



FAST-NU PESHAWAR

# Natural Language Processing

*Sami Uddin Shinwari*

*samu.shinwari@gmail.com*

*MS Data Science*

## Part I

# Stemming in Natural Language Processing

## 1 Introduction

Stemming is a process in Natural Language Processing (NLP) that reduces words to their root or base form. The goal of stemming is to normalize words to ensure that different forms of a word (e.g., "running," "runner," "ran") are analyzed as the same root word (e.g., "run"). This helps in various NLP tasks such as information retrieval, text analysis, and text preprocessing.

## 2 Key Concepts of Stemming

### 2.1 Root Form Reduction

For example, the words "connected," "connecting," and "connection" are all reduced to "connect."

### 2.2 Algorithms

- **Porter Stemmer:** One of the most widely used stemming algorithms, which uses a series of rules to iteratively trim suffixes from words.
- **Snowball Stemmer:** An improvement over the Porter Stemmer with additional rules and optimizations.
- **Lancaster Stemmer:** A more aggressive stemmer compared to Porter and Snowball, often resulting in shorter stems.
- **Lovins Stemmer:** One of the earliest stemmers, known for its large set of rules and irregular word handling.

### 2.3 Use Cases

- **Search Engines:** Enhances search by matching similar word forms.
- **Text Analysis:** Improves accuracy in sentiment analysis, topic modeling, etc.
- **Machine Learning:** Preprocessing step for text classification, clustering, and other tasks.

### 3 Example

```
from nltk.stem import PorterStemmer

stemmer = PorterStemmer()
words = ["running", "runner", "ran", "runs"]
stems = [stemmer.stem(word) for word in words]

print(stems)
# Output: ['run ', 'runner ', 'ran ', 'run ']
```

### 4 Considerations

- **Over-stemming:** When different words are reduced to the same root incorrectly (e.g., "universe" and "university" both stemmed to "univers").
- **Under-stemming:** When words that should be stemmed to the same root remain distinct.
- **Language Dependency:** Stemming rules are language-specific, requiring different stemmers for different languages.

### 5 Comparison to Lemmatization

- **Stemming** involves removing word endings to achieve the root form, often without considering whether the root is a valid word.
- **Lemmatization** involves using a vocabulary and morphological analysis of words, returning the base or dictionary form (lemma) of a word.

```
from nltk.stem import WordNetLemmatizer

lemmatizer = WordNetLemmatizer()
words = ["running", "runner", "ran", "runs"]
lemmas = [lemmatizer.lemmatize(word, pos='v') for word in words]

print(lemmas)
# Output: ['run ', 'runner ', 'run ', 'run ']
```

In summary, stemming is a crucial step in NLP for text normalization, helping to improve the performance of various text processing applications by ensuring that different forms of a word are treated similarly.