in NLP

## Overview

This document outlines two exercises in natural language processing (NLP) and data visualization, using Python libraries for tokenization, frequency analysis, and word cloud generation. Each exercise demonstrates different aspects of text processing, from breaking down text into individual words to visualizing word frequencies with creative word clouds.

## Exercise 1: Tokenization and Frequency Analysis

**Objective:**
To process and analyze a sample text using tokenization and frequency distribution.

**Libraries Used:**

- `nltk` **(Natural Language Toolkit):** Employed for text tokenization (splitting text into words) and analyzing word frequencies within a sample passage.
- `matplotlib`**:** Used for visualization, particularly to view distributions (though not directly in this exercise).

**Process Summary:**

1. **Mount Google Drive:** Access Google Drive to easily work with external files in Colab.
2. **Install and Import Libraries:** Ensure `nltk` and other relevant libraries are installed for processing and analysis.
3. **Tokenization:** Divide the sample text into tokens (individual words and symbols), enabling further analysis.
4. **Frequency Analysis:** Calculate how often each word (or token) appears in the text using a frequency distribution, revealing the most common words and their respective counts.

**Output:**
A list of tokens from the text, with each token's frequency highlighted, enabling insights into the most frequent terms within the passage.

Frequency of each token:

Artificial: 1

Intelligence: 1
(: 1
AI: 3
): 1
is: 2
a: 1
rapidly: 1
evolving: 1
field: 1
that: 1
transforming: 1
industries: 1
and: 4
daily: 1
life: 1
.: 4
It: 1
enables: 1
machines: 1
to: 2
learn: 1
from: 1
data: 1
,: 9
make: 1
decisions: 1
solve: 1
problems: 1
with: 1
minimal: 1
human: 1
intervention: 1
applications: 1
are: 1
seen: 1
in: 1
healthcare: 1
finance: 1
automotive: 1
many: 1
other: 1
sectors: 1
As: 1
technology: 1

advances: 1
the: 1
potential: 1
of: 1
continues: 1
grow: 1
raising: 1
questions: 1
about: 1
ethics: 1
safety: 1
societal: 1
impact: 1

---

## Exercise 2: Word Cloud Generation

**Objective:**
To create word clouds that visually represent word frequencies, including custom-shaped clouds for unique presentation.

**Libraries Used:**

- `wordcloud`**:** Generates word clouds based on word frequencies.
- `matplotlib`**:** Displays the word clouds visually.
- `PIL` **(Python Imaging Library):** Loads and manipulates images, enabling custom shapes for the word clouds.
- `numpy`**:** Converts images to arrays for mask application, shaping the word clouds.

**Process Summary:**

1. **Text Data Collection:** Load text data from an external source (e.g., a novel or document), forming the basis of the word cloud.
2. **Generate General Word Cloud:** Create a standard word cloud without custom shapes, using default filters (stopwords) to exclude common, unimportant words.
3. **Display Word Cloud:** Use `matplotlib` to display the word cloud, presenting word frequency visually where larger words appear more frequently in the text.
4. **Custom-Shaped Word Clouds:**
   - **Load Mask Images:** Choose images to serve as custom masks, giving the word cloud unique shapes.
   - **Generate Custom Word Clouds:** Apply the mask images to produce shaped word clouds, adding a creative visual aspect to the text analysis.

**Output:**
A series of word clouds displaying frequent terms from the text data, with different shapes (based on chosen masks) for enhanced visual appeal.



---

## Summary

These exercises offered insights into:

- Tokenizing and analyzing text data to determine word frequencies.
- Creating effective visualizations with general and custom-shaped word clouds, making it easier to interpret large text data at a glance.

By leveraging `nltk`, `wordcloud`, `PIL`, and `matplotlib`, we explored various text processing and visualization techniques to make textual data more accessible and visually engaging.