

cs294-5: Statistical Natural Language Processing

Assignment 3: Part-of-Speech Tagging

Due: Oct 17th

In this assignment, you will build the important components of a part-of-speech tagger, including a scoring model and a decoder.

Setup: The data for this assignment is in `/home/ff/cs294-5/corpora/assignment3`, or `/work/cs294-5/corpora/assignment3`. It uses the same Penn Treebank data as the first assignment, this time with the part-of-speech labels.

The starting source file for this assignment is

```
.../assignments/POSTaggerTester.java
```

Make sure you can access the source and data files (the only source update from assignment 2's code base is this java file).

The World's Worst POS Tagger: Now run the test harness, `assignments.POSTaggerTester`. You will need to run it with the command line option `-path DATA_PATH`, where `DATA_PATH` is whichever data path above you are using. This class is a fully functional, if minimalist, POS tagger. Its main method loads the standard Penn Treebank WSJ part-of-speech data set, split in the standard way into training, validation, and test sentences. The current code reads through the training data, extracting counts of which tags each word type occurs with. It also extracts a count over “unknown” words – see if you can figure out what it's doing and why. The current code then ignores the validation entirely, and gives each known word in the test data whichever tag it occurred with most frequently in the training data. Unknown words get the tag which had the most types in training. This tagger operates at about 92%, with a rather pitiful unknown word accuracy of 40%. Your job is to make a real tagger out of this one by upgrading of the placeholder components.

Part 1. A Better Sequence Model: Look at the main method – the `POSTagger` is constructed out of two components, the first of which is a `LocalTrigramScorer`. This scorer takes `LocalTrigramContexts` and produces a `Counter` mapping tags to their scores in that context. A `LocalTrigramContext` encodes a sentence, a position in that sentence, and the previous two tags. The dummy scorer ignores the previous tags, looks at the word at the current position, and returns a (log) conditional distribution over tags for that word:

$$\log P(t | w)$$

Therefore, the best-scoring tag sequence will be the one which maximizes the quantity:

$$\sum_i \log P(t_i | w_i)$$

Your first job is to upgrade the local scorer. You have a choice between building either an HMM tagger or a maximum-entropy tagger. If you choose to build a trigram HMM tagger, you will maximize the quantity

$$\sum_i \log [P(t_i | t_{i-1}, t_{i-2}) P(w_i | t_i)]$$

which means the local scorer would have to return

$$score(t_i) = \log P(t_i | t_{i-1}, t_{i-2}) P(w_i | t_i)$$

for each context. (Note that this is NOT a log distribution over tags). If you want to implement an MEMM, you will instead be maximizing the quantity

$$\sum_i \log P(t_i | t_{i-1}, t_{i-2}, \vec{w}, i)$$

which means that you will want to build a little maximum entropy model which predicts tags in these contexts, based on features of the contexts:

$$score(t_i) = \log P(t_i | t_{i-1}, t_{i-2}, \vec{w}, i)$$

(Note that this IS a log distribution over tags). You can do either. Warning: a full-blown maxent tagger will be *very* slow to train, on the order of hours per run, especially if you add too many feature schemas, so start early and give yourself plenty of time to run experiments. If you choose to build an HMM, which will train faster (but likely have lower accuracy), you should do something sensible for unknown words, using a technique like unknown word classes, suffix trees, or a simple maximum-entropy model of $P(\text{tag}|\text{UNK})$ used with Bayes' rule as part of your emission model.

Whichever model you choose to build, your local scorer should use the provided interface for training and validating. The assignment doesn't require that you use the validation data, but it's there if you want it. You should also get into the habit of not testing repeatedly on the test set, but rather use the validation set for tuning.

Part 2. A Better Decoder: With your improved scorer, your results should have gone up substantially. However, you may have noticed that the tester is now complaining about “decoder sub-optimality.” This is because of the second ingredient of the `POSTagger`, the decoder. The supplied implementation is a greedy decoder (equivalent to a beam decoder with beam size 1). Your final task in this assignment is to upgrade the greedy implementation with a Viterbi decoder. Decoders implement the `TrellisDecoder` interface, which takes a `Trellis` and produces a path. Trellises are

really just directed, acyclic graphs, whose nodes are states in a Markov model and whose arcs are transitions in that model, with weights attached. In this concrete case, those states are `State` objects, which encode a pair of preceding tags and a position in the sentence. The weights are scores from your local scorer. In this part of the assignment, however, it doesn't really matter where the Trellis came from. Take a look at the `GreedyDecoder`. It starts at the `Trellis.getStartState()` state, and walks forward greedily until it hits the dedicated end state. Your decoder will similarly return a list of states in which the first state is that start state and the last is the end state, but will instead return the sequence of least sum weight (recall that weights are log probabilities produced by your scorer and so should be added). A necessary (but not sufficient) condition for your Viterbi decoder to be correct is that the tester should show no decoder sub-optimality – these are cases where your model scored the correct answer better than the decoder's choice. As a target, accuracies of 94+ are good, and 96+ are basically state-of-the-art. Unknown word accuracies of 60+ are reasonable, 80+ are very good.

Note: if you want to write your decoder before your scorer, you can construct the `MostFrequentTagScorer` with an argument of `true`, which will cause it to restrict paths to tag trigrams seen in training – this boosts scores slightly and makes the greedy decoder suboptimal.

Write-up: For the write-up, I mainly just want you to describe what you've built. For a maxent model, you should mention what feature schemas you used and how well they worked. For an HMM model, you should discuss how you modeled unknown words. In either case, you should look through the errors and tell me if you can think of any ways you might fix them (whether you do fix them or not). Pay special attention to unknown words – in practice it's the unknown word behavior of a tagger that's most important.

Coding Tips: If you find yourself waiting on your maxent classifier, and want to optimize it, you will likely find that your code spends all of its time taking logs and exps in its inner loop. You can often avoid a good amount of this work using the `logAdd(x, y)` and `logAdd(x[])` functions in `math.SloppyMath`. Also, you'll notice that the object-heavy trellis and state representations are slow. If you want, you are free to optimize these to array-based representations. It's not required (or particularly recommended, really) but if you wanted to, you might find `util.Indexer` of use.