# Python Tutorial 3

LING-381-Language Technology and LLMs

Instructor: Hakyung Sung

September 18, 2025

- Tokenization
- Lemmatization
- Frequency calculation
- Concordance

- What is tokenization?
- What were the approaches that we've used in the tutorial?

**Exercise**

- Choose a language model from spaCy that you would like to experiment with (except for English).
- Test different sentences that you want to tokenize.

```
1 # Put your output here
```

- What is lemmatization?
- What were the approaches that we've used in the tutorial?

## Exercise (in class)

- Check the spaCy lemmatizer with the new language model that you've chosen.

```
1 # load the model into a variable
2
3 # use the previous function to lemmatize the text
4
5
```

Word frequency calculation is most often the very first step of the text analysis.

## Exercise (in class)

```
1 # Download the `freq_test.txt` file from the tutorial
2
3 # Mount the downloaded file to the colab
4
5 # Example usage:
6 freq = corpus_freq("freq_test.txt")
7 print(freq.most_common(10))
```

Concordancing is a method in corpus linguistics that involves identifying and displaying occurrences of a target linguistic item along with its surrounding context.

```
1 sample_list = (
2     "Some people drink coffee every morning. "
3     "I prefer tea, but I enjoy iced coffee in the summer."
4 ).lower().split(" ")
5 💡
6 samp_hits = concord(sample_list, ["coffee"], 4, 4)  # target must be a list
7
8 for hit in samp_hits:
9     print(hit)

[['some', 'people', 'drink'], 'coffee', ['every', 'morning.', 'i', 'prefer']]
[['but', 'i', 'enjoy', 'iced'], 'coffee', ['in', 'the', 'summer.']]
```

Let's check the exercise together!

- Done with Tutorial 2 →Work with Collocation (Bonus 2 points)
- Want to fix your previous submission Tutorial 2 →Redo and resubmit
- Want to work on Tutorial 2 →Do and Submit

Submission deadline: Tomorrow 11:59PM

- **What:** Collocations = word pairs/phrases that co-occur more than chance

## Preview: Collocation analysis

- **What:** Collocations = word pairs/phrases that co-occur more than chance
- **Why:** Studying collocations via n-gram analysis helps to (1) assess native-like fluency (e.g., *make a decision* vs. *do a decision*) and (2) extract informative n-grams for information retrieval.

## Preview: Collocation analysis

- **What:** Collocations = word pairs/phrases that co-occur more than chance
- **Why:** Studying collocations via n-gram analysis helps to (1) assess native-like fluency (e.g., *make a decision* vs. *do a decision*) and (2) extract informative n-grams for information retrieval.
- **Break down steps**:

## Preview: Collocation analysis

- **What:** Collocations = word pairs/phrases that co-occur more than chance
- **Why:** Studying collocations via n-gram analysis helps to (1) assess native-like fluency (e.g., *make a decision* vs. *do a decision*) and (2) extract informative n-grams for information retrieval.
- **Break down steps**:
  - Tokenize the text

- **What:** Collocations = word pairs/phrases that co-occur more than chance
- **Why:** Studying collocations via n-gram analysis helps to (1) assess native-like fluency (e.g., *make a decision* vs. *do a decision*) and (2) extract informative n-grams for information retrieval.
- **Break down steps**:
    - Tokenize the text
    - Generate word frequency across the entire corpus

## Preview: Collocation analysis

- **What:** Collocations = word pairs/phrases that co-occur more than chance
- **Why:** Studying collocations via n-gram analysis helps to (1) assess native-like fluency (e.g., *make a decision* vs. *do a decision*) and (2) extract informative n-grams for information retrieval.
- **Break down steps**:
  - Tokenize the text
  - Generate word frequency across the entire corpus
  - Generate word frequency within contexts

## Preview: Collocation analysis

- **What:** Collocations = word pairs/phrases that co-occur more than chance
- **Why:** Studying collocations via n-gram analysis helps to (1) assess native-like fluency (e.g., *make a decision* vs. *do a decision*) and (2) extract informative n-grams for information retrieval.
- **Break down steps**:
  - Tokenize the text
  - Generate word frequency across the entire corpus
  - Generate word frequency within contexts
  - Calculate association strength between word pairs using statistical measures (e.g., Mutual information [MI])