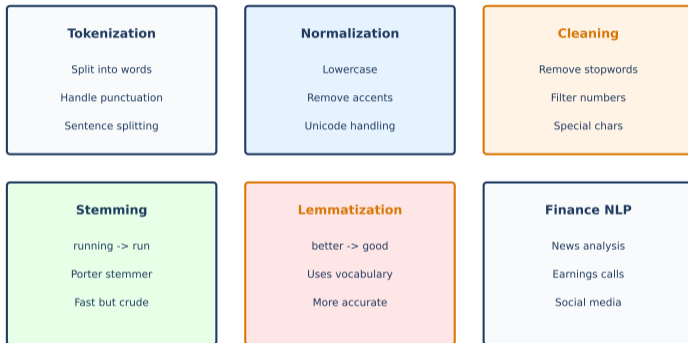


## Lesson 37 Summary: Text Preprocessing

Data Science with Python – Key Concepts

Data Science Program

## Text Preprocessing Pipeline



Clean text is essential for NLP success

## First step: split text into units

- **Word tokens:** Split on whitespace/punctuation
- **Sentence tokens:** Split on sentence boundaries
- **Subword:** BPE for rare words

---

Use `nltk.word_tokenize()` or `spacy`

## Standardize text format:

- **Lowercase:** Reduce vocabulary size
- **Unicode:** Normalize accents
- **Numbers:** Replace with tokens

---

`text.lower()` is simple but effective

## Remove low-information words:

- **Common:** the, a, is, and, of
- **Domain-specific:** Add custom stopwords
- **Caution:** Not always beneficial

---

```
from nltk.corpus import stopwords
```

# Stemming vs Lemmatization

## Reduce words to root form:

- **Stemming:** Fast, rule-based (running → run)
- **Lemmatization:** Accurate, uses vocabulary
- **Trade-off:** Speed vs accuracy

---

Lemmatization preferred for understanding

## Pattern matching for cleaning:

- **Remove URLs:** http pattern matching
- **Extract:** Emails, numbers, dates
- **Replace:** Normalize patterns

---

```
import re; re.sub(pattern, replacement, text)
```

### Essential Commands:

Task	Code
Tokenize	<code>nltk.word_tokenize(text)</code>
Lowercase	<code>text.lower()</code>
Stopwords	<code>nltk.corpus.stopwords.words('english')</code>
Stem	<code>PorterStemmer().stem(word)</code>
Lemmatize	<code>WordNetLemmatizer().lemmatize(word)</code>

---

Preprocessing quality determines model quality