# Text Similarity Using Jaccard Similarity

## Overview

*Jaccard Similarity* is a statistical measure used to gauge the similarity and diversity between two sets. In the context of text analysis, it quantifies the overlap between two texts by comparing the set of unique words present in each. The Jaccard Similarity coefficient is defined as the size of the intersection divided by the size of the union of the sample sets. ?cite?turn0search6?

## Why Use Jaccard Similarity?

Jaccard Similarity is particularly useful for:

- **Simple Text Comparison**: It provides a straightforward method to compare texts based on shared unique words.
- **Plagiarism Detection**: By measuring the overlap between documents, it can help identify potential instances of plagiarism.
- **Document Clustering**: Assists in grouping similar documents by evaluating shared content.

## Prerequisites

Before running the code, ensure you have the following:

- **Python 3.x**: The latest version of Python installed.
- **Basic Python Libraries**: No additional installations are required as the code utilizes standard Python libraries.

## Files Included

- `jaccard_similarity.py`: Contains the implementation of the Jaccard Similarity calculation between two text strings.

## Code Description

The provided code calculates the Jaccard Similarity between two text strings. Here's a step-by-step breakdown:

1. **Define the Texts**:

```
text1 = "machine learning is amazing"
text2 = "deep learning is a subset of machine learning"
```

Two sample text strings are defined for comparison.

2. **Convert Texts to Sets of Words**:

```
set1 = set(text1.split())
set2 = set(text2.split())
```

Each text is split into words, and then converted into sets to extract unique words.

3. **Calculate Intersection and Union**:

```
 intersection = len(set1.intersection(set2))
union = len(set1.union(set2))
```

The intersection and union of the two sets are computed to determine common and total unique words.

4. **Compute Jaccard Similarity**:

```
 similarity = intersection / union
print(f"Jaccard Similarity: {similarity:.4f}")
```

The Jaccard Similarity is calculated by dividing the size of the intersection by the size of the union, and the result is printed.

---

# Expected Output

When the code is executed, it will output the Jaccard Similarity between the two provided text strings. For the given examples, the output will be:

```
 Jaccard Similarity: 0.4286
```

This indicates that approximately 42.86% of the unique words are shared between the two texts.

---

# Use Cases

- **Document Comparison**: Quickly assess the similarity between two documents based on shared unique words.
- **Data Deduplication**: Identify and remove duplicate or near-duplicate records in datasets.
- **Recommendation Systems**: Compare user profiles or items to recommend similar content.

---

# Advantages

- **Simplicity**: Easy to implement and interpret.
- **Efficiency**: Computationally efficient for small to medium-sized texts.
- **Applicability**: Useful in various domains, including natural language processing and information retrieval.

---

# Considerations

- **Sensitivity to Vocabulary**: Does not account for synonyms or semantic similarities between words.
- **Ignores Word Frequency**: Treats all words equally, regardless of how often they appear in the text.
- **Case Sensitivity**: The method is case-sensitive; converting texts to lowercase before processing can mitigate this issue.

---

# References

- [Jaccard Similarity - LearnDataSci](#)
- [How to Calculate Jaccard Similarity in Python - GeeksforGeeks](#)
- [Overview of Text Similarity Metrics in Python - Medium](#)