# POS Tagging in Python

Here's a simple Python code example for POS (Part-of-Speech) tagging using the **Natural Language Toolkit (nltk)** library, which is beginner-friendly and widely used for NLP tasks:

```python
# Import necessary modules
import nltk

# Download the required resources
nltk.download('punkt')  # Tokenizer
nltk.download('averaged_perceptron_tagger')  # POS Tagger

# Input sentence
sentence = "I am learning Python for natural language processing."

# Tokenize the sentence into words
words = nltk.word_tokenize(sentence)

# Perform POS tagging
pos_tags = nltk.pos_tag(words)

# Print the POS tags
print("Word and POS Tags:")
for word, tag in pos_tags:
    print(f"{word}: {tag}")
```

## Explanation:

1. **Tokenization**: Splits the sentence into individual words (tokens) using `nltk.word_tokenize()`.

2. **POS Tagging**: Tags each token with its corresponding POS tag using `nltk.pos_tag()`.

3. **Output**: Prints each word along with its tag.

## Example Output:

For the sentence `"I am learning Python for natural language processing."`, the output might look like this:

```makefile
Word and POS Tags:
I: PRP
am: VBP
learning: VBG
Python: NNP
for: IN
natural: JJ
language: NN
processing: NN
.: .
```

Each tag (like `PRP`, `VBP`, `NN`, etc.) represents a part of speech:

- `PRP` : Pronoun

- `VBP` : Verb, present tense

- `NN` : Noun, singular

- `JJ` : Adjective

- `.` : Punctuation

Let me know if you'd like further details!

simple code for regular expression tagger

Here's a simple Python code to create a **regular expression-based tagger** for POS tagging using the **NLTK library**. This method assigns tags to words based on predefined patterns.

```python
import nltk
from nltk.tag import RegexpTagger

# Define patterns for tagging
patterns = [
    (r'.*ing$', 'VBG'),    # Gerunds (e.g., 'learning', 'running')
    (r'.*ed$', 'VBD'),     # Past tense verbs (e.g., 'played', 'hoped')
    (r'.*es$', 'VBZ'),     # Verbs ending in 'es' (e.g., 'goes', 'does')
    (r'.*ly$', 'RB'),      # Adverbs (e.g., 'quickly', 'slowly')
    (r'.*ion$', 'NN'),     # Nouns ending in 'ion' (e.g., 'information')
    (r'.*s$', 'NNS'),      # Plural nouns (e.g., 'cats', 'dogs')
```

```python
    (r'^[A-Z].*$', 'NNP'),  # Proper nouns (e.g., 'Python', 'India')
    (r'.*', 'NN')           # Default: Noun (e.g., 'book', 'apple')
]

# Create the RegexpTagger
tagger = RegexpTagger(patterns)

# Input sentence
sentence = "He is learning Python and writes code quickly."

# Tokenize the sentence into words
words = nltk.word_tokenize(sentence)

# Apply the RegexpTagger
tags = tagger.tag(words)

# Print the tagged words
print("Word and POS Tags:")
for word, tag in tags:
    print(f"{word}: {tag}")
```

## Explanation:

1. **Patterns**: A list of `(regex, tag)` pairs. These regex patterns match specific word endings or features to assign a POS tag.

2. **RegexpTagger**: Uses the defined patterns to tag words.

3. **Tokenization**: Splits the input sentence into individual words.

4. **Tagging**: Each word is tagged according to the regex patterns, or the default (`NN`) tag is applied.

## Example Output:

For the sentence `"He is learning Python and writes code quickly."`, the output might look like this:

```makefile
Word and POS Tags:
He: NN
is: NNS
learning: VBG
```

```
Python: NNP
and: NN
writes: NNS
code: NN
quickly: RB
.: NN
```

## Customization:

- You can modify the `patterns` list to fit specific tagging needs.

- This tagger is simple and may not perform as well as machine learning taggers but is useful for understanding regex-based approaches.

Let me know if you'd like further guidance!

simple code for unigram tagger

Here's a simple Python code to implement a **Unigram Tagger** using the **NLTK library**. A Unigram Tagger assigns the most probable POS tag to a word based on training data.

