

## PROJECT SPECIFICATION

## Part of Speech Tagging

## General Requirements

| CRITERIA   | MEETS SPECIFICATIONS   |
|--|--|
| Submission includes all files required for grading                     | <ul style="list-style-type: none"><li>Includes <code>HMM Tagger.ipynb</code> displaying output for all executed cells</li><li>Includes <code>HMM Tagger.html</code>, which is an HTML copy of the notebook showing the output from executing all cells</li></ul> |
| Submitted files are complete and do not include any disallowed changes | Submitted notebook has made no changes to test case assertions   |

## Baseline Tagger Implementation

| CRITERIA | MEETS SPECIFICATIONS |
|----------|----------------------|
|          |                      |

| CRITERIA   | MEETS SPECIFICATIONS  |
|--|---|
| Student correctly implements the <code>pair_counts()</code> function | <p>Emission count test case assertions all pass.</p> <ul style="list-style-type: none"> <li>• The emission counts dictionary has 12 keys, one for each of the tags in the universal tagset</li> <li>• "time" is the most common word tagged as a NOUN</li> </ul>      |
| Correct baseline MFC tagger implementation                           | <p>Baseline MFC tagger passes all test case assertions and produces the expected accuracy using the universal tagset.</p> <ul style="list-style-type: none"> <li>• &gt;95.5% accuracy on the training sentences</li> <li>• 93% accuracy the test sentences</li> </ul> |

## Calculating Tag Counts

| CRITERIA   | MEETS SPECIFICATIONS                  |
|--|---------------------------------------|
| Correct <code>unigram_counts()</code> implementation | All unigram test case assertions pass |
| Correct <code>bigram_counts()</code> implementation  | All bigram test case assertions pass  |

| CRITERIA   | MEETS SPECIFICATIONS                              |
|--|---|
| Correct<br><code>start_counts()</code><br>and<br><code>end_counts()</code><br>implementation | All start and end count test case assertions pass |

### Basic HMM Tagger Implementation

| CRITERIA   | MEETS SPECIFICATIONS  |
|--|---|
| Correct<br>HMM<br>network<br>construction        | All model topology test case assertions pass  |
| Correct<br>basic HMM<br>tagger<br>implementation | Basic HMM tagger passes all assertion test cases and produces the expected accuracy using the universal tagset. <ul style="list-style-type: none"> <li>• &gt;97% accuracy on the training sentences</li> <li>• &gt;95.5% accuracy the test sentences</li> </ul> |

## Suggestions to Make Your Project Stand Out!

Students may run their taggers on more complex datasets (for example, the `nltk.corpus.brown` or `nltk.corpus.treebank` datasets).

Students may also try more advanced HMMs:

- Using pseudocounts or interpolated smoothing to handle missing data
- Retrain the hidden markov model using Baum-Welch re-estimation (available via the `.fit()` method in Pomegranate)

---

[Student FAQ](#)