## Code:

```python
import nltk
from nltk.tag import UnigramTagger
from nltk.corpus import treebank

# Download required resources
nltk.download('treebank')  # Pre-trained corpus
nltk.download('punkt')     # Tokenizer

# Load a dataset (treebank) and split into training and testing data
tagged_sentences = treebank.tagged_sents()
train_data = tagged_sentences[:3000]  # Training data
test_data = tagged_sentences[3000:]   # Test data

# Create the UnigramTagger using training data
unigram_tagger = UnigramTagger(train_data)

# Test the UnigramTagger on some sample text
sentence = "I am learning Python for natural language processing."
```

```python
tokens = nltk.word_tokenize(sentence)

# Tag the sentence
tagged_sentence = unigram_tagger.tag(tokens)

# Print the tagged sentence
print("Tagged Sentence:")
for word, tag in tagged_sentence:
    print(f"{word}: {tag}")

# Evaluate the tagger's accuracy on the test data
accuracy = unigram_tagger.evaluate(test_data)
print(f"\nAccuracy on test data: {accuracy:.2f}")
```

## Explanation:

1. **Training Data**: The code uses the Treebank corpus (a pre-tagged dataset) to train the Unigram Tagger.

   - The training data is the first 3000 sentences.

   - The test data (remaining sentences) is used for evaluation.

2. **Unigram Tagger**:

   - It tags each word with the most frequent tag it saw during training.

   - Words not seen in the training data are tagged as `None`.

3. **Tokenization**: The sentence is tokenized into individual words before tagging.

4. **Evaluation**: The `evaluate()` method checks the tagger's accuracy on the test data.

## Example Output:

For the sentence: `"I am learning Python for natural language processing."`

You might get an output like this:

```
yaml
```

```
Tagged Sentence:
I: PRP
am: VBP
learning: VBG
Python: None
for: IN
natural: JJ
language: NN
processing: NN
.: .


Accuracy on test data: 0.86
```

## Notes:

- `None` appears for unknown words (e.g., "Python") not seen during training.

- To improve performance for unseen words, you can combine the Unigram Tagger with a **backoff tagger** (like the Default Tagger).

Let me know if you want an example with backoff!

> simple code for n gram tagger

Here's a simple Python code to implement an **N-gram Tagger** (e.g., Bigram or Trigram) using the **NLTK library**. An N-gram Tagger assigns tags to words based on the context of the previous one or two words.

## Code:

```python
import nltk
from nltk.tag import BigramTagger, TrigramTagger, UnigramTagger
from nltk.corpus import treebank

# Download required resources
nltk.download('treebank')  # Pre-tagged dataset
nltk.download('punkt')     # Tokenizer

# Load and split the dataset into training and testing sets
tagged_sentences = treebank.tagged_sents()
```

```python
train_data = tagged_sentences[:3000]  # Training data
test_data = tagged_sentences[3000:]  # Testing data

# Create an N-gram tagger (BigramTagger or TrigramTagger)
# Use backoff to handle unseen words
unigram_tagger = UnigramTagger(train_data)  # Unigram as fallback
bigram_tagger = BigramTagger(train_data, backoff=unigram_tagger)
trigram_tagger = TrigramTagger(train_data, backoff=bigram_tagger)

# Test sentence
sentence = "I am learning Python for natural language processing."
tokens = nltk.word_tokenize(sentence)

# Tag the sentence using the TrigramTagger
tagged_sentence = trigram_tagger.tag(tokens)

# Print the tagged sentence
print("Tagged Sentence:")
for word, tag in tagged_sentence:
    print(f"{word}: {tag}")

# Evaluate the tagger on the test data
accuracy = trigram_tagger.evaluate(test_data)
print(f"\nAccuracy on test data: {accuracy:.2f}")
```

## Explanation:

1. **Unigram, Bigram, Trigram**:

   - The **Unigram Tagger** tags words individually.

   - The **Bigram Tagger** considers the previous word for context.

   - The **Trigram Tagger** considers the previous two words for context.

2. **Backoff**:

   - If the Trigram Tagger cannot tag a word (unseen pattern), it uses the Bigram Tagger.

   - If the Bigram Tagger fails, it falls back to the Unigram Tagger.

3. **Training and Testing**:

- The **Treebank corpus** provides pre-tagged sentences.

- Training is done on 3000 sentences, and the remaining are used for testing.

4. **Evaluation**:

- The `evaluate()` method measures the accuracy of the tagger on the test data.

---

## Example Output:

For the sentence: `"I am learning Python for natural language processing."`

The output might look like:

```yaml
Tagged Sentence:
I: PRP
am: VBP
learning: VBG
Python: None
for: IN
natural: JJ
language: NN
processing: NN
.: .


Accuracy on test data: 0.89
```

---

## Notes:

- The tagger's accuracy improves with backoff, as it can handle unseen words.

- You can use `BigramTagger` alone if you don't need Trigrams.

Let me know if you'd like further explanation or modifications!

simple code of affix tagger

An **Affix Tagger** in NLP is a rule-based tagger that assigns POS tags based on prefixes or suffixes of words. Here's a simple implementation using NLTK:

## Code:

```python
import nltk
from nltk.tag import AffixTagger
from nltk.corpus import treebank

# Download required resources
nltk.download('treebank')  # Pre-tagged dataset
nltk.download('punkt')     # Tokenizer

# Load and split the dataset into training and testing sets
tagged_sentences = treebank.tagged_sents()
train_data = tagged_sentences[:3000]  # Training data
test_data = tagged_sentences[3000:]  # Test data

# Create an Affix Tagger
# Assign tags based on the last 2 characters (suffix) of the word
affix_tagger = AffixTagger(train_data, affix_length=-2)

# Test sentence
sentence = "I am learning Python for natural language processing."
tokens = nltk.word_tokenize(sentence)

# Tag the sentence
tagged_sentence = affix_tagger.tag(tokens)

# Print the tagged sentence
print("Tagged Sentence:")
for word, tag in tagged_sentence:
    print(f"{word}: {tag}")

# Evaluate the tagger's accuracy on the test data
```

```
accuracy = affix_tagger.evaluate(test_data)
print(f"\nAccuracy on test data: {accuracy:.2f}")
```

## Explanation:

1. **Affix Length**:

   - The `affix_length=-2` parameter considers the last two characters (suffix) of each word to assign a tag.
   - You can modify `affix_length` (e.g., `-3` for suffixes or `2` for prefixes).

2. **Training Data**:

   - The `train_data` is a collection of pre-tagged sentences from the Treebank corpus.

3. **Tagging**:

   - Words are tagged based on their suffixes.
   - Unknown words (not matching patterns) are tagged as `None`.

4. **Evaluation**:

   - The `evaluate()` method calculates the accuracy of the Affix Tagger on test data.

## Example Output:

For the sentence: `"I am learning Python for natural language processing."`

The output might look like this:

```yaml
Tagged Sentence:
I: None
am: None
learning: VBG
Python: None
for: None
natural: None
```

```
language: NN
processing: NN
.: None


Accuracy on test data: 0.43
```

## Notes:

- The Affix Tagger works best for languages where suffixes/prefixes indicate grammatical roles (e.g., verb forms in English).

- Accuracy can be improved by combining it with a backoff tagger (e.g., UnigramTagger).

Let me know if you'd like to see how to combine it with backoff!

brill tagger simple code for

The **Brill Tagger** is a rule-based tagger that improves upon an initial tagging (like Unigram or Bigram Taggers) by learning transformation rules. Here's a simple example using **NLTK's Brill Tagger**.

## Code:

```python
import nltk
from nltk.tag import UnigramTagger
from nltk.tag.brill import BrillTaggerTrainer, Template
from nltk.tag.brill_trainer import Pos, Word
from nltk.corpus import treebank

# Download necessary resources
nltk.download('treebank')
nltk.download('punkt')

# Load and split the dataset into training and testing sets
```

```python
tagged_sentences = treebank.tagged_sents()
train_data = tagged_sentences[:3000]  # Training data
test_data = tagged_sentences[3000:]  # Test data

# Step 1: Create a base tagger (Unigram Tagger)
base_tagger = UnigramTagger(train_data)

# Step 2: Define templates for the Brill Tagger
templates = [
    Template(Pos([-1])),  # Use the POS tag of the previous word
    Template(Pos([1])),   # Use the POS tag of the next word
    Template(Pos([-2])),  # Use the POS tag of the word two steps back
    Template(Pos([2])),   # Use the POS tag of the word two steps ahead
    Template(Word([-1])), # Use the previous word
    Template(Word([1]))   # Use the next word
]

# Step 3: Train the Brill Tagger
trainer = BrillTaggerTrainer(initial_tagger=base_tagger, templates=templates)
brill_tagger = trainer.train(train_data, max_rules=10)  # Use up to 10 rules

# Step 4: Test the Brill Tagger
sentence = "I am learning Python for natural language processing."
tokens = nltk.word_tokenize(sentence)

# Tag the sentence
tagged_sentence = brill_tagger.tag(tokens)

# Print the tagged sentence
print("Tagged Sentence:")
for word, tag in tagged_sentence:
    print(f"{word}: {tag}")

# Evaluate the Brill Tagger's accuracy on the test data
accuracy = brill_tagger.evaluate(test_data)
print(f"\nAccuracy on test data: {accuracy:.2f}")
```

## Explanation:

1. **Base Tagger:**

   - A Unigram Tagger is used as the initial tagger. It provides a baseline for Brill Tagger to improve upon.

2. **Templates:**

   - Templates define the rules Brill Tagger will learn.

   - For example, `Pos([-1])` means the POS tag of the previous word is used to learn a rule.

3. **Brill Tagger Training:**

   - The `train()` method learns rules from the training data. The `max_rules=10` parameter limits the number of transformation rules.

4. **Testing:**

   - The trained Brill Tagger is applied to a tokenized sentence and outputs tags for each word.

5. **Evaluation:**

   - Accuracy is calculated on the test data.

---

## Example Output:

For the sentence: `"I am learning Python for natural language processing."`

You might see:

```yaml
Tagged Sentence:
I: PRP
am: VBP
learning: VBG
Python: NNP
for: IN
natural: JJ
language: NN
processing: NN
.: .
```

```
Accuracy on test data: 0.91
```

## Notes:

- Brill Tagger's performance depends on the base tagger and templates used.

- To improve tagging for unseen words, consider adding more templates or combining with backoff taggers.

Let me know if you'd like additional guidance!

simple code for NER tagger using spacy

**2/2**

Here's a simple example of Named Entity Recognition (NER) tagging using the **spaCy** library:

## Code:

```python
import spacy

# Load the pre-trained spaCy model
nlp = spacy.load("en_core_web_sm")

# Input sentence
sentence = "Barack Obama was the 44th President of the United States and was born in Hawaii."

# Process the sentence
doc = nlp(sentence)

# Extract named entities
print("Named Entities, their labels, and explanations:")
```

```
for ent in doc.ents:
    print(f"{ent.text}: {ent.label_} ({spacy.explain(ent.label_)})")
```

## Explanation:

1. **Loading the Model**:

   - `en_core_web_sm` is a small, pre-trained spaCy model that includes NER capabilities.

2. **Processing Text**:

   - `nlp(sentence)` processes the input sentence and identifies entities.

3. **Extracting Entities**:

   - `doc.ents` gives all the named entities in the sentence.

   - Each entity has a `text` (the entity itself) and a `label_` (the entity type).

4. **Entity Labels**:

   - `spacy.explain(ent.label_)` provides a human-readable explanation of the label
     (e.g., `PERSON`, `ORG`, `GPE`).

## Example Output:

For the sentence: `"Barack Obama was the 44th President of the United States and was born in Hawaii."`

The output might look like this:

```yaml
Named Entities, their labels, and explanations:
Barack Obama: PERSON (People, including fictional)
44th: ORDINAL (Ordinal numbers)
United States: GPE (Countries, cities, states)
Hawaii: GPE (Countries, cities, states)
```

## Notes:

- **Entity Types:**

  - `PERSON` : A person's name.

  - `GPE` : Geopolitical entities (e.g., countries, cities).

  - `ORG` : Organizations (e.g., companies).

  - `DATE` , `TIME` , `ORDINAL` , etc., for dates, times, and numerical references.

- **Custom Models:**

  - For domain-specific texts, you can train your own spaCy NER model.

Let me know if you'd like to explore training a custom NER model or additional examples!